
Modal-set estimation with an application to clustering

Heinrich Jiang ¹

Samory Kpotufe
Princeton University

Abstract

We present a procedure that can estimate – with statistical consistency guarantees – any local-maxima of a density, under benign distributional conditions. The procedure estimates all such local maxima, or *modal-sets*, of any bounded shape or dimension, including usual point-modes. In practice, modal-sets can arise as dense low-dimensional structures in noisy data, and more generally serve to better model the rich variety of locally dense structures in data.

The procedure is then shown to be competitive on clustering applications, and moreover is quite stable to a wide range of settings of its tuning parameter.

1 INTRODUCTION

Mode estimation is a basic problem in data analysis. Modes, i.e. points of locally high density, serve as a measure of central tendency and are therefore important in unsupervised problems such as outlier detection, image or audio segmentation, and clustering in particular (as cluster cores). In the present work, we are interested in capturing a wider generality of *modes*, i.e. general structures (other than single-points) of locally high density, that can arise in modern data.

For example, application data in \mathbb{R}^d (e.g. speech, vision, medical imaging) often display locally high-density structures of arbitrary shape. Such an example is shown in Figure 1. While there are many quite reasonable ways of modeling such locally high-density structures, we show that the simple model investigated

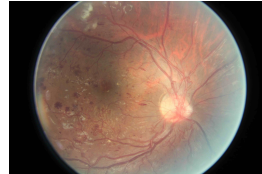


Figure 1: Medical imaging: high-resolution image of an eye with structures of blood capillaries. The aim in diabetic-retinopathy is to detect and delineate such capillary structures towards medical diagnostics.

in this work yields a practically successful procedure with strong statistical guarantees. In particular, our model simply allows the unknown density f to be locally maximal (or approximately maximal), not only at point-modes, but also on nontrivial subsets of \mathbb{R}^d . In other words, we make no a priori assumption on the nature of local maxima of the ground-truth f .

We will refer to connected subsets of \mathbb{R}^d where the unknown density f is locally maximal as *modal-sets* of f . A modal-set can be of any bounded shape and dimension, from 0-dimensional (point modes), to full dimensional surfaces, and aim to capture the rich variety of dense structures in data.

Our main contribution is a procedure, M(odal)-cores, that consistently estimates all such modal-sets from data, of general shape and dimension, with minimal assumption on the unknown f . If the ground-truth f is locally maximal at a point-mode, we return an estimate that converges to a point; if instead a modal-set is a surface, we consistently estimate that surface. Furthermore we have no a priori assumption on the number of modal-sets, besides that it is finite. Figure 2 shows a typical simulation on some arbitrary high-density structures.

The procedure builds on recent developments in topological data analysis [1, 2, 3, 4, 5, 6, 7], and works by carefully traversing a hierarchy of k -NN graphs which encode level sets of a k -NN density estimate. We show that, under mild uniform continuity assumptions on f , the Hausdorff distance between any modal-set and its estimate vanishes as $n \rightarrow \infty$ (Theorem 1, Section 2);

¹Much of this work was done when this author was at Princeton University.

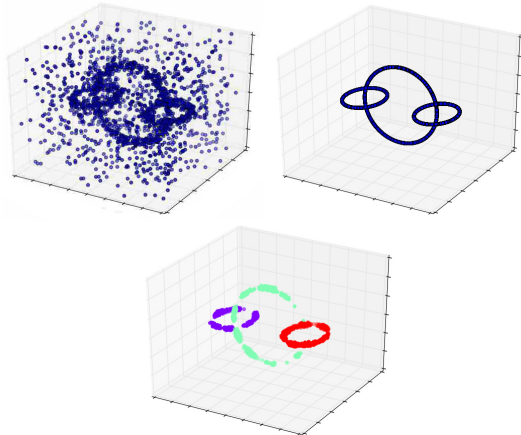


Figure 2: A first simulation on arbitrary shapes. (Top-Left) Points-cloud generated as three 1-dimensional rings + noise. (Top-Right) The 3 rings, and (Bottom) their estimate (as modal-sets) by M-cores.

the estimation rate for point-modes matches (up to $\log n$) the known minimax rates. Furthermore, under mild additional conditions on f (Hölder continuity), *false* structures (due to empirical variability) are correctly identified and pruned. We know of no other such general consistency guarantees in the context of mode estimation, especially for a practical procedure.

Finally, the present procedure is easy to implement and yields competitive scores on many clustering applications (Section 4); here, as in *mode-based clustering*, clusters are simply defined as regions of high-density of the data, and the estimated modal-sets serve as the centers of these regions, i.e., as *cluster-cores*. A welcome aspect of the resulting clustering procedure is its stability to tuning settings of the parameter k (from k -NN): it maintains high clustering scores (computed with knowledge of the ground-truth) over a wide range of settings of k , for various datasets. Such stability to tuning is of practical importance, since typically the ground-truth is unknown, making it difficult to tune the various hyperparameters of clustering procedures in practice. The performance of the present procedure seems rather robust to choices of its hyperparameters.

In the next section we put our result in context with respect to previous work on mode estimation and density-based clustering in general.

Related Work

- Much theoretical work on mode-estimation is concerned with understanding the statistical difficulty of the problem, and as such, often only considers the case of densities with single point-modes [8, 9, 10, 11, 12, 13]. The more practical case of densities with multiple

point-modes has received less attention in the theoretical literature. However there exist practical estimators, e.g., the popular *Mean-Shift* procedure (which doubles as a clustering procedure), which are however harder to analyze. Recently, [14] shows the consistency of a variant of Mean-Shift. Other recent work of [15] derives a method for pruning false-modes obtained by mode-seeking procedures. Also recent, the work of [16] shows that point-modes of a k -NN density estimate f_k approximate the true modes of the unknown density f , assuming f only has point-modes and bounded Hessian at the modes; their procedure, therefore operates on level-sets of f_k (similar to ours), but fails in the presence of more general high-density structures such as modal-sets. To handle such general structures, we have to identify more appropriate level-sets to operate on, the main technical difficulty being that local-maxima of f_k can be relatively far (in Hausdorff) from those of f , for instance single-point modes rather than more general modal-sets, due to data-variability. The present procedure handles general structures, and is consistent under the much weaker conditions of continuity (of f) on a compact domain.

