

The Laplacian K -Modes Algorithm for Clustering

Weiran Wang Miguel Á. Carreira-Perpiñán

Electrical Engineering and Computer Science, University of California, Merced

<http://eecs.ucmerced.edu>

Jun 15, 2014

Abstract

In addition to finding meaningful clusters, centroid-based clustering algorithms such as K -means or mean-shift should ideally find centroids that are valid patterns in the input space, representative of data in their cluster. This is challenging with data having a nonconvex or manifold structure, as with images or text. We introduce a new algorithm, Laplacian K -modes, which naturally combines three powerful ideas in clustering: the explicit use of assignment variables (as in K -means); the estimation of cluster centroids which are modes of each cluster's density estimate (as in mean-shift); and the regularizing effect of the graph Laplacian, which encourages similar assignments for nearby points (as in spectral clustering). The optimization algorithm alternates an assignment step, which is a convex quadratic program, and a mean-shift step, which separates for each cluster centroid. The algorithm finds meaningful density estimates for each cluster, even with challenging problems where the clusters have manifold structure, are highly nonconvex or in high dimension. It also provides centroids that are valid patterns, truly representative of their cluster (unlike K -means), and an out-of-sample mapping that predicts soft assignments for a new point.

1 Introduction

Given a dataset $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^D$, centroid-based clustering algorithms such as K -means (Bishop, 2006) and mean-shift (Fukunaga and Hostetler, 1975; Cheng, 1995; Carreira-Perpiñán, 2000; Comaniciu and Meer, 2002) estimate a representative $\mathbf{c}_k \in \mathbb{R}^D$ of each cluster k in addition to assigning data points to clusters. Besides finding meaningful clusters, we would ideally like to find centroids that are valid patterns in the input space, representative of data in their cluster. This is challenging with data having a nonconvex or manifold structure, as with images or text. Fig. 1 illustrates this with a single cluster consisting of continuously rotated digit-1 images. Since these images represent a nonconvex cluster in the high-dimensional pixel space, their mean (which averages all orientations) is not a valid digit-1 image, which makes the centroid not interpretable and hardly representative of a digit 1. Mean-shift does not work well either: to produce a single mode, a large bandwidth is required, which makes the mode lie far from the manifold; a smaller bandwidth does produce valid digit-1 images, but then multiple modes arise for the same cluster, and under mean-shift they define each a cluster.

Forcing the centroids to be exemplars is often regarded as a way to ensure the centroids are valid patterns. Although there exist exemplar-based or K -medoids clustering algorithms (Kaufman and Rousseeuw, 1990; Bishop, 2006; Hastie et al., 2009) which constrain centroids to be points from the dataset (“exemplars”) and often minimize a K -means type of objective function with a possibly non-Euclidean distance, such algorithms are typically slow because updating centroid \mathbf{c}_k requires testing all pairs of points in cluster k . Besides, the exemplars themselves are often noisy and thus not that representative of their neighborhood.

Remarkably, no algorithm that partitions the data using exactly K meaningful modes existed in the literature until the K -modes algorithm proposed by Carreira-Perpiñán and Wang (2013). This combines the idea of clustering through binary assignment variables with the idea that high-density points are representative of a cluster. Each centroid found by the K -modes algorithm is the mode of a kernel density estimate defined by data points in each cluster. As a result, the centroids average out noise or idiosyncrasies that exist in individual data points and are representative of their cluster and neighborhood. This can be seen from the

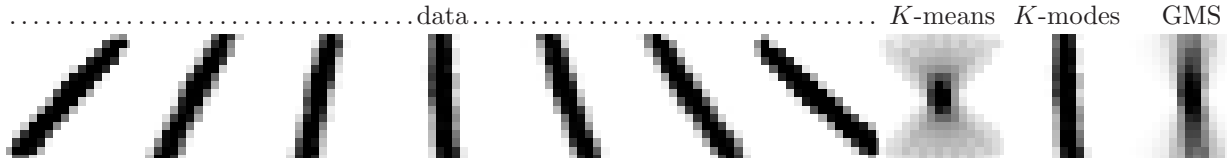


Figure 1: A cluster of 7 rotated-1 USPS digit images and the centroids found by K -means, K -modes (both with $K = 1$) and mean-shift (with σ so there is one mode). Example taken from Carreira-Perpiñán and Wang (2013).

K -modes centroid for the rotated digit-1 problem in Fig. 1. K -modes was also shown to have nice properties such as being more robust to mis-specification of the bandwidth and to outliers. Its optimization procedure is also very efficient.

One important disadvantage of K -modes is that it uses the same assignment rule as K -means (each point is assigned to its closest centroid in Euclidean distance), so it can only find convex clusters (a Voronoi tessellation). Therefore, like K -means, it cannot handle clusters with nonconvex shapes or manifold structure, unlike mean-shift or spectral clustering (Shi and Malik, 2000). The main contribution of this paper is to solve this issue, while keeping the nice properties that K -modes does have. The key idea is to modify the K -modes objective function such that the assignment rule becomes much more flexible. We then give an alternating optimization procedure to find the assignments and the modes. The resulting *Laplacian K -modes* algorithm is able to produce for each cluster a nonparametric density and a mode as valid representative (like K -modes), to separate nonconvex shaped clusters (like mean-shift and spectral clustering), and to give soft assignment of data points to each cluster. Yet, all of these merits are achieved at a reasonable computational cost, and the algorithm works well with high-dimensional data.

2 Related work

2.1 Centroids-based algorithms

Given the number of clusters K , K -means (Bishop, 2006) minimizes the objective

$$\begin{aligned} \min_{\mathbf{Z}, \mathbf{C}} \quad & \sum_{k=1}^K \sum_{n=1}^N z_{nk} \|\mathbf{x}_n - \mathbf{c}_k\|^2 \\ \text{s.t.} \quad & z_{nk} \in \{0, 1\}, \sum_{k=1}^K z_{nk} = 1, n = 1, \dots, N \end{aligned} \quad (1)$$

where $\mathbf{Z} = (z_{nk})$ are binary assignment variables (of point n to cluster k) and $\mathbf{C} = (\mathbf{c}_1, \dots, \mathbf{c}_K)$ are centroids living in \mathbb{R}^D . At an optimum, centroid \mathbf{c}_k is the mean of the points in its cluster.

Given a bandwidth $\sigma > 0$, Gaussian mean-shift (Fukunaga and Hostetler, 1975; Cheng, 1995; Carreira-Perpiñán, 2000; Comaniciu and Meer, 2002) defines a kernel density estimate (kde)

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N G(\|\mathbf{x} - \mathbf{x}_n\|/\sigma)^2 \quad (2)$$

with kernel $G(t) \propto e^{-t/2}$, and applies the iteration (started from each data point):

$$p(n|\mathbf{x}) = \frac{\exp(-\frac{1}{2}\|\mathbf{x} - \mathbf{x}_n\|/\sigma)^2}{\sum_{n'=1}^N \exp(-\frac{1}{2}\|\mathbf{x} - \mathbf{x}_{n'}\|/\sigma)^2}, \quad \mathbf{x} \leftarrow \mathbf{f}(\mathbf{x}) = \sum_{n=1}^N p(n|\mathbf{x})\mathbf{x}_n \quad (3)$$

which converges to a mode (local maximum) of p from nearly any initial \mathbf{x} (Carreira-Perpiñán, 2007). Each mode is the centroid for one cluster, which contains all data points that converge to its mode. The user parameter is the bandwidth σ , which determines the resulting number of clusters implicitly.

