

Point-Set Kernel Clustering

Kai Ming Ting, Jonathan R. Wells, and Ye Zhu

Abstract—Measuring similarity between two objects is the core operation in existing clustering algorithms in grouping similar objects into clusters. This paper introduces a new similarity measure called point-set kernel which computes the similarity between an object and a set of objects. The proposed clustering procedure utilizes this new measure to characterize every cluster grown from a seed object. We show that the new clustering procedure is both effective and efficient that enables it to deal with large scale datasets. In contrast, existing clustering algorithms are either efficient or effective. In comparison with the state-of-the-art density-peak clustering and scalable kernel k-means clustering, we show that the proposed algorithm is more effective and runs orders of magnitude faster when applying to datasets of millions of data points, on a commonly used computing machine.

Index Terms—Cluster analysis, kernel clustering, point-set kernel, data dependent kernel

1 INTRODUCTION

SIMILARITY between two objects is used as the basis in grouping objects into clusters in existing clustering algorithms. State-of-the-art clustering algorithms are density based [1], and they rely on distance measure (aka a similarity measure) between two objects to compute the density of each data point, representing each object. An early influential density based algorithm DBSCAN [2] separates core points from noise points by using a density threshold, where the former has high density and the latter has low density. Then only core points in the neighbourhood of each other are grouped into the same cluster. The key advantage of DBSCAN is that it can detect clusters of arbitrary shapes and sizes. The exact condition under which DBSCAN fails to identify all clusters has been determined recently [3]. In general terms, DBSCAN fails to identify all clusters when the clusters have hugely varying densities.

Density-peak clustering (DP) [1], which is a more recent density-based algorithm, begins by identifying density peaks that are far from each other; and then links data points that are transitively connected to a peak to form a cluster. Because it does not employ a density threshold, DP has avoided the weakness of DBSCAN mentioned above. Though DP is a stronger clustering algorithm than DBSCAN in general [4], DP has its own weaknesses. A key weakness is the requirement to find all density peaks in the first step, after the density of each point has been estimated based on a distance/similarity measure. As a result, it has difficulty identifying clusters, where each cluster has uniform density distribution.

A recent research has shown that a kernelized DBSCAN called MBSCAN [4], which employs Isolation Kernel [5], overcomes the weakness of DBSCAN and uplifts its clustering performance to the same level as DP. However, it has high computational cost because it employs the same algorithm as DBSCAN and only replaces the distance measure with Isolation Kernel. Like DP, DBSCAN and MBSCAN are unable to deal with large scale datasets because they have quadratic time complexity.

In a nutshell, high computational cost is a longstanding fundamental issue of existing density-based clustering algorithms. We contend that the root cause is due to the use of a similarity between two data points, resulting the computational cost to be at least proportional to the square of the data size, i.e., n^2 , where n is the number of data points in a given dataset. This has restricted the existing density-based clustering algorithms to small datasets only, as evidenced on the datasets used in their evaluations [1], [2], [3], [4].

In the age of big data, these density-based algorithms could not be used, despite their superior clustering capability in comparison with more traditional clustering algorithms such as k-means [6]. Although there are attempts to parallelize these algorithms or approximate the clustering outcomes through sampling, the fundamental limitation remains, i.e., their computational cost being at least proportional to n^2 . In other words, these attempts are a mitigating approach that enables some large datasets to be executed in reasonable time. But huge datasets remain out of bound for these algorithms running on a machine with a fixed number of CPUs (see Section 6.2 for a scaleup test example.)

Significance

This paper introduces the first kernel-based clustering which has runtime proportional to data size and yields clustering outcomes that are superior to those of existing clustering algorithms. The success of the new clustering is due to the use of a proposed point-set kernel which has an exact finite dimensional feature map. This enables the new clustering to (i) characterize clusters of arbitrary shapes, varied densities and sizes in a dataset; and (ii) run orders of magnitude faster than existing state-of-the-art clustering algorithms which have quadratic time cost. It is the only clustering algorithm which can process millions of data points on a commonly used machine, as far as we know.

The rest of the paper is organized as follows. The proposed point-set kernel and its properties are described in the next three sections. Section 4 presents the proposed clustering algorithm, followed by a conceptual comparison in the next section. Section 6 provides the empirical evaluation results. A discussion of related issues and conclusions are provided in the last two sections.

K. M. Ting is affiliated with the National Key Laboratory for Novel Software Technology, Nanjing University, China. E-mail: tingkm@nju.edu.cn
J. R. Wells and Y. Zhu are affiliated with the School of Information Technology, Deakin University, Geelong, Australia. E-mail: jonathan.wells.research@gmail.com, ye.zhu@ieee.org.

2 PROPOSED POINT-SET KERNEL

Rather than relying on a similarity between two data points, we propose a new similarity which measures how similar a data point $x \in \mathbb{R}^d$ is to a set of data points G , as a point-set kernel:

$$\hat{K}(x, G) = \langle \Phi(x), \hat{\Phi}(G) \rangle \quad (1)$$

and

$$\hat{\Phi}(G) = \frac{1}{|G|} \sum_{y \in G} \Phi(y) \quad (2)$$

where $\hat{\Phi}$ is the *kernel mean map*¹ of \hat{K} ; Φ is the feature map of a point-to-point kernel κ ; and $\langle a, b \rangle$ denotes a dot product between two vectors a and b .

In contrast, the point-to-point kernel (a similarity between two data points), expressed as a dot product, is given as follows [9]:

$$\kappa(x, y) = \langle \Phi(x), \Phi(y) \rangle \quad (3)$$

Notice that the summation in $\hat{\Phi}(G)$ (in Equation 2) is to be done once only as a preprocessing. Then computing $\hat{K}(x, G)$ in equation 1, based on the dot product, takes a fixed amount of time only, independent of n (the data size of G .)

Therefore, to compute the similarity of x with respect to G for all points x in G , i.e., $\hat{K}(x, G) \forall x \in G$, has a computational cost which is proportional to n only.

Also note that the use of the feature map Φ is necessary in order to achieve the stated efficiency. The alternative, which employs the point-to-point kernel/distance directly in the computation, will have a computational cost that is proportional to n^2 —the root cause of high computational cost in existing density-based algorithms.

The point-set kernel formulation assumes that the point-to-point kernel κ has a finite-dimensional feature map Φ .

Commonly used point-to-point kernels (such as Gaussian and Laplacian kernels) have two key limitations [4], [5]: they have a feature map of intractable dimensionality [10]; and their similarity is independent of a given dataset. The first limitation prevents these kernels to be used in the proposed formulation directly.

\hat{K} built from Isolation Kernel

We propose to use a recently introduced point-to-point kernel which has an exact finite dimensional feature map called Isolation Kernel [4], [5] as κ in \hat{K} . It has two characteristics, which are antitheses to the two limitations mentioned above, i.e., it has a finite-dimensional feature map; and its similarity adapts to local density of the data distribution of a given dataset.

The first characteristic enables Isolation Kernel to be used directly in the proposed point-set kernel. The exact finite-dimensional feature map is crucial in achieving the only kernel-based clustering which has runtime proportional to data size, we will propose in Section 4.

The second characteristic has a specific data dependent property: *two points in a sparse region are more similar than two points of equal inter-point distance in a dense region* [4], [5]. This characteristic is crucial for the proposed clustering algorithm to obtain good clustering outcomes.

1. Kernel mean embedding [7], [8] is an approach to convert a point-to-point kernel into a distribution kernel which measures similarity between two distributions. The proposed point-set kernel can be viewed as a special case of kernel mean embedding. Kernel mean embedding uses the same kernel mean map we have stated here.

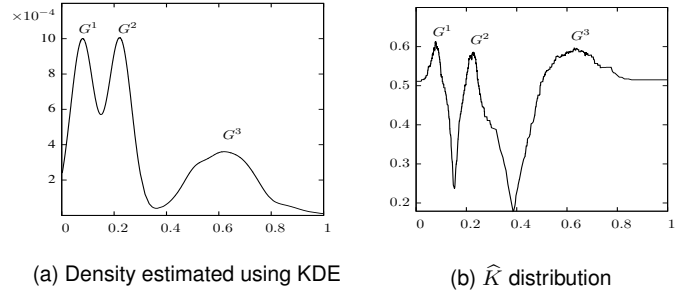


Fig. 1. Density distribution versus \hat{K} similarity distribution on a one-dimensional dataset having two dense clusters and one sparse cluster.

As the point-set kernel is constructed from a dataset D , Equations 1 and 2 are more precisely expressed as:

$$\hat{K}(x, G|D) = \langle \Phi(x|D), \hat{\Phi}(G|D) \rangle$$

and

$$\hat{\Phi}(G|D) = \frac{1}{|G|} \sum_{y \in G} \Phi(y|D)$$

where $G \subseteq D$; and Φ is the feature map of Isolation Kernel which is constructed from D (and Isolation Kernel has no functional form.)

The symbol ' $|D$ ' is dropped from the expressions in the rest of the paper for brevity.

The details of Isolation Kernel and its feature map can be found in Appendix A.

Once \hat{K} is derived from a dataset D , it is fixed and ready to be used in a kernel-based algorithm.

\hat{K} similarity distribution

The point-set kernel can be used to describe the data distribution of a dataset in terms of similarity distribution, independent of the clustering process. To do this, some sets in the dataset must be given, either as the ground truth or the clustering outcome of an algorithm.

Definition 1. Given sets $G^j, j = 1, \dots, k$ in a dataset D and the \hat{K} derived from D , the \hat{K} similarity distribution is defined as:

$$\max_j \hat{K}(x, G^j), \forall x \in \mathbb{R}^d.$$

The distribution is analogous to the density distribution estimated by, e.g., a kernel density estimator (KDE); except that sets in a dataset must be provided. In other words, the \hat{K} similarity distribution describes the data distribution in terms of the given sets in the dataset.

An example comparison of density distribution and \hat{K} similarity distribution of a same dataset is given in Figure 1.

3 PROPERTIES OF POINT-SET KERNEL FOR POINTS OUTSIDE A CLUSTER

To use the point-set kernel to grow a cluster, we need to understand the properties of the kernel for points outside the cluster.

Given a dataset D and an (expanding) cluster $C \subset D$.

Let $x \in D \setminus C$ and $x' \in D \setminus C'$; the distance between x and a set C be $\ell(x, C) \equiv \ell(x, \bar{x}_C)$, where \bar{x}_C is the preimage of $\hat{\Phi}(C|D)$; and $\rho(C) > \rho(C')$, where $\rho(C)$ denotes the average density of C .