A related line of work, which seeks more general structures than point-modes, is that of *ridge* estimation (see e.g. [17, 18]). A ridge is typically defined as a lower-dimensional structure away from which the density curves (in some but not all directions), and can serve to capture various lower-dimensional patterns apparent in point clouds. Also, related to ridge-estimation, is the area of *manifold-denoising* which seeks to understand noise-conditions under which a low-dimensional manifold can be recovered from noisy observations (see e.g. [19] for a nice overview), often with knowledge of the manifold dimension. In contrast to this line of work, the modal-sets defined here can be full-dimensional and are always local maxima of the density. Also, unlike in ridge estimation, we do not require local differentiability of the unknown f , nor knowledge of the dimension of the structure, thus targeting a different set of practical structures.

- A main application of the present work, and of mode-estimation in general, is *density-based clustering*. Such clustering was formalized in early work of [20, 21, 22], and can take various forms, each with their advantage.

In its hierarchical version, one is interested in estimating the connected components (CCs) of *all* level sets $\{f \geq \lambda\}_{\lambda > 0}$ of the unknown density f . Many recent works analyze approaches that consistently estimate such a hierarchy under quite general conditions, e.g. [1, 2, 3, 4, 5, 6, 7].

In the *flat* clustering version, one is interested in estimating the CCs of $\{f \geq \lambda\}$ for a single λ , somehow ap-

appropriately chosen [23, 24, 25, 26, 27, 28]. The popular DBSCAN procedure [29] can be viewed as estimating such single level set. The main disadvantage here is in the ambiguity in the choice of λ , especially when the levels λ of f have different numbers of clusters (CCs).

Another common flat clustering approach, most related to the present work, is *mode-based* clustering. The approach clusters points to estimated modes of f , a fixed target, and therefore does away with the ambiguity in choosing an appropriate level λ of f [30, 31, 32, 33, 34, 35, 36, 37]. As previously discussed, these approaches are however hard to analyze in that mode-estimation is itself not an easy problem. Popular examples are extensions of k -Means to categorical data [38], and the many variants of Mean-Shift which cluster points by gradient ascent to the closest mode. Notably, the recent work [39] analyzes clustering error of Mean-Shift in a general high-dimensional setting with potentially irrelevant features. The main assumption is that f only has point-modes.

2 OVERVIEW OF RESULTS

2.1 Basic Setup and Definitions

We have samples $X_{[n]} = \{X_1, \dots, X_n\}$ drawn i.i.d. from a distribution \mathcal{F} over \mathbb{R}^d with density f . We let \mathcal{X} denote the support of f . Our main aim is to estimate all local maxima of f , or *modal-sets* of f , as we will soon define.

We first require the following notions of distance between sets.

Definition 1. For $M \subset \mathcal{X}$, $x \in \mathcal{X}$, let $d(x, M) := \inf_{x' \in M} \|x - x'\|$. The **Hausdorff distance** between $A, B \subset \mathcal{X}$ is defined as $d(A, B) := \max\{\sup_{x \in A} d(x, B), \sup_{y \in B} d(y, A)\}$.

A modal set, defined below, extends the notion of a point-mode to general subsets of \mathcal{X} where f is locally maximal. These can arise for instance, as discussed earlier, in applications where high-dimensional data might be modeled as a (disconnected) manifold \mathcal{M} + ambient noise, each connected component of which induces a modal set of f in ambient space \mathbb{R}^D (see e.g. Figure 2).

Definition 2. For any $M \subset \mathcal{X}$ and $r > 0$, define the envelope $B(M, r) := \{x : d(x, M) \leq r\}$. A connected set M is a **modal-set** of f if $\forall x \in M$, $f(x) = f_M$ for some fixed f_M , and there exist $r > 0$ such that $f(x) < f_M$ for all $x \in B(M, r) \setminus M$.

Remark 1. The above definition can be relaxed to ϵ_0 -modal sets, i.e., to allow f to vary by a small ϵ_0 on M . Our results extend easily to this more relaxed definition, with minimal changes to some constants. This

is because the procedure operates on f_k , and therefore already needs to account for variations in f_k on M . This is described in the Appendix.

2.2 Estimating Modal-sets

The algorithm relies on nearest-neighbor density estimate f_k , defined as follows.

Definition 3. Let $r_k(x) := \min\{r : |B(x, r) \cap X_{[n]}| \geq k\}$, i.e., the distance from x to its k -th nearest neighbor. Define the **k -NN density estimate** as

$$f_k(x) := \frac{k}{n \cdot v_d \cdot r_k(x)^d},$$

where v_d is the volume of a unit ball in \mathbb{R}^d .

Furthermore, we need an estimate of the level-sets of f ; various recent work on cluster-tree estimation (see e.g. [6]) have shown that such level sets are encoded by subgraphs of certain *modified* k -NN graphs. Here however, we directly use k -NN graphs, simplifying implementation details, but requiring a bit of additional analysis.

Definition 4. Let $G(\lambda)$ denote the (*mutual*) **k -NN graph** with vertices $\{x \in X_{[n]} : f_k(x) \geq \lambda\}$ and an edge between x and x' iff $\|x - x'\| \leq \min\{r_k(x), r_k(x')\}$.