Both algorithms have well-known pros and cons. K -means tends to define round clusters; mean-shift can obtain clusters of arbitrary shapes and has been very popular in low-dimensional clustering applications

such as image segmentation (Comaniciu and Meer, 2002), but does not work well in high dimension due to the scarcity of data and difficulty in obtaining a good kde. Both algorithms suffer from outliers, which can move centroids outside their cluster in K -means or create singleton modes in mean-shift. Computationally, K -means is much faster than mean-shift, at $\mathcal{O}(KND)$ and $\mathcal{O}(N^2D)$ per iteration, respectively, particularly with large datasets (in fact, accelerating mean-shift has been a topic of active research, e.g. Carreira-Perpiñán (2006); Yuan et al. (2010)). Mean-shift does not require a value of K , which is sometimes convenient, but it is often desirable to force an algorithm to produce exactly K clusters (e.g. in figure-ground separation ($K = 2$) or in medical image analysis, where one may know or estimate the number of organs sought), which is not straightforward for mean-shift to achieve.

2.1.1 The K -modes algorithm

We now briefly summarize the original K -modes algorithm (Carreira-Perpiñán and Wang, 2013). Its objective function

$$\begin{aligned} \max_{\mathbf{Z}, \mathbf{C}} \quad & \sum_{n=1}^N \sum_{k=1}^K z_{nk} G\left(\left\|\frac{\mathbf{x}_n - \mathbf{c}_k}{\sigma}\right\|^2\right) \\ \text{s.t.} \quad & z_{nk} \in \{0, 1\}, \sum_{k=1}^K z_{nk} = 1, n = 1, \dots, N \end{aligned} \tag{4}$$

can be seen as the sum of a kde defined by each cluster separately. It is convenient to solve this problem with alternating optimization over \mathbf{Z} and \mathbf{C} . For fixed \mathbf{C} , the optimization over \mathbf{Z} decouples over each point, and \mathbf{x}_n is assigned to cluster $l = \arg \max_k G(\|(\mathbf{x}_n - \mathbf{c}_k)/\sigma\|^2) = \arg \min_k \|\mathbf{x}_n - \mathbf{c}_k\|$ due to the discrete constraints of the problem. For fixed \mathbf{Z} , the optimization over \mathbf{C} decouples over each cluster and we have a separate unconstrained maximization for each centroid, of the form $L(\mathbf{c}_k) = \sum_{n=1}^N z_{nk} G(\|(\mathbf{x}_n - \mathbf{c}_k)/\sigma\|^2)$, which is proportional to the cluster’s kde (this is why each centroid is truly a mode), and can be done with mean-shift updates as in eq. (3). The cost per outer iteration of this procedure is $\mathcal{O}(KND)$, which mainly comes from computing distances between data points and centroids. Since each step is strictly feasible and decreases the objective or leaves it unchanged, this converges to a local optimum in a finite number of outer-loop steps if the \mathbf{C} -step is exact.

2.2 Laplacian smoothing and learning soft assignments

Obtaining hard assignments by optimizing over a discrete cluster indicator matrix is usually difficult, because interesting objective functions are typically NP-hard. Spectral clustering algorithms (Shi and Malik, 2000; Yu and Shi, 2003) avoid this difficulty by first approximating the solution using eigenvectors of the normalized graph Laplacian. However, since the eigenvectors do not readily provide valid assignments, these algorithms need to run another clustering algorithm (usually K -means) on the eigenvectors to obtain actual partitions of the data—a post-processing step that is somewhat artificial and can introduce multiple local optima. In Laplacian K -modes, we relax \mathbf{Z} to be a stochastic matrix, so our \mathbf{Z} -step results from a convex QP and provides soft assignments of points to clusters, which may also be used as posterior probabilities.

Laplacian smoothing has also been used in combination with nonnegative matrix factorization (NMF) for clustering (Cai et al., 2011). NMF learns a decomposition of the input data matrix where both basis and coefficients are nonnegative, and tends to produce a parts-based representation of the data (Lee and Seung, 1999), though this is not always so. Cai et al. (2011) add a Laplacian smoothing term regarding the coefficient matrix to the NMF objective function, so that data points that are close in input space are encouraged to have similar representations using the basis set. K -means is then applied to the learned coefficients matrix to obtain a final partition of the data. Like spectral clustering, this algorithm does not directly optimize over the assignments, but obtains them in a post-processing step.

There has been recent work in clustering that directly optimizes over a stochastic assignment matrix. Arora et al. (2011) optimize over a stochastic matrix \mathbf{P} such that $\mathbf{P}\mathbf{P}^\top$ best approximates a rescaled similarity matrix. However, the optimization problem has multiple solutions which are related by rotations. Therefore, they propose to exploit the geometry of the problem using a rotation-based algorithm, which is straightforward for up to $K = 4$ clusters but requires an optimization procedure to computing the projection onto the probability simplex for $K > 4$ clusters. The idea of AnchorGraphs (Liu et al., 2010) is used by Yang and

Oja (2012) to approximate the affinities between data points through a two-step Data-Cluster-Data (DCD) random walk. They then minimize the generalized KL divergence between a given sparse affinity matrix and the affinities obtained from DCD. In this formulation, the matrix containing probabilities of points moving to the (augmented) cluster nodes is stochastic. These approaches are related to Laplacian K -modes in that they optimize some objective over the assignment probabilities. But our Laplacian K -modes algorithm also obtains prototypical centroids, does not have the issue of rotational equivalence of Arora et al. (2011), and makes use of the efficient projection onto the probability simplex to deal with any number of clusters.

3 Algorithm

3.1 The Laplacian K -modes algorithm

We change the assignment rule of K -modes to handle more complex shaped clusters based on two ideas: (1) the observation that *nearby data points should have similar assignments*; and (2) the use of *soft assignments*, which allows more flexibility in the clusters and simplifies the optimization. We first build a graph (e.g. k -nearest-neighbor graph) on the dataset, and let w_{mn} be an affinity (e.g. binary, heat kernel) between \mathbf{x}_m and \mathbf{x}_n . We then add to the K -modes objective function a Laplacian smoothing term $\frac{\lambda}{2} \sum_{m=1}^N \sum_{n=1}^N w_{mn} \|\mathbf{z}_m - \mathbf{z}_n\|^2$ to be minimized, where $\mathbf{z}_n = [z_{n1}, \dots, z_{nK}]^\top$ is the assignment vector of \mathbf{x}_n , $n = 1, \dots, N$, to each of the K clusters, and $\lambda \geq 0$ is a trade-off parameter. The assignments are now continuous variables, but constrained to be positive and sum to 1. Thus, z_{nk} can be considered as the probability of assigning \mathbf{x}_n to cluster k (soft assignment). Thus, the *Laplacian K -modes* objective function is:

$$\begin{aligned} \min_{\mathbf{Z}, \mathbf{C}} \quad & \frac{\lambda}{2} \sum_{m=1}^N \sum_{n=1}^N w_{mn} \|\mathbf{z}_m - \mathbf{z}_n\|^2 - \sum_{n=1}^N \sum_{k=1}^K z_{nk} G\left(\left\|\frac{\mathbf{x}_n - \mathbf{c}_k}{\sigma}\right\|^2\right) \\ \text{s.t.} \quad & \sum_{k=1}^K z_{nk} = 1, \quad n = 1, \dots, N, \\ & z_{nk} \geq 0, \quad n = 1, \dots, N, \quad k = 1, \dots, K. \end{aligned} \quad (5)$$