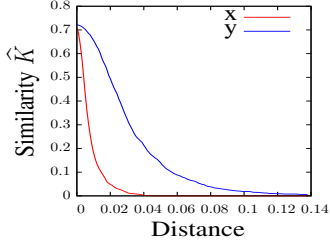


Fig. 2. Fall-off-the-cliff property of \hat{K} : x is close to the dense cluster C^2 and y is close to the sparse cluster C^3 . The x-axis: distances $\ell(x, C^2)$ and $\ell(y, C^3)$ correspond to x and y curves, respectively. C^2 and C^3 are subsets of G^2 and G^3 (target clusters), respectively, shown in Figure 1(a). Each C has the 50 highest density points in G .

Proposition 1. The point-set kernel $\hat{K}(x, C|D)$ derived from D has the following properties:

- (a) Fall-off-the-cliff property: $\hat{K}(x, C)$ decreases sharply as $\ell(x, C)$ increases.
- (b) Data dependent property: $\frac{d\hat{K}(x, C)}{dx} > \frac{d\hat{K}(x', C')}{dx'}$, if $\ell(x, C) = \ell(x', C')$ and $\rho(C) > \rho(C')$.

In other words, the rate of falling-off at x is data dependent: it is proportional to the average density of C . This property follows directly from the data dependent property of Isolation Kernel. See the proof in Appendix B.

These properties can also be understood from the implementation of Isolation Kernel. For each isolation partitioning, partitions at the boundary of C cover a limited area outside C . This contributes to the fall-off-the-cliff property. Also, the dense part of C has small partitions and the sparse part of C has large partitions. This creates different rates of falling-off mentioned above.

An example of the fall-off-the-cliff property is shown in Figure 2. It shows that the rate of falling-off is higher wrt the dense cluster than that wrt the sparse cluster.

An example of \hat{K} distribution, as stated in Definition 1, is shown in Figure 1, in comparison with the density distribution as estimated by a kernel density estimator using Gaussian kernel. Notice that the ‘valley’ between the two dense clusters is significantly sharper than that between the sparse cluster and the dense cluster. This is a direct result of the two properties mentioned above.

The point-set kernel that employs Isolation Kernel has two important influences. First, with an appropriate clustering algorithm design, the above-mentioned properties enable

- Each cluster to be expanded radially in all directions from a seed in multiple iterations, where each iteration recruits a subset of new members in the immediate neighborhood of the expanding cluster.
- Arbitrary-shaped clusters of different densities and sizes to be discovered via the above cluster expansion.

Second, in terms of runtime efficiency, the use of Isolation Kernel in the point-set kernel enables similarity between a point and a set to be computed efficiently, replacing point-to-point similarity/distance as the core operation. The latter is the root cause of high time complexity in existing algorithms.

The proposed point-set kernel clustering algorithm is designed based on these two properties. It is described in the next section.

4 CLUSTERING BASED ON POINT-SET KERNEL

The proposed clustering, called **point-set kernel clustering** or **psKC**, employs the point-set kernel \hat{K} to characterize clusters. It

Algorithm 1: point-set Kernel Clustering psKC

Input : D : dataset, τ : similarity threshold, ϱ : growth rate
Output: $G^j, j = 1, \dots, k$: k clusters, N : noise set

```

1  $k = 0$ ;
2 while  $|D| > 1$  do
3    $x_p = \operatorname{argmax}_{x \in D} \hat{K}(x, D)$ ; // Seed
4    $x_q = \operatorname{argmax}_{x \in D \setminus \{x_p\}} \hat{K}(x, \{x_p\})$ ;
5    $\gamma = (1 - \varrho) \times \hat{K}(x_q, \{x_p\})$ ;
6   if  $\gamma \leq \tau$  then
7     Terminate while-do loop;
8   end
9    $k++$ ;
10   $G_0^k = \{x_p, x_q\}$ ; // Initial cluster  $k$ 
11  for  $(i = 1; \gamma > \tau; i++)$  do
12     $G_i^k = \{x \in D \mid \hat{K}(x, G_{i-1}^k) > \gamma\}$ ;
13     $\gamma = (1 - \varrho)\gamma$ ;
14  end
15   $G^k = G_{i-1}^k$ ; // Cluster  $k$  grown
16   $D = D \setminus G^k$ ;
17 end
18  $N = D$ ;
19 return  $G^j, j = 1, \dots, k; N$ ;

```

identifies all members of each cluster by first locating the seed of the dataset. Then, it expands its members in the cluster’s local neighbourhood which grows at a set rate (ϱ) incrementally; and it stops growing when all unassigned points have similarity wrt the cluster falling below a threshold (τ). The process repeats for the next cluster using the remaining points in the given dataset D , yet to be assigned to any clusters found so far, until D is empty or no point can be found which has similarity more than τ . All remaining points after the clustering process are noise as they are less than the set threshold for each of the clusters discovered.

The psKC procedure is shown in Algorithm 1.

Here we formally define the cluster which is grown from a seed, according to psKC.

Definition 2. A \hat{K}_ϱ^τ -expanded cluster grows from a seed x_p selected from D , using $\hat{K}(\cdot, \cdot)$ with similarity threshold $\tau < 1$ and growth rate $\varrho \in (0, 1)$, is defined recursively as:

$$G_i = \{x \in D \mid \hat{K}(x, G_{i-1}) > \gamma_i > \tau\}$$

where $x_q = \operatorname{argmax}_{x \in D \setminus \{x_p\}} \hat{K}(x, \{x_p\})$; $G_0 = \{x_p, x_q\}$; $\gamma_i = (1 - \varrho)\gamma_{i-1}$; and $\gamma_0 = (1 - \varrho)\hat{K}(x_q, \{x_p\})$.

The number of iterations required to expand each cluster up to the threshold τ is shown in lines 11-14 in Algorithm 1 which are the key part of the algorithm. The rest of the algorithm is for initialization and house keeping only.

Let G^j be \hat{K}_ϱ^τ -expanded cluster j from D .

The number of \hat{K}_ϱ^τ -expanded clusters in D is discovered automatically by repeating the above cluster growing process on G^k from $D \setminus \{G^j, j = 1, \dots, k-1\}$.

Definition 3. After discovering all \hat{K}_ϱ^τ -expanded clusters G^j in D , noise is defined as

$$N = \{x \in D \mid \forall j \hat{K}(x, G^j) \leq \tau\}.$$

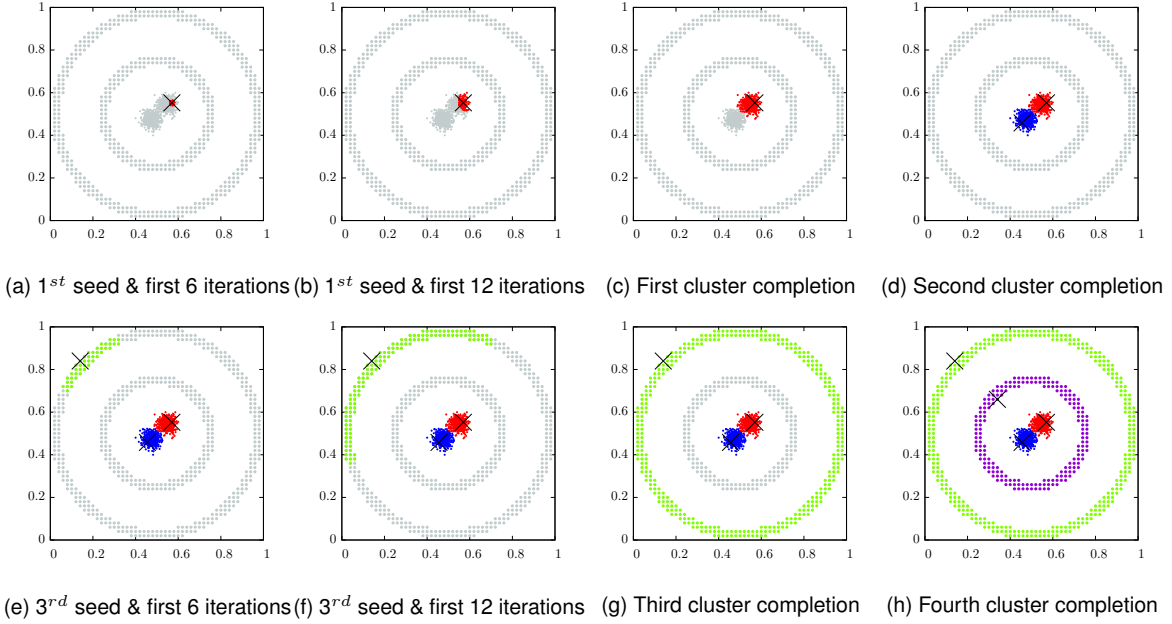


Fig. 3. psKC clustering outcome of the Ring-G dataset. Crosses indicate the seed points. Gray points are data points in D which are yet to be clustered. The number of iterations f refers to that in lines 11-14 in Algorithm 1. f used to recruit all members of each cluster is: 20, 19, 20 & 20.

A post-processing is applied to all clusters produced by psKC to ensure that the following objective is achieved:

$$\Gamma(D) = \max_{G^1, \dots, G^k} \sum_{j=1}^k \sum_{x \in G^j} \langle \Phi(x), \hat{\Phi}(G^j) \rangle. \quad (4)$$

This post-processing reexamines all points which have the lowest similarity wrt cluster G^j if they could be reassigned to other cluster to maximize the total similarity.

A demonstration of the clustering process of psKC is shown in Figure 3, based on a two-dimensional dataset which has data points distributed in two concentric rings and two Gaussian clusters. Figure 3(a) shows the progression of identifying the first seed and growing the first cluster to include all members found in the first 6 iterations; followed by including all members found in the first 12 iterations in Figure 3(b); and all members in the first cluster are found in Figure 3(c)—this is when $\hat{K}(x, G^1) < \tau$ for all x in D excluding G^1 . The same process in identifying the third cluster, i.e., the outer ring, is shown in the first three subfigures in the second row of Figure 3.

The time complexity of psKC is given in Table 1. Line 12 in the for-loop iteration in Algorithm 1 is the most expensive part. \hat{K} in line 12 needs to be evaluated $n - 2$ times in the very first iteration; and any subsequent iterations take monotonically decreasing number of \hat{K} evaluations. In the worst case, $f(\tau, \varrho)$ is the maximum number of iterations.