$G(\lambda)$ can be viewed as approximating the λ -level set of f_k , hence approximates the λ -level set of f (implicit in the connectedness result in the Appendix).

Algorithm 1 (M-cores) estimates the modal-sets of the unknown f . It is based on various insights described below. A basic idea, used for instance in point-mode estimation [16], is to proceed top-down on the level sets of f_k (i.e. on $G(\lambda)$, $\lambda \rightarrow 0$), and identify new modal-sets as they appear in separate CCs of a level λ .

Here however, we have to be more careful: the CCs of $G(\lambda)$ (modal-sets of f_k for some λ) might be singleton points (since f_k might take unique values over samples $x \in X_{[n]}$) while the modal-sets to be estimated might be of any dimension and shape. Fortunately, if a data point x , locally maximizes f_k , and belongs to some modal-set M of f , then the rest of $M \cap X_{[n]}$ must be at a nearby level; Algorithm 1 therefore proceeds by checking a nearby level ($\lambda - 9\beta_k\lambda$) from which it picks a specific set of points as an estimate of M . The main parameter here is β_k which is worked out explicitly in terms of k and requires no a priori knowledge of distributional parameters. The essential algorithmic parameter is therefore just k , which, as we will show, can be chosen over a wide range (w.r.t. n) while ensuring statistical consistency.

Definition 5. Let $0 < \delta < 1$. Define $C_{\delta,n} := 16 \log(2/\delta) \sqrt{d \log n}$, and define $\beta_k = 4 \frac{C_{\delta,n}}{\sqrt{k}}$.

If $\delta = 1/n$, then $C_{\delta,n} \approx \sqrt{d} \cdot (\log n)^{3/2}$. We note that the above definition of β_k is somewhat conservative (needed towards theoretical guarantees). Since $C_{\delta,n}$ behaves like a constant, it turns out to have little effect in implementation.

A further algorithmic difficulty is that a level $G(\lambda)$ might have too many CCs w.r.t. the ground truth. For example, due to variability in the data, f_k might have more modal-sets than f , inducing too many CCs at some level $G(\lambda)$. Fortunately, it can be shown that the nearby level $\lambda - 9\beta_k\lambda$ will likely have the right number of CCs. Such lookups down to lower-level act as a way of *pruning false modal-sets*, and trace back to earlier work [3] on pruning cluster-trees. Here, we need further care: we run the risk of over-estimating a given M if we look too far down (aggressive pruning), since a CC at lower level might contain points *far outside* of a modal-set M . Therefore, the main difficulty here is in figuring out how far down to look and yet not over-estimate *any* M (to ensure consistency). In particular our lookup *distance* of $9\beta_k\lambda$ is adapted to the level λ unlike in aggressive pruning.

Finally, for clustering with M-cores, we can simply assign every data-point to the closest estimated modal-set (acting as cluster-cores).

Algorithm 1 M-cores (estimating modal-sets).

```

Initialize  $\widehat{\mathcal{M}} := \emptyset$ . Define  $\beta_k = 4 \frac{C_{\delta,n}}{\sqrt{k}}$ .
Sort the  $X_i$ 's in decreasing order of  $f_k$  values (i.e.
 $f_k(X_i) \geq f_k(X_{i+1})$ ).
for  $i = 1$  to  $n$  do
    Define  $\lambda := f_k(X_i)$ .
    Let  $A \equiv$  CC of  $G(\lambda - 9\beta_k\lambda)$  containing  $X_i$ . (i)
    if  $A$  is disjoint from all cluster-cores in  $\widehat{\mathcal{M}}$  then
        Add  $\widehat{M} := \{x \in A : f_k(x) > \lambda - \beta_k\lambda\}$  to  $\widehat{\mathcal{M}}$ .
    end if
end for
return  $\widehat{\mathcal{M}}$ . // Each  $\widehat{M} \in \widehat{\mathcal{M}}$  is a cluster-core
    estimating a modal-set of the unknown  $f$ .
    
```

2.3 Consistency Results

Our consistency results rely on the following mild assumptions.

Assumption 1. f is continuous with compact support \mathcal{X} . Furthermore f has a finite number of modal-sets all in the interior of its support \mathcal{X} .

We will express the convergence of the procedure explicitly in terms of quantities that characterize the be-

havior of f at the boundary of every modal set. The first quantity has to do with how *salient* a modal-set, i.e, whether it is sufficiently *separated* from other modal sets. We start with the following definition of *separation*.

Definition 6. Two sets $A, A' \subset \mathcal{X}$ are r -separated, if there exists a set S such that every path from A to A' crosses S and $\sup_{x \in B(S,r)} f(x) < \inf_{x \in A \cup A'} f(x)$.

In simple terms, the definition says that there is a sufficiently wide valley in the density f between A and A' . For modal-set estimation, we also need such valleys to be sufficiently deep, which will be captured by the *curvature* of f at any modal-set.

The next quantities characterize the *change* in f in a neighborhood of a modal set M . The existence of a proper such neighborhood A_M , and appropriate functions u_M and l_M capturing smoothness and curvature, follow from the above assumptions on f . This is stated in the next proposition.

Proposition 1. Let M be a modal-set of f . Then there exists a CC A_M of some level-set $\mathcal{X}^{\lambda_M} := \{x : f(x) \geq \lambda_M\}$, containing M , such that the following holds.