We can rewrite this objective in matrix form:

$$\begin{aligned} \min_{\mathbf{Z}, \mathbf{C}} \quad & \lambda \text{tr}(\mathbf{Z}^\top \mathbf{L} \mathbf{Z}) - \text{tr}(\mathbf{B}^\top \mathbf{Z}) \\ \text{s.t.} \quad & \mathbf{Z} \mathbf{1}_K = \mathbf{1}_N, \quad \mathbf{Z} \geq \mathbf{0} \end{aligned} \quad (6)$$

where $\mathbf{L} = \mathbf{D} - \mathbf{W}$ is the graph Laplacian for the affinity matrix $\mathbf{W} = (w_{nm})$ and degree matrix $\mathbf{D} = \text{diag}(\sum_{n=1}^N w_{mn})$, $\mathbf{B} = (b_{nk})$ is an $N \times K$ matrix containing data-centroid affinities $b_{nk} = G(\|\mathbf{x}_n - \mathbf{c}_k\|/\sigma)^2$, $n = 1, \dots, N$, $k = 1, \dots, K$, $\mathbf{1}_K$ is a K dimensional vector of 1s and \geq means elementwise comparison. Other variations of the graph Laplacian can also be used (e.g. the normalized Laplacian), see von Luxburg (2007). The constraint on \mathbf{Z} shows it is a stochastic matrix. We can obtain a hard clustering if desired by assigning each point to the cluster with highest assignment value.

3.1.1 Special cases of the hyperparameters (λ, σ)

In Laplacian K -modes, in addition to K there are two user parameters: λ controls the smoothness of the assignment, and σ controls the smoothness of the kde defined on each cluster. Consider first the case of $\lambda = 0$, where Laplacian K -modes becomes the original K -modes algorithm. Carreira-Perpián and Wang (2013) already noted that the K -modes algorithm has two interesting limit cases: it becomes K -means when $\sigma \rightarrow \infty$, and a form of K -medoids when $\sigma \rightarrow 0$, since the centroids are driven towards data points. In both cases the assignments are hard (1-out-of- K coding). The case when $\lambda \rightarrow \infty$ makes the first term in eq. (5) dominant and forces all connected points to have identical assignments, which is not interesting for the purpose of clustering. Therefore, the most interesting behavior of the algorithm is for intermediate λ . Finally, another interesting special case of Laplacian K -modes corresponds to $\lambda > 0$ and $\sigma \rightarrow \infty$, which we call *Laplacian K -means*, and which seems to be a new algorithm as well.

Algorithm 1 Accelerated gradient projection for the \mathbf{Z} step.

Input: Initial $\mathbf{Z}_0 \in \mathbf{R}^{N \times K}$, $s = \frac{1}{2\lambda M}$, $M =$ largest eigenvalue of graph Laplacian \mathbf{L} .

- 1: Set $\mathbf{Y}_1 = \mathbf{Z}_0$, $t_1 = 1$, $\tau = 1$.
- 2: **repeat**
- 3: Compute gradient at \mathbf{Y}_τ : $\mathbf{G}_\tau = 2\lambda\mathbf{L}\mathbf{Y}_\tau - \mathbf{B}$
- 4: $\mathbf{Z}_\tau =$ simplex projection of each row of $\mathbf{Y}_\tau - s\mathbf{G}_\tau$
- 5: $t_{\tau+1} = (1 + \sqrt{1 + 4t_\tau^2})/2$
- 6: $\mathbf{Y}_{\tau+1} = \mathbf{Z}_\tau + (\frac{t_\tau - 1}{t_{\tau+1}})(\mathbf{Z}_\tau - \mathbf{Z}_{\tau-1})$
- 7: $\tau = \tau + 1$
- 8: **until** convergence

Output: \mathbf{Z}_τ is the solution of the \mathbf{Z} -step.

3.2 Optimization procedure for Laplacian K -modes

To solve (5), we use alternating optimization over \mathbf{C} and \mathbf{Z} , which takes advantage of the problem’s structure.

C-step For fixed \mathbf{Z} , we are only concerned with the second term of (5) which is the K -modes objective. Therefore, our step over \mathbf{C} is identical to that of K -modes: it decouples over clusters and we apply mean-shift to solve for each \mathbf{c}_k separately. The cost of this step is $\mathcal{O}(KND)$.

Z-step Unlike in K -modes, our \mathbf{Z} -step no longer decouples, which means we have to solve for NK variables all together. Since the graph Laplacian \mathbf{L} is positive semidefinite, the problem over \mathbf{Z} is a convex quadratic program (QP). While we could apply a standard QP algorithm, such as an interior point method, we provide here an algorithm that is very simple (no parameters to set), efficient and that scales well to real problems where the number of points N or the number of clusters K is very large. The solution is based on the gradient proximal algorithm used by Beck and Teboulle (2009). Their general framework solves convex problems of the form $\min_{\mathbf{x}} f(\mathbf{x}) = g(\mathbf{x}) + h(\mathbf{x})$, where g is convex and has Lipschitz continuous gradient (with constant L), and h is convex but not necessarily differentiable. The gradient proximal algorithm iteratively updates the variables by first taking a gradient step of the first function and then projecting it with the second function, i.e., $\mathbf{x}_{\tau+1} = \arg \min_{\mathbf{y}} \frac{L}{2} \|\mathbf{y} - (\mathbf{x}_\tau - \frac{1}{L} \nabla g(\mathbf{x}_\tau))\|^2 + h(\mathbf{y})$. It can be proven that the algorithm converges in objective function value with rate $\mathcal{O}(1/\tau)$ (where τ is the iteration counter) with a *constant stepsize* $\frac{1}{L}$, and using Nesterov’s acceleration scheme improves the rate to $\mathcal{O}(1/\tau^2)$.

To apply this framework to our \mathbf{Z} -step, we make the identification that g is our smooth quadratic objective function, which has continuous gradient with $L = 2\lambda M$ being the (smallest) Lipschitz constant, where M is the largest eigenvalue of \mathbf{L} , and h is the indicator function of the probability simplex. Consequently, our proximal step is computing the Euclidean projection of the gradient step onto the probability simplex. Note that computing the Euclidean projection onto the K -dimensional simplex is itself a quadratic program. Fortunately, there exists an efficient algorithm which computes the exact projection with $\mathcal{O}(K \log K)$ time complexity (Duchi et al., 2008; Wang and Carreira-Perpiñán, 2013).

We provide the accelerated gradient projection algorithm for our \mathbf{Z} -step in Algorithm 1. Notice the graph Laplacian is sparse and its largest eigenvalue M can be obtained efficiently (e.g. by power iterations). Therefore the constant stepsize s can be easily determined right after constructing the graph Laplacian. Compared to a pure gradient projection algorithm, the additional computational effort of the acceleration scheme in maintaining an auxiliary sequence \mathbf{Y} (lines 5–6 of Algorithm 1) is minimal, and we clearly observe an improved convergence behavior in experiments.