Proposition 2. *psKC performs a maximum number of iterations f in the cluster expansion process for each cluster, which is independent of data size, that is a function of τ and ϱ :*

$$f(\tau, \varrho) = \left\lceil \frac{\log \tau}{\log(1 - \varrho)} \right\rceil.$$

Proof. By setting the initial $\gamma = 1$ to its maximum value before the iteration in lines 11-14, the maximum number of iterations f is reached when $(1 - \varrho)^f = \tau$. As it is independent of the data size, this provides the proof. \square

This is a very useful property in practice because one can use it to bound the number of iterations. For example, the settings of $\tau \in [10^{-5}, 10^{-1}]$ and $\varrho \in [0.1, 0.26]$, that we have used in the experiments reported in Section 6, bound $f \in [8, 109]$.

In comparison, the time complexity of k-means is super-polynomial in the worst case which requires $f'(n) = 2^{\Omega(\sqrt{n})}$ iterations [11], where its time complexity is $O(f'(n)nk d)$ [12]. Thus, the key difference is the number of iterations, i.e., psKC controls $f(\tau, \varrho)$ via parameter setting which is independent of n ; k-means has no such mechanism and $f'(n)$ is a function of n . In practice, psKC has its iterations bounded by the search ranges of τ and ϱ , as we have shown earlier, giving it a linear time complexity. Though $f'(n)$ is typically a small fraction of n in practice, no proof has been provided thus far.

TABLE 1

Time complexities of psKC and k-means.
 t and ψ are parameters of Isolation Kernel; d is number of dimensions.

1	Build Isolation Kernel (IK)	$dt\psi$
2	Mapping D of n points using feature map of IK	$ndt\psi$
3	$f(\tau, \varrho)$ iterations in lines 11-14 for k clusters	$f(\tau, \varrho)nkdt\psi$
Total time cost for psKC		$f(\tau, \varrho)nkdt\psi$
k-means, $f'(n)$ is typically a small fraction of n in practice		$f'(n)nk d$

5 CONCEPTUAL DIFFERENCES IN COMPARISON WITH EXISTING CLUSTERING ALGORITHMS

The proposed clustering algorithm psKC is unique among the existing clustering algorithms in two key aspects:

- It is the only clustering algorithm which utilizes the **point-set kernel that employs Isolation Kernel**.
- The algorithmic design is unique: the main step in the procedure employs the proposed point-set kernel to grow every cluster from a seed.

Both the point-set kernel and the procedure are crucial in achieving a clustering algorithm which is both effective and efficient.

The finite dimensional feature map of Isolation Kernel and its use in the point-set kernel enable the algorithm to achieve its full potential: runtime proportional to data size (n)—a level unable to be achieved by existing effective clustering algorithms such as DP [1], and even less effective but efficient algorithms such as scalable kernel k-means [13]. They have at least n^2 runtimes.

Existing clustering algorithms have the following features:

- (a) Density-based clustering algorithms such as DP and DB-SCAN rely on point based density estimation that in turn relies on point-to-point distance calculations. This is the key reason why they both have n^2 runtime.
- (b) Kernel k-means [13], [14]) is weaker than psKC for two reasons. First, the clustering algorithm employed, i.e., k-means, produces clusters that are limited to globular shape and approximately the same size [15]. These limitations are inherited in feature space with the use of kernel. That is the reason why kernel k-means could not identify clusters of arbitrary shapes and different sizes (see Table 2 for details). Second, the use of a kernel which has intractable dimensionality and is data independent. The intractable dimensionality issue necessitates a kernel functional approximation [10] in order to produce a finite-dimensional feature map. This approximation reduces the quality of the final clustering outcome. The data independence issue (often ignored/unrecognized) also leads to poorer clustering outcomes.

It is possible to view the centers in kernel k-means as a kind of kernel mean map (defined in Equation 2) because each cluster center is defined as the average position of all points in a cluster in the feature space. However, both limitations mentioned above have constricted the types of clusters it can detect.

In a nutshell, many existing clustering algorithms compute point-to-point distance/kernel to derive the required similarity to characterize clusters. Even in algorithms that can be considered to have used point-set kernel, their cluster characterizations have severe limitations. In contrast, the proposed algorithm utilizes the point-set kernel built from Isolation Kernel. This allows clusters of arbitrary shapes, sizes and densities to be characterized successfully and efficiently.

A more detailed comparison with kernel k-means is instructive. Kernel k-means uses the following objective function:

$$\operatorname{argmin}_{G^1, \dots, G^k} \frac{1}{n} \sum_{j=1}^k \sum_{x \in G^j} \| \Phi(x) - \hat{\Phi}(G^j) \|_2^2 \quad (5)$$

The difference in objective functions between psKC and kernel k-means is small (between Equations 4 and 5), i.e., the latter uses a squared distance; and the former employs a dot product.

The key differences are the algorithms used to achieve these objectives. The k-means procedure iterates over two main steps: the center computation for each group; and the reassignment of every point to its nearest center. The simplicity is k-means's advantage in time complexity over other existing algorithms. But it also gives rise to its shortcomings. First, the initialization is a random grouping of clusters, and the number of clusters is specified by the user. The initial random grouping could lead to (a) counter-intuitive clustering outcomes even when the number of clusters is selected correctly; and (b) the clustering results are unstable from one run to the next. This means that k-means can

easily get trap in local minimum. Second, the procedure always produces spherical clusters of similar size in the feature space of a kernel. Third, it recomputes distances to centers for many points which do not change the assignment in the previous iteration. This constitutes a large expense of the computational cost which adds no value but is necessary in the procedure. Fourth, the use of kernel mean map to characterize each center makes k-means sensitive to outliers. In other words, the above shortcomings of kernel k-means are algorithmic, not due to a specific kernel used.

The psKC procedure does not have these shortcomings. First, the algorithm is deterministic, given a kernel function and the user-specified parameters. This resolves the instability issue and often leads to better clustering outcomes. The only randomisation is due to the Isolation Kernel. The use of most similar points in D as seeds is much more stable, even with different initializations of Isolation Kernel, compare with random groupings of clusters which can change wildly from one run to the next. Second, psKC enables detection of clusters of arbitrary shape, of different sizes and densities. Third, psKC commits each point to a cluster once it is assigned; and most points which are similar to the cluster never need to be reassigned. This is possible because of the use of a seed to grow a cluster. Points which are similar to a cluster grown from the seed will not be similar to another cluster if the points are less similar to the seeds of other clusters in the first place. The sequential determination of seeds (as opposed to the parallel determination of centers in k-means) makes that possible. As a result, psKC avoids many unnecessary recomputations in k-means mentioned earlier. Fourth, in contrast to k-means, the use of kernel mean map, not as a cluster center, but as a medium to expand a growing cluster in psKC makes it robust to outliers. This is because outliers are least similar to a growing cluster; they are thus unlikely to be included in the cluster.

There are two other important differences: **psKC is not an optimization algorithm**; whereas k-means is an expectation-and-maximization optimization method; and yet the clustering outcome of psKC is already close to the final maximization objective. The post-processing, literally tweaks at the edges, by reexamining those lowest similarity points wrt each cluster for possible better reassignments.

The relationship of spectral clustering and kernel k-means is given in Appendix C.

We show in the next section that the proposed clustering algorithm psKC is both highly efficient and producing good clustering outcomes in an empirical evaluation.

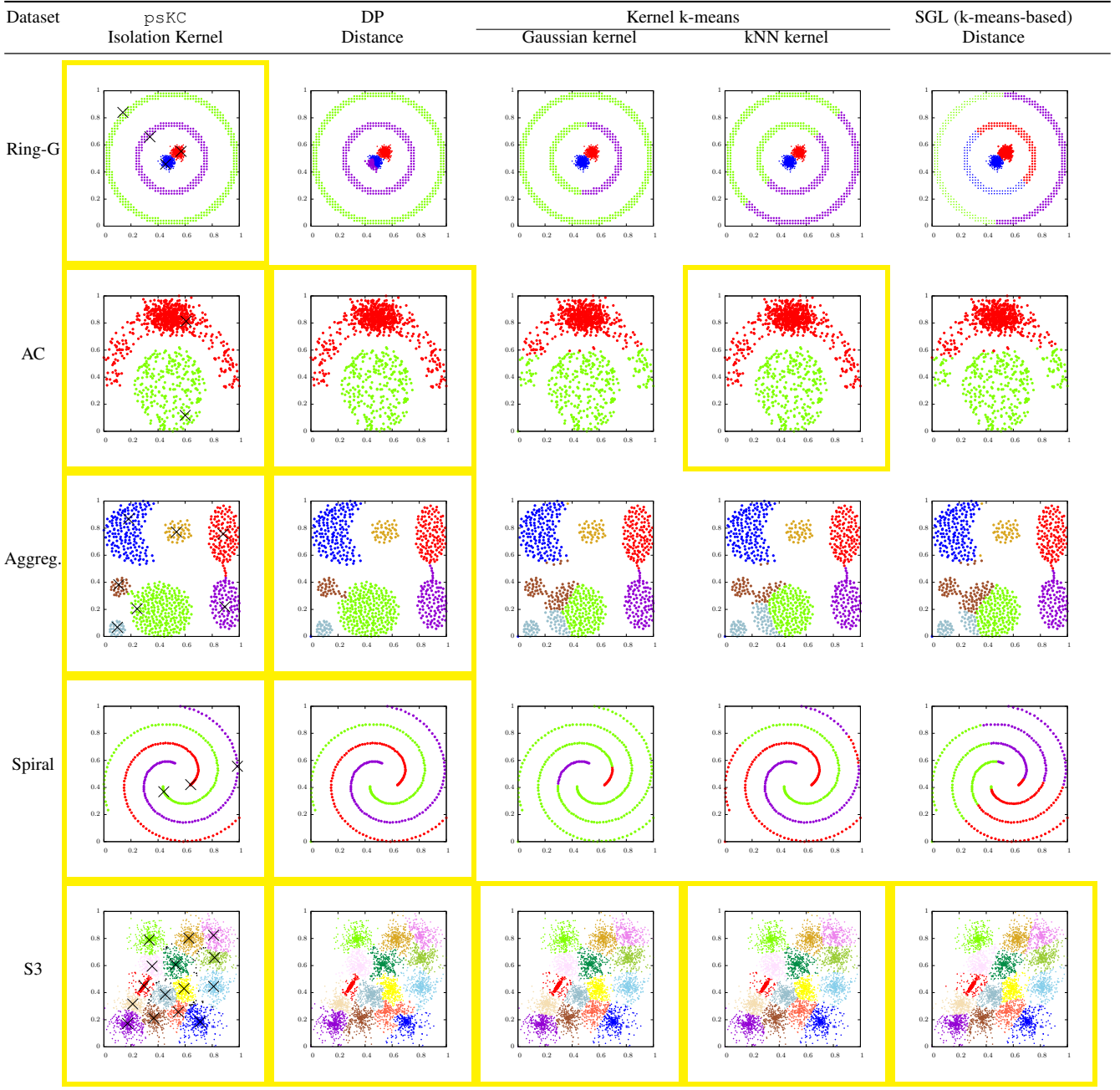
6 EMPIRICAL EVALUATION

We empirically compare psKC with DP [1], scalable kernel k-means [13] which employs Gaussian kernel, kernel k-means which employs adaptive kernel [14] and a recent k-means-based algorithm SGL [16] in terms of clustering outcomes and runtimes.