- A_M isolates M by a valley: A_M does not intersect any other modal-set; and A_M and $\mathcal{X}^{\lambda_M} \setminus A_M$ are r_s -separated (by some S_M) for some $r_s > 0$ independent of M .
- A_M is full-dimensional: A_M contains an envelope $B(M, r_M)$ of M , for some $r_M > 0$.
- f is both *smooth* and has *curvature* around M : there exist functions u_M and l_M , increasing and continuous on $[0, r_M]$, $u_M(0) = l_M(0) = 0$, such that $\forall x \in B(M, r_M)$,

$$l_M(d(x, M)) \leq f_M - f(x) \leq u_M(d(x, M)).$$

Finally, our consistency guarantees require the following admissibility condition on $k = k(n)$. This condition results, roughly, from needing the density estimate f_k to properly approximate the behavior of f in the neighborhood of a modal-set M . In particular, we intuitively need f_k values to be smaller for points far from M than for points close to M , and this should depend on the smoothness and curvature of f around M (as captured by u_M and l_M).

Definition 7. k is **admissible** for a modal-set M if (we let u_M^{-1} denote the inverse of u_M):

$$\begin{aligned} & \max \left\{ \left(\frac{24 \sup_{x \in \mathcal{X}} f(x)}{l_M(\min\{r_M, r_s\}/2)} \right)^2, 2^{7+d} \right\} \cdot C_{\delta,n}^2 \\ & \leq k \leq \frac{v_d \cdot f_M}{2^{2+2d}} \left(u_M^{-1} \left(f_M \frac{C_{\delta,n}}{2\sqrt{k}} \right) \right)^d \cdot n. \end{aligned}$$

Remark 2. The admissibility condition on k , although seemingly opaque, allows for a wide range of settings of k . For example, suppose $u_M(t) = ct^\alpha$ for some $c, \alpha > 0$. These are polynomial tail conditions common in mode estimation, following e.g. from Hölder assumptions on f . Admissibility then (ignoring $\log(1/\delta)$), is immediately seen to correspond to the wide range

$$C_1 \cdot \log n \leq k \leq C_2 \cdot n^{2\alpha/(2\alpha+d)},$$

where C_1, C_2 are constants depending on M , but independent of k and n . It's clear then that even the simple choice $k = \Theta(\log^2 n)$ is always admissible for any M for n sufficiently large.

Main theorems. We then have the following two main consistency results for Algorithm 1. Theorem 1 states a rate (in terms of l_M and u_M) at which any modal-set M is approximated by some estimate in $\widehat{\mathcal{M}}$; Theorem 2 establishes *pruning* guarantees.

Theorem 1. Let $0 < \delta < 1$. The following holds with probability at least $1 - \delta$, simultaneously for all modal-sets M of f . Suppose k is admissible for M . Then there exists $\widehat{M} \in \widehat{\mathcal{M}}$ such that the following holds. Let l_M^{-1} denote the inverse of l_M .

$$d(M, \widehat{M}) \leq l_M^{-1} \left(\frac{8C_{\delta,n}}{\sqrt{k}} f_M \right),$$

which goes to 0 as $C_{\delta,n}/\sqrt{k} \rightarrow 0$.

If k is admissible for all modal-sets M of f , then $\widehat{\mathcal{M}}$ estimates all modal-sets of f at the above rates. These rates can be instantiated under the settings in Remark 2: suppose $l_M(t) = c_1 t^{\alpha_1}$, $u_M(t) = ct^\alpha$, $\beta_1 \geq \beta$; then the above bound becomes $d(M, \widehat{M}) \lesssim k^{-1/2\alpha_1}$ for admissible k . As in the remark, $k = \Theta(\log^2 n)$ is admissible, simultaneously for all M (for n sufficiently large), and therefore all modal-sets of f are recovered at the above rate. In particular, taking large $k = O(n^{2\alpha/(2\alpha+d)})$ optimizes the rate to $O(n^{-\alpha/(2\alpha_1\alpha+\alpha_1d)})$. Note that for $\alpha_1 = \alpha = 2$, the resulting rate ($n^{-1/(4+d)}$) is tight (see e.g. [12] for matching lower-bounds in the case of point-modes $M = \{x\}$).

Finally, Theorem 2 (pruning guarantees) states that any estimated modal-set in $\widehat{\mathcal{M}}$, at a sufficiently high level (w.r.t. to k), corresponds to a *true* modal-set of f at a similar level. Its proof consists of showing that if two sets of points are wrongly disconnected at level λ , they remain connected at nearby level $\lambda - 9\beta_k\lambda$ (so are reconnected by the procedure). The main technicality is the dependence of the nearby level on the empirical λ ; the proof is less involved and given in the Appendix.

Theorem 2. Let $0 < \delta < 1$. There exists $\lambda_0 = \lambda_0(n, k)$ such that the following holds with probability at least $1 - \delta$. All modal-set estimates in $\widehat{\mathcal{M}}$ chosen at level $\lambda \geq \lambda_0$ can be injectively mapped to modal-sets $\{M : \lambda_M \geq \min_{\{x \in \mathcal{X}_{[n]} : f_k(x) \geq \lambda - \beta_k\lambda\}} f(x)\}$, provided k is admissible for all such M .

In particular, if f is Hölder-continuous, (i.e. $\|f(x) - f(x')\| \leq c\|x - x'\|^\alpha$ for some $0 < \alpha \leq 1$, $c > 0$) then $\lambda_0 \xrightarrow{n \rightarrow \infty} 0$, provided $C_1 \log n \leq k \leq C_2 n^{2\alpha/(2\alpha+d)}$, for some C_1, C_2 independent n .

Remark 3. Thus with little additional smoothness ($\alpha \approx 0$) over uniform continuity of f , any estimate above level $\lambda_0 \rightarrow 0$ corresponds to a true modal-set of f . We note that these pruning guarantees can be strengthened as needed by implementing a more aggressive pruning: simply replace $G(\lambda - 9\beta_k\lambda)$ in the procedure (on line (i)) with $G(\lambda - 9\beta_k\lambda - \tilde{\epsilon})$ using a pruning parameter $\tilde{\epsilon} \geq 0$. This allows $\lambda_0 \rightarrow 0$ faster. However the rates of Theorem 1 (while maintained) then require a larger initial sample size n . This is discussed in the Appendix.