Each iteration of Algorithm 1 costs $\mathcal{O}(NK\rho + NK \log K)$, where ρ is the neighborhood size in constructing \mathbf{L} (or the number of nonzero entries in each row). The first term accounts for computing the gradient and the second term accounts for projecting each row of \mathbf{Z} onto the probability simplex. Notice how, although it is solving a large QP, the cost per iteration of our \mathbf{Z} -step is independent of the input dimensionality D (in contrast, the \mathbf{C} -step has time complexity $\mathcal{O}(KND)$). Despite its sublinear convergence rate, the algorithm has a clear advantage in its simplicity: it does not require any line search or costly matrix operation, and it is very easy to implement.

3.2.1 Convergence properties

In the **C**-step, each mean-shift update increases the density of the cluster kde (or leaves it unchanged) and its convergence rate to a mode is linear in general (Carreira-Perpiñán, 2007). In the **Z**-step, the accelerated gradient projection converges theoretically at $\mathcal{O}(1/\tau^2)$ rate where τ is the iteration counter, although this algorithm seems to perform much better than the theoretical guarantee in practice (Beck and Teboulle, 2009). We alternate the **C** and **Z** steps until a convergence criterion is satisfied (e.g. the change to the variables is below some threshold). Notice both steps use iterative procedures, so the number of iterations depends on the convergence accuracy. Let ϵ_1 and ϵ_2 be the optimization precision for the **C**-step and **Z**-step, respectively, then one alternation of our algorithm costs $(D \log(1/\epsilon_1) + \rho/\sqrt{\epsilon_2})KN$. We found empirically that moderate accuracy and few iterations suffice for good clustering results, and our algorithm scales well due to its low per-iteration cost. In an efficient implementation, both steps can be inexact (e.g. each could run for a fixed, small number of iterations). Since the **Z**-step algorithm is feasible, exiting it early produces valid assignments.

3.2.2 Homotopy algorithm

As with K -means and K -modes, the Laplacian K -modes objective function has local optima, which are caused by the nonlinear, kde term. One strategy to find a good optimum consists of first finding a good optimum for K -means and then run a homotopy algorithm initialized there. We can construct a homotopy by varying continuously λ from 0 and σ from ∞ , which corresponds to K -means, to their target values (λ^*, σ^*) . In practice, we follow this path approximately, by running some iterations of the fixed- (λ, σ) Laplacian K -modes algorithm for each value of (λ, σ) . As is well known with homotopy techniques, this tends to find better optima than starting directly at the target value (λ^*, σ^*) . A good optimum for K -means can be obtained by picking the best of several random restarts, or by using the K -means++ initialization strategy, which has approximation guarantees (Arthur and Vassilvitskii, 2007).

3.2.3 Setting the hyperparameters

We believe that, in an unsupervised setting, the user should be able to explore different scenarios and so the algorithm should have a small number of intuitive hyperparameters to control this—rather than automatically guessing, say, the number of clusters, which often is not uniquely defined for a dataset. Laplacian K -modes has 3 intuitive hyperparameters, which allow a user to explore different scenarios: more or less clusters (K), different scales (σ), and degree of membership smoothness (λ). Typical values are around $\lambda = 1$, which allows some amount of propagation among neighbors and thus nonconvex clusters, and σ obtained from a kde bandwidth formula (Wand and Jones, 1994), such as the average distance to the 7th nearest neighbor gives a reasonable density (Zelnik-Manor and Perona, 2005). (An even better option is to use an adaptive kde, where each point has a different bandwidth, obtained using “entropic affinities” (Hinton and Roweis, 2003; Vladymyrov and Carreira-Perpiñán, 2013). Here, the bandwidth of each data point is computed so as to produce an effective number of neighbors k set by the user.) These hyperparameters values usually produce good clustering results and provide a starting point for improvement. In the homotopy algorithm, the path of σ or λ values should be followed slowly so we end in a good minimum. In practice, one changes the parameter geometrically in as many steps as one can afford computationally.

In a supervised setting, the hyperparameters can be selected with a validation set using the out-of-sample mapping for Laplacian K -modes (section 3.3).

3.3 Out-of-sample problem

We now consider the out-of-sample problem, that is, given an unseen test point $\mathbf{x} \in \mathbb{R}^D$, we wish to find a meaningful assignment $\mathbf{z}(\mathbf{x})$ to the clusters found during training. A natural and efficient way to do this is to solve a problem of the same form as (5) with a dataset consisting of the original training set augmented with \mathbf{x} , but keeping **Z** and **C** fixed to the values obtained during training (this avoids having to solve for all

Table 1: Comparison of properties of different clustering algorithms.

	K -means	K -medoids	Mean-shift	Spectral clustering	K -modes	Laplacian K -modes
Centroids	likely invalid	“valid”	“valid”	N/A	valid	valid
Nonconvex clusters	no	depends	yes	yes	no	yes
Density	no	no	yes	no	yes	yes
Assignment	hard	hard	hard	hard	hard	soft
Cost per iteration	KND	KN^2D	N^2D	$N^2 \sim N^3$	KND	$(D \log(1/\epsilon_1) + \rho/\sqrt{\epsilon_2})KN$

points again). After dropping constant terms, this is equivalent to the following problem:

$$\begin{aligned} \min_{\mathbf{z}} \quad & \lambda \sum_{n=1}^N w_n \|\mathbf{z} - \mathbf{z}_n\|^2 - \sum_{k=1}^K z_k G\left(\left\|\frac{\mathbf{x} - \mathbf{c}_k}{\sigma}\right\|^2\right) \\ \text{s.t.} \quad & z_k \geq 0, k = 1, \dots, K, \sum_{k=1}^K z_k = 1 \end{aligned}$$

where w_n is the affinity between test point \mathbf{x} and training point \mathbf{x}_n . The above problem can be further reduced to the following quadratic program:

$$\min_{\mathbf{z}} \quad \frac{1}{2} \|\mathbf{z} - (\bar{\mathbf{z}} + \gamma \mathbf{q})\|^2 \quad \text{s.t.} \quad \mathbf{z}^\top \mathbf{1}_K = 1, \mathbf{z} \geq \mathbf{0} \quad (7)$$

where the expressions for $\bar{\mathbf{z}}$, $\mathbf{q} = [q_1, \dots, q_K]^\top$ and γ are as follows:

$$\bar{\mathbf{z}} = \sum_{n=1}^N \frac{w_n}{\sum_{n'=1}^N w_{n'}} \mathbf{z}_n, \quad q_k = \frac{G(\|\mathbf{x} - \mathbf{c}_k\|/\sigma)^2}{\sum_{k'=1}^K G(\|\mathbf{x} - \mathbf{c}_{k'}\|/\sigma)^2}, \quad \gamma = \frac{\sum_{k=1}^K G(\|\mathbf{x} - \mathbf{c}_k\|/\sigma)^2}{2\lambda \sum_{n=1}^N w_n}.$$

Thus, the out-of-sample solution is the projection of the K -dimensional vector $\bar{\mathbf{z}} + \gamma \mathbf{q}$ onto the probability simplex. The computational cost is $\mathcal{O}(ND)$, dominated by the cost of $\bar{\mathbf{z}}$, since the simplex projection costs $\mathcal{O}(K \log K)$.