Commonly used benchmark datasets as well as color images are used in the experiments. As all these datasets can be visualized, we present the clustering outcomes of the algorithms under comparison by showing their segmented images or two-dimensional plots for visual inspection. The runtime is measured in terms of CPU seconds (and include the GPU seconds when GPU is used.) Color images are represented in the CIELAB color space². Other experimental settings can be found in Appendix D.

2. In three dimensions as defined by the International Commission on Illumination (<http://www.cie.co.at/>).

TABLE 2
Artificial datasets: Clustering outcomes of psKC , DP, two versions of kernel k-means and one recent k-means-based SGL [16].
The results with yellow frames indicate good clustering outcomes; and those without have poor clustering outcomes.



We present the results in four subsections: The first reports the clustering outcomes; the second presents the runtime comparison; the third provides the result of a stability analysis; and the last examines the effect of the post-processing.

6.1 Clustering outcomes

The clustering outcomes of artificial datasets, image segmentation and finding groups of handwritten digits are presented in three separate sections below.

6.1.1 Artificial datasets

A comparison of clustering outcomes of five clustering algorithms on five benchmark datasets is provided in Table 2. It is interesting

to note that DP did well in four benchmark datasets, but it did poorly on the Ring-G dataset, i.e., DP successfully identified one ring cluster and one Gaussian cluster; but split the second Gaussian cluster into two parts, where one part is joined with the second ring cluster. This is because DP has a weakness in identifying the correct peaks when (some) clusters are uniformly distributed and have varied densities. In this case, three out of the four peaks are identified in the two Gaussian clusters before other points are assigned to one of the peaks to form clusters.

All three k-means-based algorithms are weaker than DP as they did poorly on at least three out of the five datasets, i.e., Ring-G, Aggregation and Spiral. This is because of their use of k-means which has weaknesses in detecting clusters that have non-globular

TABLE 3

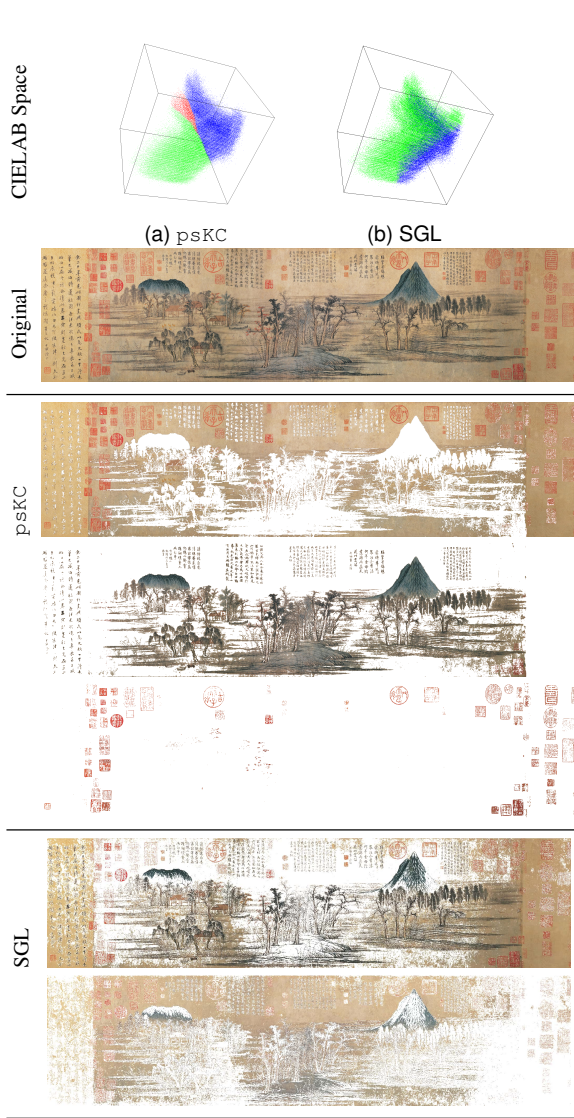
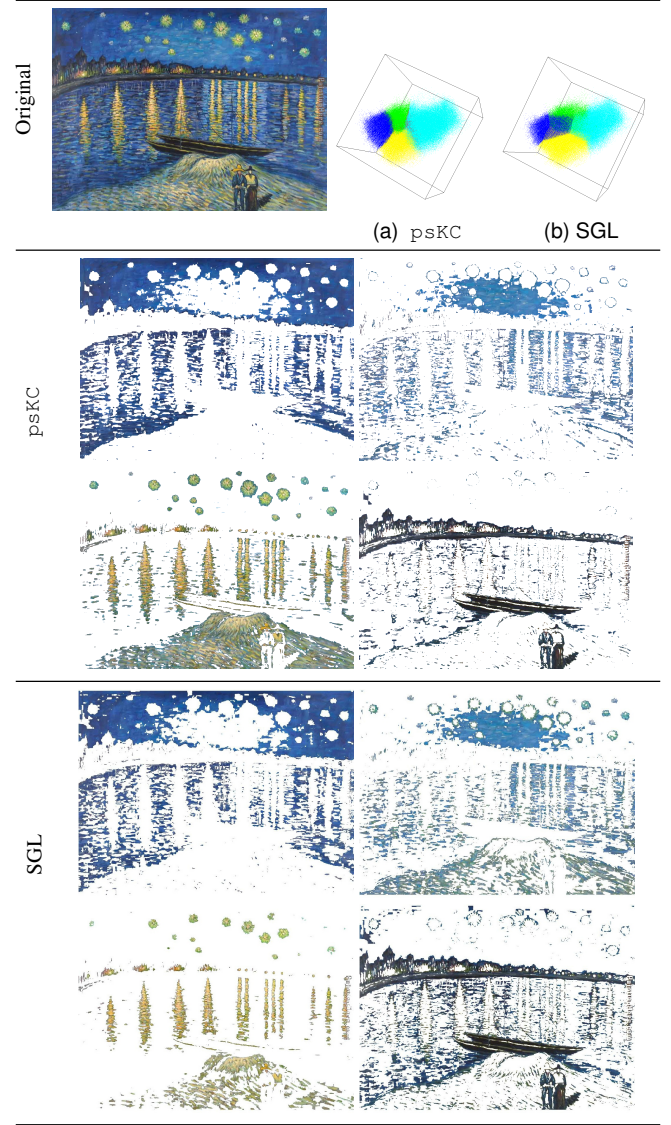
Clustering outcomes of psKC & SGL on Zhao Mengfu's Autumn Colors.

TABLE 4

Clustering outcomes of psKC and SGL on Vincent van Gogh's Starry Night over the Rhone 2.

shapes [17]. The use of a kernel in both kernel k-means transfers these fundamental weaknesses from input space to feature space. The results in Table 2 show that there is no guarantee that they can detect clusters of non-globular shapes in input space.

psKC is the only algorithm that did well in all five datasets; and it is the only algorithm that successfully identified all four clusters in the Ring-G dataset. This is a direct result of the proposed cluster identification procedure which employs the point-set kernel. Other algorithms failed to correctly identify the four clusters because their algorithmic design which must determine all density peaks/centers before individual points can be assigned to one of the peaks/centers.

The S3 dataset is the easiest to cluster. All five algorithms have good clustering outcomes though the outcomes differ slightly locally. The AC dataset is the second easiest. All algorithms, except kernel k-means with Gaussian kernel, have produced the perfect clustering outcome. Both Aggregation and Spiral were successfully clustered by psKC and DP; but both versions of kernel k-means failed to separate all clusters correctly.

6.1.2 Image segmentation

Here we examine the ability of the five clustering algorithms in dealing with images of high resolution.

Out of the four contenders of psKC , only SGL could complete the run in reasonable time.

Table 3 shows the clustering outcomes of psKC and SGL on Zhao Mengfu's Autumn Colors. psKC separates the background (plus red stamps) of this painting from the landscape, producing two distinctive clusters. The red stamps can be extracted as a cluster on its own if a different parameter setting is used.

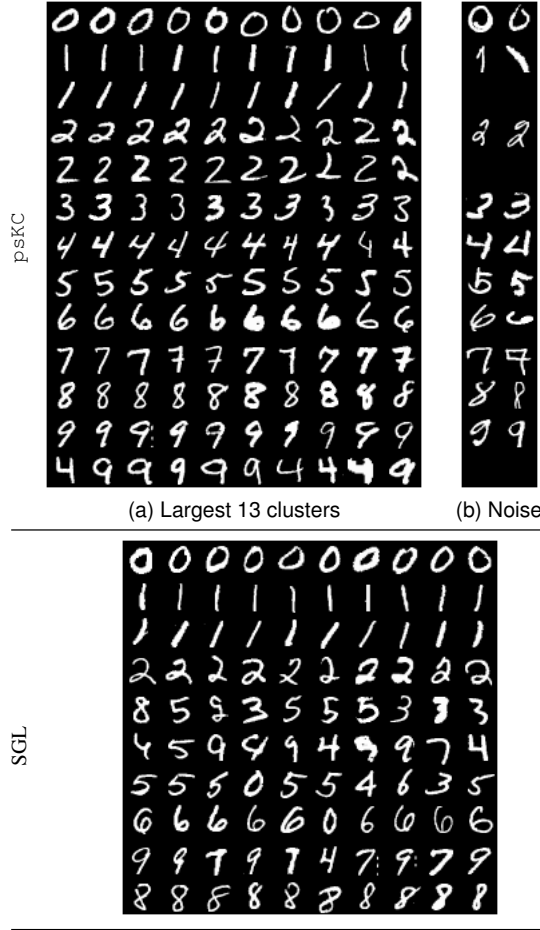
In contrast, SGL produces a clustering outcome which has less clear distinction between the background and landscape, and different portions of red stamps appear in separate clusters.

psKC and SGL have similar clustering outcomes on the van Gogh's painting, shown in Table 4.