3 ANALYSIS OVERVIEW

The bulk of the analysis is in establishing Theorem 1. The key technicalities are in bounding distances from estimated cores to an unknown number of modal-sets of general shape, dimension and location.

The analysis considers each modal-set M of f separately, and only combines results in the end into the uniform consistency statement of Theorem 1. The following notion of *distance* from the sample $X_{[n]}$ to a modal-set M will be crucial.

Definition 8. For any $x \in \mathcal{X}$, let $r_n(x) := d(\{x\}, X_{[n]})$, and $r_n(M) := \sup_{x \in M} r_n(x)$.

We require a notion of a region \mathcal{X}_M that *isolates* a modal-set M from other modal-sets. In other words, \mathcal{X}_M containing only points *close* to M but far from other modes. To this end, let S_M denote the separating set from Proposition 1.

Definition 9. $\mathcal{X}_M := \{x : \exists \text{ a path } \mathcal{P} \text{ from } x \text{ to } x' \in M \text{ such that } \mathcal{P} \cap S_M = \emptyset\}$.

For each M , define $\hat{x}_M := \operatorname{argmax}_{x \in \mathcal{X}_M \cap X_{[n]}} f_k(x)$, a local maximizer of f_k on the modal-set M . The analysis (concerning each M) proceeds in the following steps:

- *Isolation of M* : when processing \hat{x}_M , the procedure picks an estimate \widehat{M} that contains no point from (or close to) modal-sets other than M .

- *Integrality of M* : the estimate \widehat{M} picks all of the envelope $B(M, r_n(M)) \cap X_{[n]}$.
- *Consistency of \widehat{M}* : it can then be shown that $\widehat{M} \rightarrow M$ in Hausdorff distance. This involves two directions: the first direction (that points of M are close to \widehat{M}) follows from integrality; the second direction is to show that points in \widehat{M} are close to M .

The following gives an upper-bound on the distance from a modal-set to the closest sample point. It follows from Bernstein-type VC concentration on masses of balls. The proof is given in the Appendix.

Lemma 1 (Upper bound on r_n). *Let M be a modal-set with density f_M and suppose that k is admissible. With probability at least $1 - \delta$,*

$$r_n(M) \leq \left(\frac{2C_{\delta,n}\sqrt{d\log n}}{n \cdot v_d \cdot f_M} \right)^{1/d}.$$

From Lemma 1 above, for any modal-set M of f , for n sufficiently large, there is a sample point in \mathcal{X}_M , i.e., $\mathcal{X}_M \cap X_{[n]} \neq \emptyset$. The next lemma conditions on the existence of such a sample.

The proofs for Lemma 2 and 3 are given in the Appendix but we give proof sketches here.

Lemma 2 (Isolation). *Let M be a modal-set and k be admissible for M . The following holds with probability at least $1 - \delta$. Suppose there exists a sample point in \mathcal{X}_M , and define $\hat{x}_M := \operatorname{argmax}_{x \in \mathcal{X}_M \cap X_{[n]}} f_k(x)$. When \hat{x}_M is processed in Algorithm 1, an estimate \widehat{M} is added to $\widehat{\mathcal{M}}$ satisfying $\widehat{M} \subset \mathcal{X}_M$.*

Proof sketch. First we choose $\bar{r} > 0$ depending on the smoothness and decay around M . It suffices to show that (i) $B(M, \bar{r})$ and $\mathcal{X} \setminus \mathcal{X}_M$ are disjoint in the k -NN graph when \hat{x}_M is being processed in Algorithm 1 and (ii) $\hat{x}_M \in B(M, \bar{r})$. To show (i), we use the k -NN bounds to justify that the k -NN graph does not contain any points from $B(S_M, \bar{r})$ or $\mathcal{X}_M \setminus B(M, \bar{r})$. Thus, any path from \hat{x}_M to $\mathcal{X}_M \setminus B(M, \bar{r})$ must contain an edge with length greater than \bar{r} . We then show there is no such edge. To show (ii), we argue that \hat{x}_M cannot be far away from M and is in fact within distance \bar{r} from M . \square

The above Lemma 2 establishes the existence of an estimate \widehat{M} of M containing no point from other modes. The next lemma establishes that such an estimate \widehat{M} contains all of $M \cap X_{[n]}$.

Lemma 3 (Integrality). *Let M be a modal-set with density f_M , and suppose k is admissible for M . The*

following holds with probability at least $1 - \delta$. First, \mathcal{X}_M does contain a sample point. Define $\hat{x}_M := \operatorname{argmax}_{x \in \mathcal{X}_M \cap X_{[n]}} f_k(x)$. When processing \hat{x}_M in Algorithm 1, suppose we add \widehat{M} to $\widehat{\mathcal{M}}$, then $B(M, r_n(M)) \cap X_{[n]} \subseteq \widehat{M}$.

Proof sketch. Let $A := B(M, r_n(M))$ and $\lambda_0 := (1 - \frac{C_{\delta,n}}{\sqrt{k}})^2 f_M$. We begin by showing that $A \cap X_{[n]}$ is connected in $G(\lambda_0)$ as follows. Let $r_\lambda := (k/(nv_d\lambda_0))^{1/d}$ which is the k -NN radius corresponding to a point with f_k value λ_0 , and $r_o := (k/(2nv_d f_M))^{1/d}$ which will be smaller than the k -NN radius of any sample point. Next, we show that $B(x, r_o)$ contains a sample point for any $x \in B(A, r_\lambda)$. Now, for any two points $x, x' \in A \cap X_{[n]}$, we use this fact to argue for the existence of a sequence of sample points starting with x and ending with x' such that the distance between adjacent points is less than r_o and all the points in the sequence lie in $B(A, r_o)$. We then show that each pair of adjacent points is an edge in $G(\lambda_0)$ and thus $A \cap X_{[n]}$ is connected in $G(\lambda_0)$. Finally, we argue that $\lambda \leq \lambda_0$ and thus $A \cap X_{[n]}$ is connected in $G(\lambda)$ as well. \square