The solution has an intuitive interpretation, consisting of the linear combination of two terms, each a valid assignment vector (having positive elements that sum to 1). The Laplacian term, $\bar{\mathbf{z}}$, is the weighted average of the neighboring training points’ assignments, and results in nonconvex clusters. The kde term, \mathbf{q} , assigns a point based on its distances (posterior probabilities) to the centroids, and results in convex clusters. These two distinct assignment rules are combined using a weight γ to give the final assignment. Essentially, \mathbf{x} is assigned to cluster k with high probability if its nearby points are assigned to it (\bar{z}_k is large) or if it is close to \mathbf{c}_k (q_k is large). Although defined variationally, the out-of-sample mapping is just as useful as a closed-form expression: computationally it does not require an iterative procedure, and the interpretation above also illuminates the meaning of the Laplacian in the training objective (1). In fact, iterating the out-of-sample mapping sequentially over the training points gives another (slower) way to solve the \mathbf{Z} -step, i.e., alternating optimization over $\mathbf{z}_1, \dots, \mathbf{z}_N$.

Finally, Table 1 compares Laplacian K -modes with other popular clustering algorithms (see section 3.2 for the complexity analysis).

4 Experiments

4.1 Illustrative experiments

4.1.1 Spirals dataset

We first demonstrate the power of Laplacian smoothing. The 2D dataset in Fig. 2 consist of 5 spirals where each spiral contains 400 points (denoted by \circ). The natural way of partitioning this dataset into $K = 5$

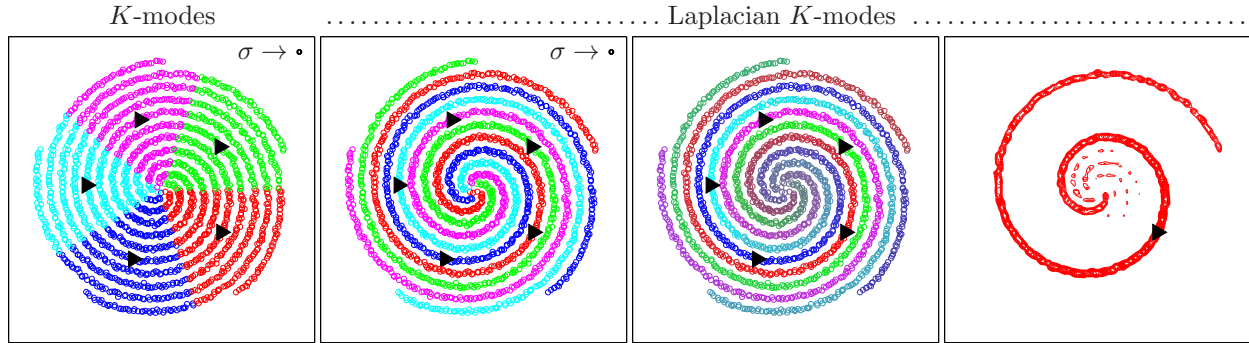


Figure 2: Synthetic dataset of 5-spirals. From left to right: K -modes clustering ($\lambda = 0$, $\sigma = 0.2$, the circle at the top right corner has a radius of σ); Laplacian K -modes clustering ($\lambda = 100$, $\sigma = 0.2$); Laplacian K -modes assignment probabilities; contours of the kde of the “red” cluster.

groups is to assign points of each spiral into a separate cluster. Due to the nonconvex shape of the spirals, the ideal result can not be possibly achieved by K -modes (plot 1, we color each point differently according to the cluster it is assigned to) or K -means (not shown, result similar to K -modes), even though the K -modes centroids (denoted by \blacktriangleright) are lying on each spiral and are valid representatives of the dataset. We then build a 5-nearest-neighbor graph on this dataset using heat kernel weighting and run Laplacian K -modes using the K -means result as initialization. We achieve a perfect separation of the spirals and one valid centroid for each spiral in a few steps of our alternating optimization scheme, as shown in plot 2. We show the assignment probabilities \mathbf{Z} in plot 3, where each data point \mathbf{x}_n is colored using a mixture of the 5 clusters’ colors with its assignment probability \mathbf{z}_n being the mixing coefficient. We show the contours of the kde defined on the “red” cluster in plot 4. In spite of the nonconvex, 1D manifold nature of the cluster, the kde is localized to the cluster and represents its shape and density well. This illustrates why using soft assignments gives more flexibility: the weights in the kde of each point vary, which allows for a more flexible kde. It is obvious that running mean-shift on this dataset with the same σ will result in a large number of modes and therefore clusters. In contrast, the number of modes is fixed in Laplacian K -modes and the algorithm will track one of the major modes in each cluster.

It is interesting to notice that because the kernel width σ we use is quite small, only a small proportion of data points are close enough to centroids to have nonzero affinity. This implies that the \mathbf{B} matrix in (5) of the main paper is quite sparse. Nonetheless, we achieve good assignment probabilities using the graph Laplacian, which propagates the sparse “label” information in \mathbf{B} throughout the graph. This also partly explains the success of Laplacian smoothing in spectral clustering (Shi and Malik, 2000) and semi-supervised learning algorithms (Zhu et al., 2003; Belkin et al., 2006).

4.1.2 Noisy two moons

We demonstrate Laplacian K -modes on the “two-moons” dataset in Fig. 3. The dataset has two nonconvex, interleaved clusters (each has 400 points) and we add many outliers (200 points) around them. The “moons” cannot be perfectly separated by either K -means (results shown in plot 1) or K -modes, since both define Voronoi tessellations. This problem is also difficult for hierarchical clustering because, as is well known, its major problem is that it creates connections between different clusters as the merging occurs. We build a 5-nearest-neighbor graph on this dataset using heat kernel weighting, and run Laplacian K -modes from the K -means initialization. We run the homotopy version and reduce σ from 5 to 0.1 in 10 steps while fixing $\lambda = 1$. The hard partition obtained, along with the two centroids and kde’s for each cluster at $\sigma = 0.1$ are given in plot 2. Even with the heavy noise and outliers, the “inliers” are still perfectly separated, the modes lie in high density areas and we obtain a good density estimate for each cluster. We show the assignment probabilities \mathbf{Z} in plot 3, colored using the same scheme as in the spirals example, where each data point \mathbf{x}_n is colored using a mixture of the clusters’ colors (red or blue) with its assignment probability \mathbf{z}_n being the mixing coefficient. The assignment is certain near the centroids (purer color) and less crisp at boundaries and outliers (mixed color). Finally, the out-of-sample mapping in input space is shown in plot 4, where we

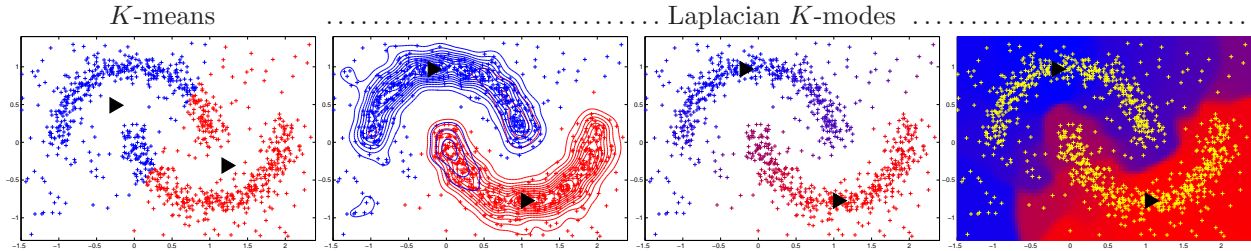


Figure 3: Synthetic dataset of 2-moons. We denote data points by $+$ and centroids by \blacktriangleright . We run Laplacian K -modes in homotopy and show results at final parameter value ($\lambda = 1$ and $\sigma = 0.1$). From left to right: K -modes clustering ($\lambda = 0$, $\sigma \rightarrow \infty$); Laplacian K -modes clustering and contours of kde of each cluster; Laplacian K -modes assignment probabilities for training set (used as mixing coefficients for coloring each training point); out-of-sample mapping in input space, colored in the same way as assignment probabilities of plot 3 (training points are now plotted in yellow).