None of DP and kernel k-means can complete in reasonable time on both images. Each of these images has a total of more than 1 million pixels. DP was unable to load the dataset on a machine with 256GB of main memory because of high memory

TABLE 5

Clustering outcome of psKC and SGL on the MNIST70k data set. The columns from left to right show the most similar digit to the least similar digit from each of the 10 equal-frequency bins in each cluster, sorted by \hat{K} in Definition 1 for psKC & by distance for SGL.



requirement using Matlab. Scalable kernel k-means took > 4 days (we terminated the run as it took too long to complete.)

6.1.3 Finding groups in a collection of handwritten images

Here we show the groups of handwritten digits found by psKC on the MNIST70k dataset which has 70,000 images.

The largest 13 clusters produced by psKC cover 92% of the images in the dataset. Table 5 shows examples of these clusters in (a). Subfigure (b) in Table 5 shows the sample noise images which do not belong to large clusters. Notice that the noise digits have different handwritten styles from those in the large clusters.

It is interesting to note that each of digits 1 & 2 has been grouped into two clusters, where each cluster has its own written style, e.g., digit 1 has a vertical written style in the first cluster; and a slant style in the second cluster.

Also note that digits 4 & 9 are grouped into three clusters. In addition of the two pure clusters, the third cluster consists of both digits of 1:3 proportion. This is in contrast to the result produced by a kNN-graph based clustering algorithm RCC³ (see Fig.3 in [18]), where both digits 4 & 9 have been grouped into a single

3. RCC was unable to complete this task in 3 days on our machine; whereas psKC ($\psi = 5000$) took < 25 minutes. See a brief description of RCC and its clustering outcomes of the artificial datasets in Appendix D.

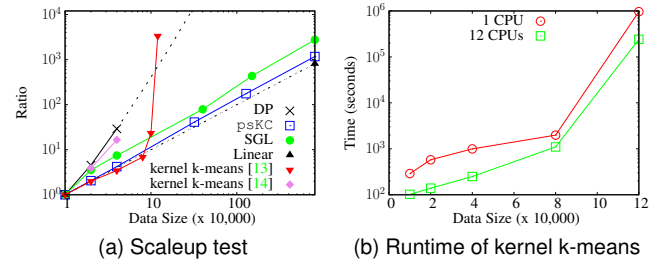


Fig. 4. Scaleup test on MNIST8M. The base in computing the ratio is the runtime at 10k points. DP crashed at 320K points. ‘Linear’ indicates the runtime of a linear-time algorithm in (a).

cluster. Also notice that the digits grouped in the third cluster have different written styles to those in the two pure clusters of 4 & 9.

Table 5 also shows the clustering outcome of SGL. Although SGL produces pure clusters for digits 0,1,2 and 8, it has mixed up multiple digits into individual clusters for other digits. Setting $k = 13$ produced worse result.

DP and the two versions of kernel k-means [13], [14] were unable to complete this task in reasonable time.

Summary: Only psKC and SGL could complete the clustering of all datasets/images used in the experiments. However, SGL produces poorer clustering outcomes than psKC on all datasets/images, except two on which they have similar outcomes.

6.2 Scaleup test

The result of a scaleup test using the MNIST8M dataset, which has a total of 8.1 million data points with 784 dimensions, is given in Figure 4(a). As expected, psKC has runtime linear to data set size, as the data size increases from 10k to 8.1 million (a factor of 810), its runtime increases by a factor of 1,171 (including both feature mapping (GPU) and clustering (CPU)). In contrast, as the data size increases from 10k to 40k (a factor of 4), DP’s runtime increases by a factor of 29; and kernel k-means using kNN kernel [14] by a factor of 17. Note that beyond 40k points, DP took too long to run and kernel k-means using kNN kernel had memory error; and the dotted line of DP in Figure 4(a) is a projected line beyond 40k points. The k-means-based SGL is the only contender which runs linearly, but it is still slower than psKC .

Kernel k-means [13] has two main components⁴: Nyström approximation [10] (to produce a finite dimensional feature map from a kernel of intractable dimensionality) and k-means.

Using a supercomputer Cray XC40 system with 1632 compute nodes, each has two 2.3GHz 16-core Haswell processors and 128GB of DRAM, the authors [13] reported an experiment using scalable kernel k-means with $s = 20$ on the MNIST8M dataset as follows:

The parallelization reduces the runtime of Nyström, but it increases the runtimes of PCA and k-means. Thus, the net speedup is significantly less, e.g., increasing the number of compute nodes 16 times from 8 to 128, the net speedup is less than 4 times. It is even less when the number of target dimensions s is increased. See Table 4 in [13] for details. This shows that parallelization

4. Dimensionality reduction using PCA is often an intermediary between feature transformation and k-means in implementations. PCA is used in order to produce the desired rank-restricted Nyström approximation, as required in their formulation (see Algorithm 2 in [13]). However, some exposition may have omitted PCA, though it is used in the code, e.g., [14].

alone has diminishing payoff as the complexity of the problem increases. In addition, parallelization needs to increase massively as the data size increases in order to complete in reasonable time.

Our experimental result in Figure 4(a) shows the same behaviour on a single-CPU machine. Scalable kernel k-means has similar runtime ratios as psKC up to 80k. But its runtime began to dramatically increase at 120k—the runtime increase is now more than 3000 times on a 12-fold increase in data size!

Even with 12 CPUs, as shown in Figure 4(b), scalable kernel k-means took more than 240,000 seconds to complete the dataset of 120k points, while the 1-CPU machine took close to 1 million seconds (more than 11 days). This is a speedup of 4 times on a 12-fold increase in the number of CPUs. In other words, the parallelization works well in scalable kernel k-means only if the number of CPUs is sufficiently large such that each CPU works on a small data set. Otherwise, a dramatic increase in runtime is expected, as shown in Figure 4(a).

In contrast, the algorithmic advantage of psKC , together with the use of the proposed point-set kernel, allows it to run on a standard machine of single-CPU (for clustering) and GPU (for feature mapping in preprocessing). This enables the clustering to be run on a commonly available machine (with both GPU and CPU) to deal with large scale datasets.

In terms of real time: on the dataset with 40k data points, psKC took 73 seconds which consists of 58 GPU seconds for feature mapping and 15 CPU seconds for clustering. In contrast, DP took 541 seconds. The gap in runtime widens as data size increases: To complete the run on 8.1 million points, DP is projected to take 379 years! That would be 12 billion seconds which is six orders of magnitude slower than psKC 's 20 thousand seconds (less than 6 hours). The widening gap is apparent in Figure 4(a). SGL took 269 seconds on 40k points, and more than 200 thousand seconds (58 hours) on 8.1 million points.

For the dataset of 120k data points, psKC took 243 seconds; whereas scalable kernel k-means took close to a million seconds, both on a one-CPU machine.

As it is, there is no opportunity for DP to do feature mapping (where GPU could be utilised). While it is possible for kernel k-means to make use of GPU as in psKC , scalable kernel k-means's main restriction is PCA which has no efficient parallel implementation, to the best of our knowledge. The clustering procedures of both DP and psKC could potentially be parallelized; but this does not change their time complexities.

6.3 Stability analysis

k-means, used in kernel k-means and SGL, is a randomised algorithm; and its clustering outcome is influenced by the initial random centers. While the psKC procedure is deterministic, the Isolation Kernel employed is generated based on random samples. Here we examine the stability of these algorithms.

Figure 5 shows the stability of the clustering outcomes in terms of F1 score [19], [20] over 10 trials, presented in box plots.

The results show that kernel k-means produced clustering which have variance much higher than those produced by psKC on the middle 50% results (showed as the box—small (large) box has low (high) variance). Kernel k-means produced wild outliers (see the three points outside the box) as shown on S3. Despite having its best result (the top outlier) is better than all other results, its two worst results (the bottom two outliers) are significantly worse than all other results. Though SGL has lower variance than kernel

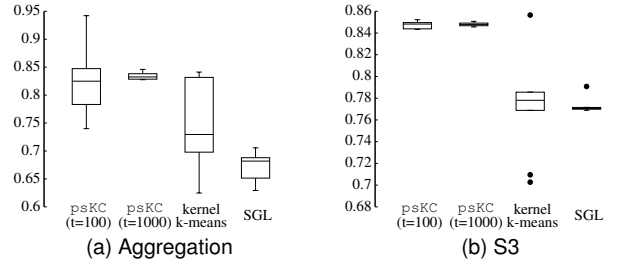


Fig. 5. Stability of clustering, presented in box plot based on 10 trials of the same parameter setting for each algorithm. The y-axis is F1 score.

TABLE 6
The effect of post-processing. Objective function $\Gamma(D)$ is Eq.4.

	psKC $\Gamma(D)$	post-processing		
		$\Gamma(D)$	#points re-assigned	time (%)
S3	1,492	1,492	0	2%
Autumn Colors	294,906	297,562	3,074	18%

k-means (because SGL chooses the best outcome from a number of k-means trials internally), its median F1 results (shown as the line inside the box) are worse than those of kernel k-means.

Overall, psKC ($t = 100$) produces higher F1 scores than kernel k-means and SGL. In addition, Figure 5 also shows that the variance can be significantly reduced by using a higher t at the cost of longer runtime. For example, on S3, psKC ($t = 100$) took 1.0 seconds, and psKC ($t = 1000$) took 9.2 seconds.

6.4 The effect of post-processing

Table 6 shows the effect of post-processing of psKC on two datasets. As mentioned in Section 5, the clustering outcome of psKC has already achieved a good approximation to the optimal maximization objective. The post-processing provides a final tweak to refine the approximation for the points having the lowest similarity only. On the small dataset S3, the post-processing did not make any re-assignment. On the large dataset Autumn Colors, the post-processing improved the objective function $\Gamma(D)$ by re-assigning three thousand points out of a million points, and it spent 18% of the total runtime of 10 seconds.

7 DISCUSSION

7.1 Isolation Kernel versus Gaussian Kernel

To get the clustering outcomes of psKC we showed here, it is crucial that the point-set kernel employs Isolation Kernel [4], [5] which is data dependent. Employing a Gaussian Kernel, which is data independent, psKC will perform poorly on datasets with clusters of non-globular shape, different data sizes and/or densities. This is because its similarity measurement is independent of data distribution. See the clustering outcomes of psKC_g which employs Gaussian kernel in Table 8 in Appendix D.