Combining isolation and integrality, we obtain:

Corollary 1 (Identification). *Suppose the conditions of Lemmas 2 and 3 hold for modal-set M . Define $\hat{f}_M := \max_{x \in \mathcal{X}_M \cap X_{[n]}} f_k(x)$. With probability at least $1 - \delta$, there exists $\widehat{M} \in \widehat{\mathcal{M}}$ such that $B(M, r_n(M)) \cap X_{[n]} \subseteq \widehat{M} \subseteq \{x \in \mathcal{X}_M \cap X_{[n]} : f_k(x) \geq \hat{f}_M - \beta_k \hat{f}_M\}$.*

Proof. By Lemma 2, there exists $\widehat{M} \in \widehat{\mathcal{M}}$ which contains only points in \mathcal{X}_M with maximum f_k value of \hat{f}_M . Thus, we have $\widehat{M} \subseteq \{x \in \mathcal{X}_M \cap X_{[n]} : f_k(x) \geq \hat{f}_M - \beta_k \hat{f}_M\}$. By Lemma 3, $B(M, r_n(M)) \cap X_{[n]} \subseteq \widehat{M}$. \square

Here, we give a sketch of the proof for Theorem 1 which can be found in the Appendix.

Proof sketch of Theorem 1. Let $\tilde{r} = l_M^{-1} \left(\frac{8C_{\delta,n}}{\sqrt{k}} f_M \right)$. There are two directions to show: $\max_{x \in \widehat{M}} d(x, M) \leq \tilde{r}$ and $\sup_{x \in M} d(x, \widehat{M}) \leq \tilde{r}$.

For the first direction, by Corollary 1 we have $\widehat{M} \subseteq \{x \in \mathcal{X}_M : f_k(x) \geq \hat{f}_M - \beta_k \hat{f}_M\}$ where $\hat{f}_M := \max_{x \in \mathcal{X}_M \cap X_{[n]}} f_k(x)$. Thus, it suffices to show

$$\inf_{x \in B(M, r_n(M))} f_k(x) \geq \sup_{\mathcal{X}_M \setminus B(M, \tilde{r})} f_k(x) + \beta_k \hat{f}_M. \quad (1)$$

Using known upper and lower bounds on f_k in terms of f , we can lower bound the LHS by approximately $\hat{f}_M - u_M(r)$ (for some $r < \tilde{r}$) and upper bound the first term on the RHS by approximately $\hat{f}_M - l_M(\tilde{r})$.

The remaining difficulty is carefully choosing an appropriate r .

For the other direction, by Corollary 1, \widehat{M} contains all sample points in $B(M, r_n(M))$. Lemma 1 and the admissibility of k implies that $r_n(x) \leq \tilde{r}$ which easily gives us the result. \square

4 EXPERIMENTS

4.1 Practical Setup

The analysis prescribes a setting of $\beta_k = O(1/\sqrt{k})$. Throughout the experiments we simply fix $\beta_k = 2/\sqrt{k}$, and let our choice of k be the essential parameter. As we will see, M-cores yields competitive and stable performance for a wide-range of settings of k . The implementation can be done efficiently and is described in the Appendix.

A Python/C++ implementation of the code at [40].

4.2 Qualitative Experiments on General Structures

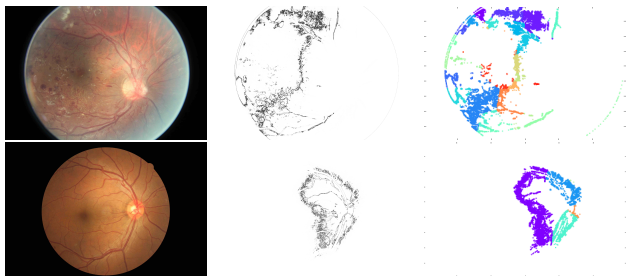


Figure 3: Diabetic Retinopathy: (Top 3 figures) An unhealthy eye, (Bottom 3 figures) A healthy eye. In both cases, shown are (1) original image, (2) a filter applied to the image, (3) modal-sets (structures of capillaries) estimated by M-cores on the corresponding filtered image. The unhealthy eye is characterized by a proliferation of damaged capillaries, while a healthy eye has visually fewer capillaries. The analysis task is to automatically discover the higher number of capillary-structures in the unhealthy eye. M-cores discovers 29 structures for unhealthy eye vs 6 for healthy.

We start with a qualitative experiment highlighting the flexibility of the procedure in fitting a large variety of high-density structures. For these experiments, we use $k = \frac{1}{2} \cdot \log^2 n$, which is within the theoretical range for admissible values of k (see Theorem 1 and discussion of Remark 2).

We consider a medical imaging problem. Figure 3 displays the procedure applied to the Diabetic Retinopathy detection problem [41]. While this is by no means

an end-to-end treatment of this detection problem, it gives a sense of M-cores’ versatility in fitting real-world patterns. In particular, M-cores automatically estimates a reasonable number of clusters, independent of shape, while pruning away (most importantly in the case of the healthy eye) false clusters due to noisy data. As a result, it correctly picks up a much larger number of clusters in the case of the unhealthy eye.

4.3 Clustering applications

We now evaluate the performance of M-cores on clustering applications, where for **clustering**: we assign every point $x_i \in X_{[n]}$ to $\operatorname{argmin}_{\widehat{M} \in \widehat{\mathcal{M}}} d(x_i, \widehat{M})$, i.e. to the closest estimated modal-set.