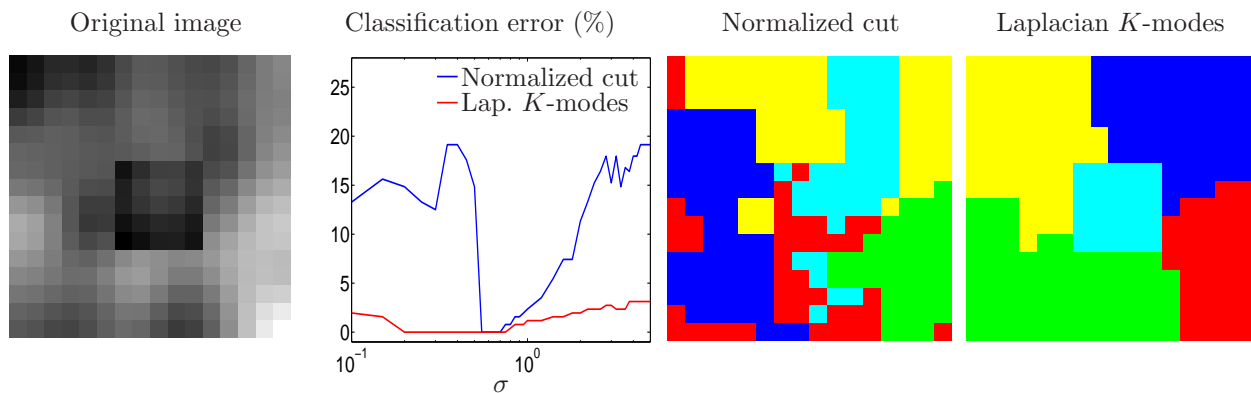


Figure 4: Occluder segmentation result: original image; classification errors over the range of σ ; segmentations of normalized cut and Laplacian K -modes at $\sigma = 0.2$.

compute out-of-sample assignments for a fine grid and color each grid point using the same scheme as in plot 3. We see clearly that the assignment rule is very different from the hard assignment of K -means. The mapping at each point combines the average assignment of nearby training points and the assignment from centroids, being able to model a complex shape.

4.2 Figure-ground segmentation

We consider the problem of segmenting an occluder from a textured background in a grayscale image. This problem has been shown to be difficult for spectral clustering (Chennubhotla and Jepson, 2003; Carreira-Perpiñán and Zemel, 2005), because of the intensity gradients between the occluder and the background (and within the background itself), which cause many graph edges to connect them, see the example in fig. 4. We formalize it as a clustering problem and partition the pixels into $K = 5$ clusters. We use for each pixel its 2D location and intensity value as features, and build a graph where each pixel is connected to the eight nearby pixels, with edges weighted using a heat kernel of width σ (σ 's value equals that of the kde bandwidth for Laplacian K -modes). The goal is to have one of the clusters extract the occluder, which can then be separated from the background. To measure the performance, we choose the cluster that overlaps most with the occluder as positive prediction (the rest of the pixels are considered as background/negative prediction) and compute the classification error. Fig. 4 shows normalized cut (Yu and Shi, 2003) performs well for a narrow range of σ , while Laplacian K -modes (with fixed $\lambda = 0.1$) has a much more stable performance when using the same graph: the range of good σ values that produce a perfect segmentation is much wider. This shows the difference between our algorithm and spectral clustering: even though both algorithms impose smoothness on assignments, the graph Laplacian is only a regularization term in our model, and the kde

term makes our algorithm more robust to the graph construction. As a more powerful algorithm, Laplacian K -modes has inherited the robustness properties from the original K -modes algorithm (Carreira-Perpiñán and Wang, 2013).

4.3 Clustering analysis

We report clustering statistics in datasets with known pattern class labels (which the algorithms did not use): (1) MNIST (LeCun et al., 1998), which contains 28×28 grayscale handwritten digit images (we randomly sample 200 of each digit); (2) COIL-20 (Nene et al., 1996), which contains 32×32 grayscale images of 20 objects viewed from varying angles; (3) the NIST Topic Detection and Tracking (TDT2) corpus, which contains on-topic documents of different semantic categories (documents appearing in more than one category are removed and only the largest 30 categories are kept). Statistics of the datasets are collected in table 2. Datasets (2) and (3) are the same as used by Cai et al. (2011), and we also use the same features: pixel values for (1) and (2), and TFIDF for (3).

We compare the following algorithms: K -means, initialized randomly; K -modes, a special case of Laplacian K -modes with $\lambda = 0$; Gaussian mean-shift (GMS), we search for σ that produces exactly K modes; Normalized cut (NCut), one typical spectral clustering algorithm, we use the implementation of Yu and Shi (2003); Graph regularized NMF (GNMF) proposed by Cai et al. (2011); Data-Cluster-Data random walk (DCD) proposed by Yang and Oja (2012); and Laplacian K -modes, initialized from K -means. K is set to the number of classes in the ground truth. All the datasets are normalized to have unit norm per sample.

Several algorithms use the graph Laplacian: for NCut, GNMF, and Laplacian K -modes, we build a 5-nearest-neighbor graph and use a binary weighting scheme to compute the graph Laplacian (as in Cai et al., 2011); for DCD, we find that it achieves better performance using a graph built with a larger neighborhood size, so we let DCD select its optimal size in $\{5, 10, 20, 30\}$. We run each algorithm with 20 random restarts, letting them use optimal values for their respective hyperparameters (if they have any) based on a grid search, and report the best performance from different random restarts. Notice we do not use the out-of-sample mapping here because its does not exist for all algorithms we compare with. Clustering accuracy (ACC) and normalized mutual information (NMI), two widely used criteria (Cai et al., 2011; Arora et al., 2011), are used for evaluation. The results are given in table 3 (N/A means our GMS code ran out of memory).

It is clear that algorithms using Laplacian smoothing are in general superior than algorithms not using it, which demonstrates the importance of the graph Laplacian in separating nonconvex and manifold clusters. GMS performs poorly for the reasons described earlier. On all datasets, Laplacian K -modes achieves the best or close to best performance under both criteria. We find there exists a wide range of hyperparameters with which our algorithm gives very competitive performance. We are able to further improve our performance on COIL-20 using the homotopy technique described earlier: we fix λ at 0.01, and decrease σ from 0.45 to 0.1 gradually in 7 steps, initializing the algorithm for the current σ value from the solution for the previous σ value. This improved result is shown in parenthesis in table 3.