In addition to poor clustering outcomes, the use of Gaussian kernel in point-set kernel has high computational cost. This is because it has a feature map having intractable dimensionality.

7.2 Issue with running DP on a subsample

Our results show that density-based algorithm such as DP is a stronger clustering algorithm than kernel k-means, in terms of clustering outcomes. But DP is one of the most computationally

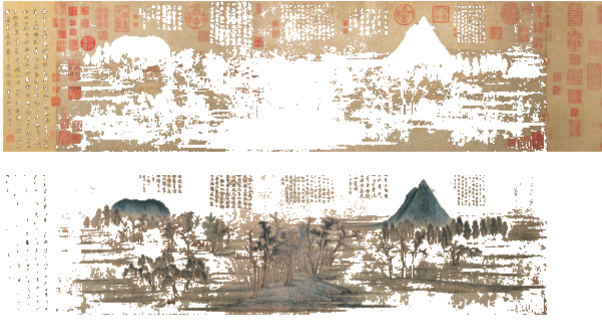


Fig. 6. Clustering outcome of DP on the reduced resolution of the Chinese painting: Autumn Colors.

expensive algorithms: it needs large memory space and its runtime is proportional to the square of data size (n^2).

It is possible to run DP on an image of high resolution by reducing its resolution first. However, the effect of this reduction can be counterproductive. For example, we have attempted to reduce the resolution of the Chinese painting: Autumn Colors, shown in Table 3. The resolution must be reduced from 1 million pixels to 90,000 pixels for DP to run on a machine with 256GB memory. As a result of this reduction, salient features of the painting were destroyed, e.g., some brush writings have become less obvious. With reduced density of pixels of brush writing, it is no more exhibited as a cluster of interest in its own right in the CIELAB space. Thus, DP is unable to identify the brush writing (together with the painting using the same brush color) as a cluster. An example clustering outcome of DP is shown in Figure 6. Note that the background and brush writing/painting are not clearly separated into two different clusters by DP, unlike psKC 's clustering outcomes shown in Table 3. In a nutshell, reducing the resolution would add another issue to DP's existing weakness in identifying clusters of uniform distribution as exemplified by the Ring-G dataset shown in Table 2.

7.3 Relation to Support Vector Clustering

In the framework of support vector machines, Support Vector Clustering (SVC) [21] describes a way to build a minimal enclosing sphere in feature space (induced by a kernel function) to capture the majority of the points in the given dataset. This approach is computationally expensive. Approximation methods have been developed to reduce its computational demands. For example, a subsampling method is proposed to reduce the learning cost in constructing the minimal enclosing sphere, and then employ Voronoi cells to assign points to different clusters [22]. This subsampling approximation method suffers from the same issue we described in the last subsection.

The only connection between psKC and SVC is the use of a kernel. Still, psKC employs a data dependent kernel called Isolation Kernel. Like all existing kernel-based methods, SVC employs a data independent kernel such as Gaussian kernel.

7.4 Some algorithmic weakness cannot be overcome by changing the distance function

A previous work [4] has shown that the clustering outcome of DBSCAN [2] can be improved by simply replacing Euclidean distance with Isolation Kernel. However, the resultant algorithm called MBSCAN [4], though performs better than DBSCAN on

the Ring-G and S3 datasets shown in Table 8 in Appendix D, MBSCAN still performs poorer than psKC on the S3 dataset.

This is another example of the limitation of existing clustering procedures: DBSCAN has some inherent algorithmic weakness that cannot be overcome by using a data dependent kernel.

7.5 Current approaches to runtime issue

A considerable amount of research has been delegated to mitigate the persistent longstanding runtime issue of using a point-to-point distance/kernel. For example, various indexing techniques [23], [24], [25] have been explored to reduce the n^2 runtime. While some have claimed to have achieved runtime proportional to n , this often comes with the cost of reduced task-specific performance [23], [24], [25].

For the point-to-point kernel-based methods, an alternative to indexing is to use kernel functional approximation [10], [26] to speed up their runtimes, e.g., the scalable kernel k-means [13] mentioned above.

Both approaches of indexing and kernel functional approximation reduce runtime by sacrificing task-specific performance. In contrast, psKC has its runtime proportional to n , **needing neither indexing nor kernel functional approximation**. As a result, its clustering outcomes are not compromised.

Parallelization is a way to distribute the workload, with/without any of the two approaches mentioned above. A massive parallelization enables an algorithm to run on a large dataset that was impossible with a single CPU. Given a fixed number of CPUs, it sets a limit on the data size that it can handle, as shown in Section 6.2 with scalable kernel k-means. A massive parallelization could also be used to scale up psKC , like any other algorithms. In other words, what we have presented in this paper is more fundamental than and orthogonal to parallelization.

7.6 Linear-time clustering algorithms

k-means [6] is the clustering algorithm of choice in many applications [15], including text clustering, graph clustering and hierarchical clustering (where k-means is the core clustering algorithm). This is because it often has linear runtime in practice [15], though its time complexity is superpolynomial in the worst case which requires $2^{\Omega(\sqrt{n})}$ iterations [11]. Being a linear-time kernel clustering algorithm, where the maximum number of iterations is fixed for some range of parameter settings, psKC offers an alternative to k-means in these and many other applications.

8 CONCLUSIONS

We show that the proposed clustering psKC outclasses DP, DBSCAN and two versions of kernel k-means as well as other variants such as SGL and RCC in terms of both clustering outcomes and runtime efficiency. We identify that the two root causes of shortcomings of existing clustering algorithms are (i) the use of *data independent* point-to-point distance/kernel (where the kernel has a feature map with intractable dimensionality) to compute the required similarity; and (ii) the algorithmic designs that constrict the types of clusters that they can identify. For example, kernel k-means can be interpreted as already using a point-set kernel. Yet, because its use is restricted to represent cluster centers, only clusters of globular shape and approximately the same size can be detected in feature space; and this does not guarantee that clusters of non-globular shape and different sizes in input space can be

detected. This is also the cause of its sensitivity to outliers. Both root causes have led to poor clustering outcomes.

The first root cause is the source of the longstanding runtime issue in density-based clustering algorithms that has prevented them from handling large scale datasets.

We address these root causes by using a *data dependent point-set kernel* and a new clustering algorithm which utilizes the point-set kernel as a medium to grow and characterize clusters—enabling clusters of arbitrary shapes and sizes to be detected and they are insensitive to outliers. As a result, psKC is the only clustering algorithm which is both effective and efficient—a quality which is all but nonexistent in current clustering algorithms. It is also the only kernel-based clustering that has runtime proportional to data size.

We show that the number of iterations in psKC is independent of data size. The parameter search ranges in psKC can be set to bound the number of iterations; this gives the linear time complexity in practice, enabling psKC to deal with large datasets.

ACKNOWLEDGEMENT

Kai Ming Ting is supported by Natural Science Foundation of China (62076120).