We compare M-cores to two common density-based clustering procedures, DBSCAN and Mean-Shift, as implemented in the *sci-kit-learn* package. Mean-Shift clusters data around point-modes, i.e. local-maxima of f , and is therefore most similar to M-cores in its objective.

Clustering scores. We compute two established scores which evaluate a clustering against a labeled ground-truth. The *rand-index*-score is the 0-1 accuracy in grouping pairs of points, (see e.g. [42]); the *mutual information*-score is the (information theoretic) mutual-information between the distributions induced by the clustering and the ground-truth (each cluster is a mass-point of the distribution, see e.g. [43]). For both scores we report the *adjusted* version, which adjusts the score so that a random clustering (with the same number of clusters as the ground-truth) scores near 0 (see e.g. [42], [43]).

Datasets. Phonemes [44], and UCI datasets: Glass, Seeds, Iris, and Wearable Computing. They are described in the table below.

Dataset	n	d	Labels	Description
Phonemes	4509	256	5	Log-periodograms of spoken phonemes
Glass	214	7	6	Properties of different types of glass
Seeds	210	7	3	Geometric measurements of wheat kernels
Iris	150	4	3	Various measurements over species of flowers
Wearable	10000	12	5	4 sensors on a human body, recording body posture and activity

Results. Figure 4 reports the performance of the procedures for each dataset. Rather than reporting the performance of the procedures under *optimal-tuning*, we report their performance *over a range* of hyperparameter settings, mindful of the fact that optimal-

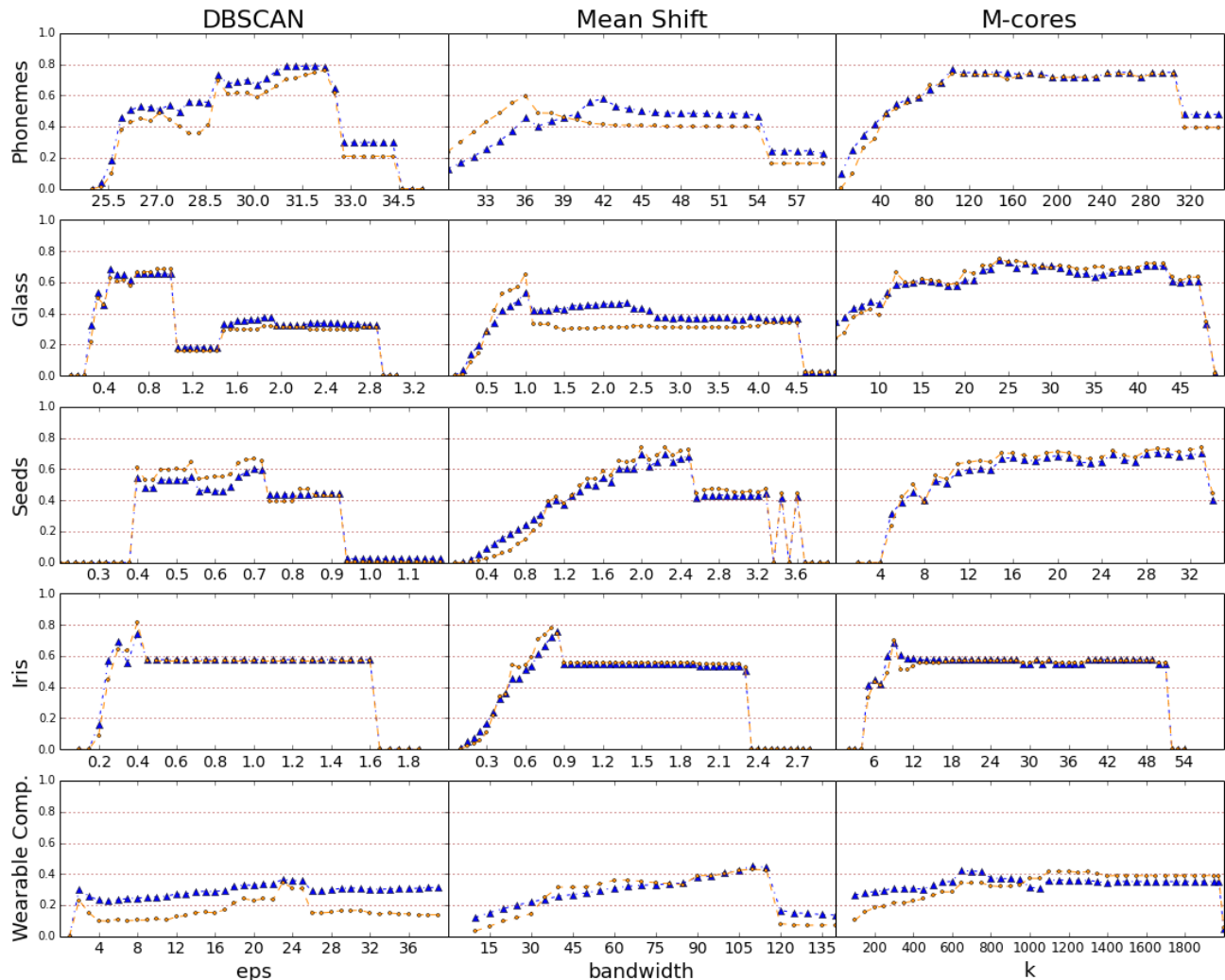


Figure 4: Comparison on real datasets (along the rows) across different hyperparameter settings for each algorithm (along the columns). The hyperparameters being tuned are displayed at the bottom of the figure for each clustering algorithm. Scores: the blue line with triangular markers is Adjusted-Mutual-Information, and the dotted red line is Adjusted-Rand-Index.

tuning is hardly found in practice (this is a general problem in clustering given the lack of ground-truth to guide tuning).