Table 2: Statistics (size, dimensionality, # of classes) of three real world datasets.

dataset	N	D	K
MNIST	2000	784	10
COIL-20	1440	1024	20
TDT2	9394	36771	30

Table 3: Clustering accuracy and normalized mutual information (%) on 3 datasets.

	dataset	K -means	K -modes	GMS	NCut	GNMF	DCD	Laplacian K -modes
ACC	MNIST	58.2	59.2	15.9	65.5	66.2	69.4	70.5
	COIL-20	66.5	67.2	27.2	79.0	75.3	71.5	81.0 (81.5)
	TDT2	68.9	70.0	N/A	88.4	88.6	55.1	91.4
NMI	MNIST	53.3	53.6	6.51	66.9	64.9	65.6	68.8
	COIL-20	75.3	75.9	38.9	88.0	87.5	77.6	87.3 (88.0)
	TDT2	75.3	75.8	N/A	83.7	83.7	68.6	88.8

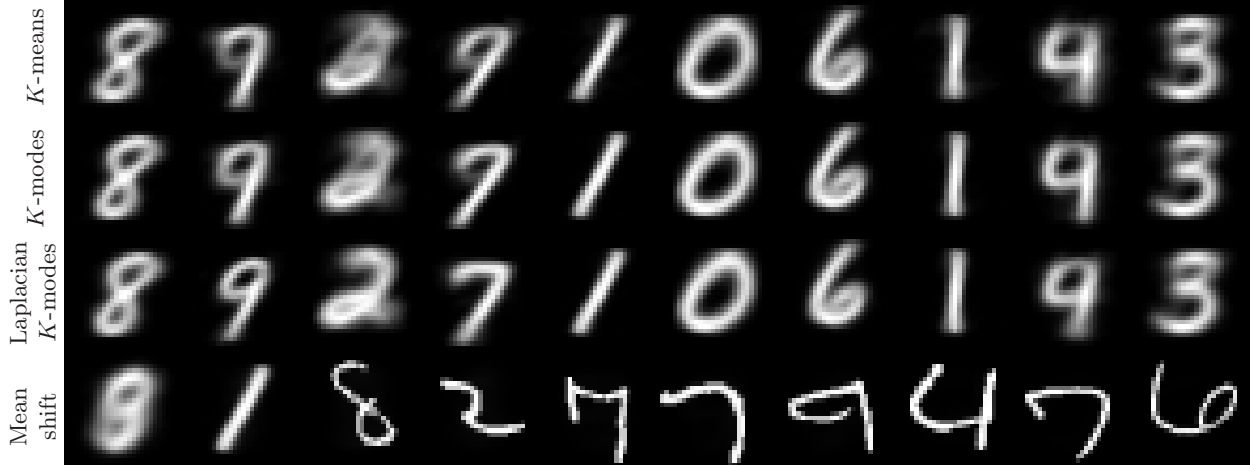


Figure 5: Centroids found by different algorithms on MNIST. *First row:* K -means ($\lambda = 0$, $\sigma \rightarrow \infty$, ACC: 55.2%, NMI: 50.2%). *Second row:* K -modes ($\lambda = 0$, $\sigma = 0.35$, ACC: 56.0%, NMI: 50.6%). *Third row:* Laplacian K -modes ($\lambda = 0.07$, $\sigma = 0.35$, ACC: 70.5%, NMI: 68.8%). *Fourth row:* mean-shift ($\sigma = 0.2485$).

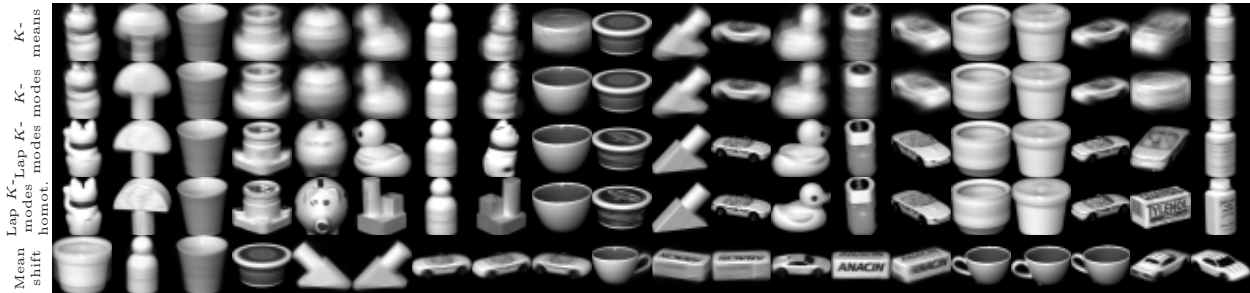


Figure 6: Centroids found by different algorithms on COIL-20. *First row:* K -means ($\lambda = 0$, $\sigma \rightarrow \infty$, ACC: 64.8%, NMI: 73.5%). *Second row:* K -modes ($\lambda = 0$, $\sigma = 0.3$, ACC: 65.5%, NMI: 73.0%). *Third row:* Laplacian K -modes ($\lambda = 0.01$, $\sigma = 0.1$, ACC: 73.2%, NMI: 83.8%). *Fourth row:* Laplacian K -modes with homotopy ($\lambda = 0.01$, ACC: 81.5%, NMI: 88.0%) in σ . *Fourth row:* mean-shift ($\sigma = 0.1591$).

Another key advantage of Laplacian K -modes is that the centroids are interpretable patterns of the dataset. We show the centroids (each as an image) found by centroid-based algorithms (using optimal hyperparameters) on MNIST in Fig. 5 and COIL-20 in Fig. 6, all using the K -means initialization. For such high-dimensional problems with small K , GMS tends to have a majority of centroids associated with very few points that are outliers with unusual patterns, and we do not show them here. Not surprisingly, some K -means centroids are blurry images consisting of an average of digits/objects of different identity and style. This implies some centroids lie between different branches of data manifolds, thus in low-density areas and not prototypical. The K -modes centroids have cleaner shapes, but the identities of the different centroids somewhat overlap, so that some classes are represented by multiple centroids, at the expense of other classes, which may be represented by no centroids. This is because K -modes only takes into account the kde's, and classes with higher sample density will receive more centroids. However, this will result in a class being abruptly partitioned, and the smoothness term in Laplacian K -modes prevents this, so its centroids not only have prototypical shapes, but also cover more digit/object identities. On COIL-20, starting from the K -means solution, and by slowly varying the hyperparameter σ , the centroids gradually move to high-density areas under the homotopy algorithm, being more prototypical, while the assignments separate different objects better.

Applying our algorithm at an intermediate σ achieves just the right amount of smoothing. This is clearly seen from the centroids obtained on MNIST. It allows the centroids to look like valid digit images, but at

the same time to average out noise, unusual strokes or other idiosyncrasies of the dataset images (while not averaging digits of different identities or different styles, as K -means does). This yields centroids that are more representative even than individual images of the dataset.