REFERENCES

- [1] A. Rodriguez and A. Laio, “Clustering by fast search and find of density peaks,” *Science*, vol. 344, no. 6191, pp. 1492–1496, 2014.
- [2] M. Ester, H.-P. Kriegel, J. S., and X. Xu, “A density-based algorithm for discovering clusters in large spatial databases with noise,” in *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96)*. AAAI Press, 1996, pp. 226–231.
- [3] Y. Zhu, K. M. Ting, and M. J. Carman, “Density-ratio based clustering for discovering clusters with varying densities,” *Pattern Recognition*, vol. 60, pp. 983–997, 2016.
- [4] X. Qin, K. M. Ting, Y. Zhu, and V. C. S. Lee, “Nearest-neighbour-induced isolation similarity and its impact on density-based clustering,” in *Proceedings of The Thirty-Third AAAI Conference on Artificial Intelligence*, 2019, pp. 4755–4762.
- [5] K. M. Ting, Y. Zhu, and Z.-H. Zhou, “Isolation kernel and its effect on SVM,” in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2018, pp. 2329–2337.
- [6] J. MacQueen, “Some methods for classification and analysis of multivariate observations,” in *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Statistics*. Berkeley, Calif.: University of California Press, 1967, pp. 281–297.
- [7] A. Smola, A. Gretton, L. Song, and B. Schölkopf, “A hilbert space embedding for distributions,” in *Algorithmic Learning Theory*, M. Hutter, R. A. Servedio, and E. Takimoto, Eds. Springer Berlin Heidelberg, 2007, pp. 13–31.
- [8] K. Muandet, K. Fukumizu, B. Sriperumbudur, and B. Schölkopf, “Kernel mean embedding of distributions: A review and beyond,” *Foundations and Trends in Machine Learning*, vol. 10 (1–2), pp. 1–141, 2017.
- [9] N. Cristianini and J. Shawe-Taylor, *An Introduction to Support Vector Machines and Other Kernel Based Learning Methods*. Cambridge University Press, 2000.
- [10] C. K. I. Williams and M. Seeger, “Using the Nyström method to speed up kernel machines,” in *Advances in Neural Information Processing Systems 13*, T. K. Leen, T. G. Dietterich, and V. Tresp, Eds., 2001, pp. 682–688.
- [11] D. Arthur and S. Vassilvitskii, “How slow is the k-means method?” in *Proceedings of the Twenty-second Annual Symposium on Computational Geometry*, 2006, pp. 144–153.
- [12] J. A. Hartigan and M. A. Wong, “Algorithm as 136: A k-means clustering algorithm,” *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, vol. 28, no. 1, pp. 100–108, 1979.
- [13] S. Wang, A. Gittens, and M. W. Mahoney, “Scalable kernel k-means clustering with nyström approximation: Relative-error bounds,” *Journal of Machine Learning Research*, vol. 20, no. 12, pp. 1–49, 2019.
- [14] D. Marin, M. Tang, I. B. Ayed, and Y. Boykov, “Kernel clustering: density biases and solutions,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 1, pp. 136–147, 2019.
- [15] C. C. Aggarwal, *Data Mining: The Textbook*. Springer, 2015.
- [16] Z. Kang, X. Z. Zhiping Lin, and W. Xu, “Structured graph learning for scalable subspace clustering: From single view to multiview,” *IEEE Transactions on Cybernetics*, 2021. [Online]. Available: [10.1109/TCYB.2021.3061660](https://doi.org/10.1109/TCYB.2021.3061660)
- [17] P.-N. Tan, M. Steinbach, A. Karpatne, and V. Kumar, *Introduction to Data Mining*. Second Edition, Pearson, 2018.
- [18] S. A. Shah and V. Koltun, “Robust continuous clustering,” *Proceedings of the National Academy of Sciences*, vol. 114, no. 37, pp. 9814–9819, 2017.
- [19] B. Larsen and C. Aone, “Fast and effective text mining using linear-time document clustering,” in *Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining*, 1999, pp. 16–22.
- [20] R. M. Aliguliyev, “Performance evaluation of density-based clustering methods,” *Information Sciences*, vol. 179, no. 20, pp. 3583–3602, 2009.
- [21] A. Ben-Hur, D. Horn, H. T. Siegelmann, and V. Vapnik, “Support vector clustering,” *Journal of Machine Learning Research*, vol. 2, pp. 125–137, 2001.
- [22] K. Kim, Y. Son, and J. Lee, “Voronoi cell-based clustering using a kernel support,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, no. 4, pp. 1146–1156, 2015.
- [23] A. Andoni and I. Razenshteyn, “Optimal data-dependent hashing for approximate near neighbors,” in *Proceedings of the Forty-seventh Annual ACM Symposium on Theory of Computing*, 2015, pp. 793–801.
- [24] E. Schubert, J. Sander, M. Ester, H. P. Kriegel, and X. Xu, “DbSCAN revisited, revisited: Why and how you should (still) use dbSCAN,” *ACM Transaction on Database Systems*, vol. 42, no. 3, pp. 19:1–19:21, 2017.
- [25] T. Liu, C. Rosenberg, and H. A. Rowley, “Clustering billions of images with large scale nearest neighbor search,” in *Proceedings of the IEEE Workshop on Applications of Computer Vision*, 2007.
- [26] T. Yang, Y.-F. Li, M. Mahdavi, R. Jin, and Z.-H. Zhou, “Nyström method vs random fourier features: A theoretical and empirical comparison,” in *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1*, ser. NIPS’12. USA: Curran Associates Inc., 2012, pp. 476–484.
- [27] F. T. Liu, K. M. Ting, and Z.-H. Zhou, “Isolation forest,” in *Proceedings of the IEEE International Conference on Data Mining*, 2008, pp. 413–422.
- [28] I. S. Dhillon, Y. Guan, and B. Kulis, “Kernel k-means: Spectral clustering and normalized cuts,” in *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2004, pp. 551–556.
- [29] F. R. Bach and M. I. Jordan, “Learning spectral clustering,” in *Advances in Neural Information Processing Systems 16*, S. Thrun, L. K. Saul, and B. Schölkopf, Eds., 2004, pp. 305–312.
- [30] C. E. Priebe, Y. Park, J. T. Vogelstein, J. M. Conroy, V. Lyzinski, M. Tang, A. Athreya, J. Cape, and E. Bridgeford, “On a two-truths phenomenon in spectral graph clustering,” *Proceedings of the National Academy of Sciences*, vol. 116, no. 13, pp. 5995–6000, 2019.
- [31] I. S. Dhillon, Y. Guan, and B. Kulis, “Weighted graph cuts without eigenvectors a multilevel approach,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 29, no. 11, pp. 1944–1957, 2007.
- [32] M. Tang, D. Marin, I. Ben Ayed, and Y. Boykov, “Kernel cuts: Kernel and spectral clustering meet regularization,” *International Journal of Computer Vision*, vol. 127, no. 5, p. 477–511, May 2019.
- [33] P. Fränti and S. Sieranoja, “K-means properties on six clustering benchmark datasets,” *Applied Intelligence*, vol. 48, no. 12, pp. 4743–4759, 2018.
- [34] G. Loosli, S. Canu, and L. Bottou, “Training invariant support vector machines using selective sampling,” in *Large Scale Kernel Machines*, L. Bottou, O. Chapelle, D. DeCoste, and J. Weston, Eds. Cambridge, MA: MIT Press, 2007, pp. 301–320.
- [35] E. Frank, M. A. Hall, and I. H. Witten, *The WEKA Workbench. Online Appendix for ‘Data Mining: Practical Machine Learning Tools and Techniques’*, 4th ed. Morgan Kaufmann, 2016.
- [36] M. J. Black and A. Rangarajan, “On the unification of line processes, outlier rejection, and robust statistics with applications in early vision,” *International journal of computer vision*, vol. 19, no. 1, pp. 57–91, 1996.
- [37] M. Brito, E. Chavez, A. Quiroz, and J. Yukich, “Connectivity of the mutual k-nearest-neighbor graph in clustering and outlier detection,” *Statistics & Probability Letters*, vol. 35, no. 1, pp. 33–42, 1997.
- [38] B. Larsen and C. Aone, “Fast and effective text mining using linear-time document clustering,” in *Proceedings of the 20th International Conference on World Wide Web*, 2011, pp. 577–586.

APPENDIX A

ISOLATION KERNEL

We provide the pertinent details of Isolation Kernel in this section. Other details can be found in [4], [5].

Let $D = \{x_1, \dots, x_n\}, x_i \in \mathbb{R}^d$ be a dataset sampled from an unknown probability density function $x_i \sim F$. Let $\mathbb{H}_\psi(D)$ denote the set of all partitionings H that are admissible under D , where each H covers the entire space of \mathbb{R}^d . Each of the ψ isolating partitions $\theta[z] \in H$ isolates one data point z from the rest of the points in a random subset $\mathcal{D} \subset D$, and $|\mathcal{D}| = \psi$, where $z \in \mathcal{D}$ has the equal probability of being selected from D .

Definition 4. For $x, y \in \mathbb{R}^d$, Isolation Kernel of x and y wrt D is defined to be the expectation taken over the probability distribution on all partitionings $H \in \mathbb{H}_\psi(D)$ that both x and y fall into the same isolating partition $\theta[z] \in H, z \in \mathcal{D}$:

$$\kappa_\psi(x, y | D) = \mathbb{E}_{\mathbb{H}_\psi(D)}[\mathbb{1}(x, y \in \theta[z] | \theta[z] \in H)] \quad (6)$$

where $\mathbb{1}(\cdot)$ is an indicator function.

In practice, κ_ψ is constructed using a finite number of partitionings $H_i, i = 1, \dots, t$, where H_i is created using $\mathcal{D}_i \subset D$:

$$\begin{aligned} \kappa_\psi(x, y | D) &= \frac{1}{t} \sum_{i=1}^t \mathbb{1}(x, y \in \theta | \theta \in H_i) \\ &= \frac{1}{t} \sum_{i=1}^t \sum_{\theta \in H_i} \mathbb{1}(x \in \theta) \mathbb{1}(y \in \theta) \end{aligned} \quad (7)$$

where θ is a shorthand for $\theta[z]$; and $\kappa_\psi(x, y | D) \in [0, 1]$.

Isolation Kernel is positive semi-definite as Eq 7 is a quadratic form, i.e., it defines a Reproducing Kernel Hilbert Space.

aNNE Implementation: As an alternative to using trees [27] in its first implementation of Isolation Kernel [5], a nearest neighbour ensemble (aNNE) has been used instead [4].

Like the tree method, the nearest neighbour method also produces each H model which consists of ψ isolating partitions θ , given a subsample of ψ points. Rather than representing each isolating partition as a hyper-rectangle, it is represented as a cell in a Voronoi diagram, where the boundary between two points is the equal distance from these two points.

H , being a Voronoi diagram, is built by employing ψ points in \mathcal{D} , where each isolating partition or Voronoi cell $\theta \in H$ isolates one data point from the rest of the points in \mathcal{D} . The point which determines a cell is regarded as the cell centre.

Given a Voronoi diagram H constructed from a sample \mathcal{D} of ψ points, the Voronoi cell centred at $z \in \mathcal{D}$ is:

$$\theta[z] = \{x \in \mathbb{R}^d \mid z = \underset{z \in \mathcal{D}}{\operatorname{argmin}} \|x - z\|_2\}.$$

Note that the boundaries of a Voronoi diagram are derived implicitly to be equal distance between any two points in \mathcal{D} ; and they need not be derived explicitly in realising Isolation Kernel.

A.1 Kernel distributions in contour plots

Figure 7 compares the contour plots of of Isolation Kernel and Gaussian Kernel. Notice that each contour line of Isolation Kernel, which denotes the same similarity to the centre (red point), is elongated along the sparse region and compressed along the dense region. In contrast, Gaussian kernel (or any data independent kernel) has the same symmetrical contour lines around the centre point, independent of data distribution (as shown in Figure 7(b)).

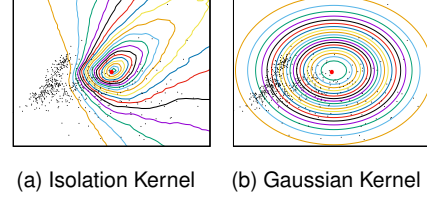


Fig. 7. Contour plots of Isolation Kernel and Gaussian Kernel on a real-world dataset. Black dots are data points.

A.2 Feature map of Isolation Kernel

Definition 5. Feature map of Isolation Kernel. For point $x \in \mathbb{R}^d$, the feature mapping $\Phi : x \rightarrow \{0, 1\}^{t \times \psi}$ of κ_ψ is a vector that represents the partitions in all the partitioning $H_i \in \mathbb{H}_\psi(D)$, $i = 1, \dots, t$; where x falls into only one of the ψ Voronoi cells in each partitioning H_i .

Given H_i , $\Phi_i(x)$ is a ψ -dimensional binary column vector representing all Voronoi cells $\theta_j \in H_i, j = 1, \dots, \psi$; where x falls into only one of the ψ Voronoi cells. The j -component of the vector is: $\Phi_{ij}(x) = \mathbb{1}(x \in \theta_j | \theta_j \in H_i)$. Given t partitionings, $\Phi(x)$ is the concatenation of $\Phi_1(x), \dots, \Phi_t(x)$.

A commonly used kernel such as Gaussian kernel is defined with a function; and its feature map is derived from this function. As a result, the kernel can be computed without a feature map.

Instead, Isolation Kernel, having no functional form, is defined based on a finite-dimensional feature map, defined from samples of a given dataset. The samples, as well as the resultant partitions, are products of randomized processes. Once they are determined in the preprocessing step, the feature map defines the Isolation Kernel *exactly* in its subsequent use in an algorithm.