5 Conclusion

Our Laplacian K -modes algorithm enjoys some of the best properties of a range of clustering algorithms. It is nonparametric and allows the user to work with a kernel density estimate that produces exactly K clusters (as in K -means and K -modes), even in high dimension (unlike in mean-shift), and which can be nonconvex (as in mean-shift and spectral clustering). It also finds centroids that are valid patterns and lie in high-density areas, are representative of their cluster and neighborhood, yet they average out noise or idiosyncrasies that exist in individual data points. Computationally, our current alternating optimization scheme is simple, efficient and scales well. Experiments demonstrate the superior performance of Laplacian K -modes compared to well-known algorithms.

References

- R. Arora, M. Gupta, A. Kapila, and M. Fazel. Clustering by left-stochastic matrix factorization. In L. Getoor and T. Scheffer, editors, *Proc. of the 28th Int. Conf. Machine Learning (ICML 2011)*, pages 761–768, Bellevue, WA, June 28 – July 2 2011.
- D. Arthur and S. Vassilvitskii. **k-means++**: The advantages of careful seeding. In *Proc. of the 18th ACM-SIAM Symposium on Discrete Algorithms (SODA 2007)*, pages 1027–1035, New Orleans, LA, Jan. 7–9 2007.
- A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sciences*, 2(1):183–202, 2009.
- M. Belkin, P. Niyogi, and V. Sindhwani. Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *J. Machine Learning Research*, 7:2399–2434, Nov. 2006.
- C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer Series in Information Science and Statistics. Springer-Verlag, Berlin, 2006.
- D. Cai, X. He, J. Han, and T. S. Huang. Graph regularized nonnegative matrix factorization for data representation. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 33(8):1548–1560, Aug. 2011.
- M. Á. Carreira-Perpiñán. Mode-finding for mixtures of Gaussian distributions. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 22(11):1318–1323, Nov. 2000.
- M. Á. Carreira-Perpiñán. Acceleration strategies for Gaussian mean-shift image segmentation. In C. Schmid, S. Soatto, and C. Tomasi, editors, *Proc. of the 2006 IEEE Computer Society Conf. Computer Vision and Pattern Recognition (CVPR’06)*, pages 1160–1167, New York, NY, June 17–22 2006.
- M. Á. Carreira-Perpiñán. Gaussian mean shift is an EM algorithm. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 29(5):767–776, May 2007.
- M. Á. Carreira-Perpiñán and W. Wang. The K -modes algorithm for clustering. Unpublished manuscript, arXiv:1304.6478, Apr. 23 2013.
- M. Á. Carreira-Perpiñán and R. S. Zemel. Proximity graphs for clustering and manifold learning. In L. K. Saul, Y. Weiss, and L. Bottou, editors, *Advances in Neural Information Processing Systems (NIPS)*, volume 17, pages 225–232. MIT Press, Cambridge, MA, 2005.
- Y. Cheng. Mean shift, mode seeking, and clustering. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 17(8):790–799, Aug. 1995.

- C. Chennubhotla and A. Jepson. EigenCuts: Half-lives of EigenFlows for spectral clustering. In S. Becker, S. Thrun, and K. Obermayer, editors, *Advances in Neural Information Processing Systems (NIPS)*, volume 15, pages 705–712. MIT Press, Cambridge, MA, 2003.
- D. Comaniciu and P. Meer. Mean shift: A robust approach toward feature space analysis. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 24(5):603–619, May 2002.
- J. Duchi, S. Shalev-Shwartz, Y. Singer, and T. Chandra. Efficient projections onto the ℓ_1 -ball for learning in high dimensions. In A. McCallum and S. Roweis, editors, *Proc. of the 25th Int. Conf. Machine Learning (ICML'08)*, pages 272–279, Helsinki, Finland, July 5–9 2008.
- K. Fukunaga and L. D. Hostetler. The estimation of the gradient of a density function, with application in pattern recognition. *IEEE Trans. Information Theory*, IT-21(1):32–40, Jan. 1975.
- T. J. Hastie, R. J. Tibshirani, and J. H. Friedman. *The Elements of Statistical Learning—Data Mining, Inference and Prediction*. Springer Series in Statistics. Springer-Verlag, second edition, 2009.
- G. Hinton and S. T. Roweis. Stochastic neighbor embedding. In S. Becker, S. Thrun, and K. Obermayer, editors, *Advances in Neural Information Processing Systems (NIPS)*, volume 15, pages 857–864. MIT Press, Cambridge, MA, 2003.
- L. Kaufman and P. J. Rousseeuw. *Finding Groups in Data: An Introduction to Cluster Analysis*. Wiley Series in Probability and Mathematical Statistics. John Wiley & Sons, 1990.
- Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proc. IEEE*, 86(11):2278–2324, Nov. 1998.
- D. D. Lee and H. S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, Oct. 21 1999.
- W. Liu, J. He, and S.-F. Chang. Large graph construction for scalable semi-supervised learning. In J. Fürnkranz and T. Joachims, editors, *Proc. of the 27th Int. Conf. Machine Learning (ICML 2010)*, Haifa, Israel, June 21–25 2010.
- S. A. Nene, S. K. Nayar, and H. Murase. Columbia object image library (COIL-20). Technical Report CUCS-005-96, Dept. of Computer Science, Columbia University, Feb. 1996.
- J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 22(8):888–905, Aug. 2000.
- M. Vladymyrov and M. Á. Carreira-Perpiñán. Entropic affinities: Properties and efficient numerical computation. In S. Dasgupta and D. McAllester, editors, *Proc. of the 30th Int. Conf. Machine Learning (ICML 2013)*, pages 477–485, Atlanta, GA, June 16–21 2013.
- U. von Luxburg. A tutorial on spectral clustering. *Statistics and Computing*, 17(4):395–416, Dec. 2007.
- M. P. Wand and M. C. Jones. *Kernel Smoothing*. Number 60 in Monographs on Statistics and Applied Probability. Chapman & Hall, London, New York, 1994.
- W. Wang and M. Á. Carreira-Perpiñán. Projection onto the probability simplex: An efficient algorithm with a simple proof, and an application. Unpublished manuscript, arXiv:1309.1541, Sept. 3 2013.
- Z. Yang and E. Oja. Clustering by low-rank doubly stochastic matrix decomposition. In J. Langford and J. Pineau, editors, *Proc. of the 29th Int. Conf. Machine Learning (ICML 2012)*, pages 831–838, Edinburgh, Scotland, June 26 – July 1 2012.
- S. X. Yu and J. Shi. Multiclass spectral clustering. In *Proc. 9th Int. Conf. Computer Vision (ICCV'03)*, pages 313–319, Nice, France, Oct. 14–17 2003.
- X. Yuan, B.-G. Hu, and R. He. Agglomerative mean-shift clustering. *IEEE Trans. Knowledge and Data Engineering*, 24(2):209–219, Feb. 2010.

- L. Zelnik-Manor and P. Perona. Self-tuning spectral clustering. In L. K. Saul, Y. Weiss, and L. Bottou, editors, *Advances in Neural Information Processing Systems (NIPS)*, volume 17, pages 1601–1608. MIT Press, Cambridge, MA, 2005.
- X. Zhu, Z. Ghahramani, and J. Lafferty. Semi-supervised learning using Gaussian fields and harmonic functions. In T. Fawcett and N. Mishra, editors, *Proc. of the 20th Int. Conf. Machine Learning (ICML'03)*, pages 912–919, Washington, DC, Aug. 21–24 2003.