APPENDIX B

POINT-SET KERNEL: DATA DEPENDENT PROPERTY

The data dependent property of point-set kernel follows directly from the data dependent property of Isolation Kernel. Its Lemma [4] is re-stated as follows:

Lemma 1. $\forall x, y \in \mathcal{X}_T$ (dense region) and $\forall x', y' \in \mathcal{X}_S$ (sparse region) such that $\forall z \in \mathcal{X}_S, z' \in \mathcal{X}_T, \rho(z) > \rho(z')$, the nearest neighbour-induced Isolation Kernel κ has the characteristic that for $\ell_p(x - y) = \ell_p(x' - y')$ implies $\kappa(x, y | D) < \kappa(x', y' | D)$.

Let \bar{x}_C and $\bar{x}_{C'}$ be the preimages of $\hat{\Phi}(C|D)$ and $\hat{\Phi}(C'|D)$, respectively; and $\rho(C) > \rho(C') \equiv \rho(\bar{x}_C) > \rho(\bar{x}_{C'})$.

Further let two reference points $x \in C$ and $x' \in C'$ such that $\ell(x, C) = \ell(x', C') \equiv \ell_p(x - \bar{x}_C) = \ell_p(x' - \bar{x}_{C'})$. We then have reference values: $\kappa(x, \bar{x}_C) < \kappa(x', \bar{x}_{C'})$; or for tiny distance $\ell(x, C)$: $\kappa(x, \bar{x}_C) \approx \kappa(x', \bar{x}_{C'})$.

From Lemma 1, because κ decreases at a faster rate in \mathcal{X}_T than that in \mathcal{X}_S as Δx increases (i.e., the probability of two points of equal inter-point distance falling into the same Voronoi cell is a monotonically decreasing function wrt the density of the cell, see the proof of the Lemma in [4]), the following holds:

$$\begin{aligned} \kappa(x + \Delta x, \bar{x}_C) &< \kappa(x' + \Delta x, \bar{x}_{C'}) \text{ \& } \\ [\kappa(x, \bar{x}_C) - \kappa(x + \Delta x, \bar{x}_C)] &> [\kappa(x', \bar{x}_{C'}) - \kappa(x' + \Delta x, \bar{x}_{C'})]. \end{aligned}$$

Because $\kappa(x, \bar{x}_C) \equiv \hat{K}(x, C)$,

$$\frac{\hat{K}(x, C) - \hat{K}(x + \Delta x, C)}{\Delta x} > \frac{\hat{K}(x', C') - \hat{K}(x' + \Delta x, C')}{\Delta x}.$$

This gives $\frac{d\hat{K}(x, C)}{dx} > \frac{d\hat{K}(x', C')}{dx}$, if $\ell(x, C) = \ell(x', C')$ and $\rho(C) > \rho(C')$.

APPENDIX C

KERNEL K-MEANS VERSUS SPECTRAL CLUSTERING

The close relationship between kernel k-means and spectral clustering is well established [28], [29]. Kernel k-means and many spectral clustering methods use the same k-means algorithm. In order to make them scalable for large datasets, kernel k-means usually employs a kernel functional approximation (e.g., Nyström [10]) as a preprocessing. An example is scalable kernel k-means [13], and it has a direct comparison with a spectral clustering that employs the same Nyström, utilizing a small set of ‘landmark’ points. SGL [16] is another example that employs a small set of landmark points to deal with large datasets, though the problem formulation differs. It has close relation to spectral clustering because of the use of the idea of graph embedding.

Spectral clustering can be viewed more broadly as embedding a graph into Euclidean space via an eigenvalue decomposition, before a clustering algorithm is applied. The embedding could be done in several ways; so as the clustering (though typically k-means, a clustering based on Gaussian Mixture model has also been explored, e.g., [30].)

Because of the eigenvalue decomposition, spectral clustering is usually more computationally expensive than kernel k-means. Converting a spectral clustering problem into a weighted kernel k-means problem is one way to reduce its time complexity [31]. A recent connection between spectral clustering and kernel k-means is made via regularisation [32].

APPENDIX D

EXPERIMENTAL SETTINGS & ADDITIONAL RESULTS

This section provides the details of the datasets/images used, the experimental settings and additional results.

Datasets used

- 1) [Artificial benchmark datasets](#): all benchmark datasets are from [33], except the Ring-G dataset which is our creation, and the AC dataset was first used in [14].
- 2) [Starry Night over the Rhone 2](#) (932 x 687)
- 3) [Zhao Mengfu's Autumn Colors](#) (2,005 x 500)
- 4) [MNIST data sets](#): images of handwritten digits, each is represented with 784 dimensions:
MNIST8M has 8.1 million images [34].
MNIST70k is the source dataset (having 70,000 images) from which the MNIST8M dataset is generated.

Experimental settings. We search parameters in each of the algorithms, i.e., DP, scalable kernel k-means, kernel k-means and psKC ; and report their best clustering outcomes after the search.

We implemented psKC in C++. Scalable kernel k-means is implemented in Scala as part of the Spark framework [13]; DBSCAN is implemented in Java as part of the WEKA framework [35]; RCC is implemented in Python [18]; and DP, DP_{ik} , k-means and kNN kernel are implemented in Matlab [1], [14].

The parameter search ranges used in the experiments are:

- DP: ϵ (the bandwidth used for density estimation) is in $[0.001m, 0.002m, \dots, 0.4m]$ where m is the maximum pairwise distance.
- kernel k-means [14]: k in kNN kernel is in $[0.01n, 0.02n, \dots, 0.99n]$; and the number of dimensions used is 100.
- Scalable kernel k-means [13]: σ in $[0.1, 0.25, 0.5, 1, \dots, 16, 24, 32]$; $s = 100$ (target #dimensions of the PCA step) and $c = 400$ (#dimensions of Nyström’s output), except for data set less than 400 points then it is $s = 20$ and $c = 200$.

- psKC : ψ in $[2, 4, 6, 8, 16, 24, 32]$, $\tau = 0.1$ and $\varrho = 0.1$ for the images, except MNIST70k: $\psi = 5000$, $\tau = 0.001$. For artificial datasets, ψ in $[55, 70, 128, 256, 512, 1024]$, τ in $[0.1, 0.2, 0.5, 1, 2, 5, 10, 20, 50, 100, 800] \times 10^{-4}$, $\varrho = 0.1$ (except in Ring-G, where $\varrho = 0.26$ was used.)
- psKC_g : $\gamma = 2^i$ where i in $[1, 2, 3, \dots, 16]$, $\tau = 0.1$ and ϱ in $[0.1, 0.01, 0.001, \dots, 1 \times 10^{-10}]$.
- DBSCAN: ϵ in $[0.001m, 0.002m, \dots, 0.999m]$ and MinPts in $[2, 3, \dots, 30]$, where m is the maximum pairwise distance and MinPts is the density threshold.
- MBSCAN: $\psi = 32$, $t = 200$, ϵ in $[0.001m, 0.002m, \dots, 0.999m]$ and MinPts in $[2, 3, \dots, 30]$, where m is the maximum pairwise distance and MinPts is the density threshold.
- DP_{ik} : For DP, ϵ is in $[0.001m, 0.002m, \dots, 0.4m]$ where m is the maximum pairwise distance. Setting for Isolation Kernel: ψ in $[2, 4, 6, 8, 16, 24, 32]$.
- RCC: k is in $[1, 2, \dots, 200]$ and τ is in $[1, 1.1, \dots, 20]$.
- SGL [16]: $\#anchor$ in $[25, 30, 33, 36, 40, 50]$, α in $[0.1, 1, 10, 50]$ and β in $[0.0001, 0.001]$

All algorithms were given the advantage by setting the number of clusters to the true cluster number, except DBSCAN, MBSCAN, RCC and psKC which find the number of clusters automatically through their parameter settings. $t = 100$ is the default setting for Isolation Kernel for all experiments, except those in Section 6.3.

The experiments were run on a Linux CPU machine: AMD 16-core CPU with each core running at 2.0 GHz and 32 GB RAM. The feature space conversion was executed on a machine having GPU: 2 x GTX 1080 Ti with each card having 12 GB RAM.

Guide for psKC parameter tuning: In addition to $\varrho = 0.1$ which can usually be set as default, the following two rules may be used to tune psKC :

1. If clusters are joined but need to be split, then increase ψ .
2. If clusters are split but need to be joined, then decrease τ .

Parameter settings used by psKC in to produce the clustering outcomes shown in Tables 2 to 5 are listed in Table 7.

Source code of psKC is available at

<https://github.com/IsolationKernel/Codes/tree/main/IDK>

Additional results: Table 8 shows the clustering outcomes of other clustering algorithms, in addition to those used in the main text. The additional algorithms are DBSCAN [2], MBSCAN [4], RCC⁵ [18], k-means [6], SGL [16], DP_{ik} and psKC_g . Both MBSCAN and DP_{ik} employ Isolation Kernel. psKC_g employs the following point-set kernel, where κ_g is the Gaussian kernel:

$$\hat{K}_g(x, G) = \frac{1}{|G|} \sum_{y \in G} \kappa_g(x, y)$$

The result in Table 8 shows that MBSCAN which employs IK has the best clustering outcomes. It is the only one among these algorithms that correctly cluster four out of the five datasets. But, MBSCAN is still worse than psKC in terms of clustering outcomes and time complexity.

5. Robust continuous clustering [18] optimises a continuous and differentiable objective based on the duality between robust estimation and line processes [36]. Its objective is applied on the connectivity of a mutual k -nearest neighbours (m - k NN) graph [37] and is optimised using a linear least-squares solver. RCC has the ability to detect arbitrarily shaped clusters based on the m - k NN graph, but its computational complexity is dominated by the m - k NN graph construction which costs n^2 [38]. Thus, although the solver used in RCC scales to large datasets, RCC is still unable to deal with large datasets.

TABLE 7

The psKC parameters used to generate the clustering outcomes in Tables 2, 3, 4 & 5. The random seed was set to 42, $t = 100$ and $\sigma = 0.1$ for all experiments, except $\sigma = 0.26$ for Ring-G.

Table	Dataset	ψ	τ	Table	Dataset	ψ	τ
2	Ring-G	128	2×10^{-3}	3	Autumn Colors	8	0.1
	AC	256	2×10^{-4}	4	Red stamps (Class 7)	24	0.1
	Aggreg	128	1×10^{-2}		Starry Night	16	
	Spiral	512	1×10^{-4}	5	MNIST70k	5,000	1×10^{-3}
	S3	70	8×10^{-2}				

TABLE 8

Artificial datasets: Clustering outcomes of k-means, DBSCAN, psKC_g which employs Gaussian Kernel, MBSCAN & DP_{ik} which employ Isolation Kernel and RCC. The results with yellow frames indicate good clustering outcomes; and those without have poor clustering outcomes.