For M-cores we vary the parameter k . For DBSCAN and Mean-Shift, we vary the main parameters, respectively eps (choice of level-set), and $bandwidth$ (used in density estimation). M-cores yields competitive performance across the board, with stable scores over a large range of values of k (relative to sample size). Such stable performance to large changes in k is quite desirable, considering that proper tuning of hyperparameters remains a largely open problem in clustering.

Conclusion

We presented a theoretically-motivated procedure which can consistently estimate modal-sets, i.e. non-trivial high-density structures in data, under benign distributional conditions. This procedure is easily implemented and yields competitive and stable scores in clustering applications.

Acknowledgements

We are grateful to Sanjoy Dasgupta for useful discussions in the beginning of this project. Finally, part of this work was presented at the Dagstuhl 2016 seminar on Unsupervised Learning, and we are thankful for the useful feedback from participants.

References

- [1] W. Stuetzle and R. Nugent. A generalized single linkage method for estimating the cluster tree of a density. *Journal of Computational and Graphical Statistics*, 19(2):397–418, 2010.
- [2] K. Chaudhuri and S. Dasgupta. Rates for convergence for the cluster tree. In *Advances in Neural Information Processing Systems*, 2010.
- [3] S. Kpotufe and U. von Luxburg. Pruning nearest neighbor cluster trees. In *International Conference on Machine Learning*, 2011.
- [4] A. Rinaldo, A. Singh, R. Nugent, and L. Wasserman. Stability of density-based clustering. *Journal of Machine Learning Research*, 13:905–948, 2012.
- [5] S. Balakrishnan, S. Narayanan, A. Rinaldo, A. Singh, and L. Wasserman. Cluster trees on manifolds. In *Advances in Neural Information Processing Systems*, pages 2679–2687, 2013.
- [6] K. Chaudhuri, S. Dasgupta, S. Kpotufe, and U. von Luxburg. Consistent procedures for cluster tree estimation and pruning. *Arxiv*, 2014.
- [7] Justin Eldridge, Yusu Wang, and Mikhail Belkin. Beyond hartigan consistency: Merge distortion metric for hierarchical clustering. *Conference on Learning Theory*, 2015.
- [8] Emanuel Parzen et al. On estimation of a probability density function and mode. *Annals of mathematical statistics*, 33(3):1065–1076, 1962.
- [9] Herman Chernoff. Estimation of the mode. *Annals of the Institute of Statistical Mathematics*, 16(1):31–41, 1964.
- [10] William F Eddy et al. Optimum kernel estimators of the mode. *The Annals of Statistics*, 8(4):870–882, 1980.
- [11] Luc Devroye. Recursive estimation of the mode of a multivariate density. *Canadian Journal of Statistics*, 7(2):159–167, 1979.
- [12] Aleksandr Borisovich Tsybakov. Recursive estimation of the mode of a multivariate distribution. *Problemy Peredachi Informatsii*, 26(1):38–45, 1990.
- [13] Christophe Abraham, Gérard Biau, and Benoît Cadre. On the asymptotic properties of a simple estimate of the mode. *ESAIM: Probability and Statistics*, 8:1–11, 2004.
- [14] Ery Arias-Castro, David Mason, and Bruno Pelletier. On the estimation of the gradient lines of a density and the consistency of the mean-shift algorithm. *Unpublished Manuscript*, 2013.
- [15] Christopher Genovese, Marco Perone-Pacifico, Isabella Verdinelli, and Larry Wasserman. Non-parametric inference for density modes. *arXiv preprint arXiv:1312.7567*, 2013.
- [16] Sanjoy Dasgupta and Samory Kpotufe. Optimal rates for k-nn density and mode estimation. In *Advances in Neural Information Processing Systems*, pages 2555–2563, 2014.
- [17] Umut Ozertem and Deniz Erdogmus. Locally defined principal curves and surfaces. *The Journal of Machine Learning Research*, 12:1249–1286, 2011.
- [18] Christopher R Genovese, Marco Perone-Pacifico, Isabella Verdinelli, Larry Wasserman, et al. Non-parametric ridge estimation. *The Annals of Statistics*, 42(4):1511–1545, 2014.
- [19] Sivaraman Balakrishnan, Alessandro Rinaldo, Don Sheehy, Aarti Singh, and Larry A Wasserman. Minimax rates for homology inference. In *AISTATS*, volume 9, pages 206–207, 2012.
- [20] JW Carmichael, J Alan George, and RS Julius. Finding natural clusters. *Systematic Zoology*, pages 144–150, 1968.
- [21] John A Hartigan. *Clustering algorithms*. John Wiley & Sons, Inc., 1975.
- [22] J.A. Hartigan. Consistency of single linkage for high-density clusters. *Journal of the American Statistical Association*, 76(374):388–394, 1981.
- [23] P. Rigollet and R. Vert. Fast rates for plugin estimators of density level sets. *Bernoulli*, 15(4):1154–1178, 2009.
- [24] A. Singh, C. Scott, and R. Nowak. Adaptive hausdorff estimation of density level sets. *Annals of Statistics*, 37(5B):2760–2782, 2009.
- [25] M. Maier, M. Hein, and U. von Luxburg. Optimal construction of k-nearest neighbor graphs for identifying noisy clusters. *Theoretical Computer Science*, 410:1749–1764, 2009.
- [26] A. Rinaldo and L. Wasserman. Generalized density clustering. *Annals of Statistics*, 38(5):2678–2722, 2010.
- [27] I. Steinwart. Adaptive density level set clustering. In *24th Annual Conference on Learning Theory*, 2011.
- [28] Bharath K Sriperumbudur and Ingo Steinwart. Consistency and rates for clustering with dbscan. In *International Conference on Artificial Intelligence and Statistics*, pages 1090–1098, 2012.
- [29] Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm for discovering clusters in large spatial databases

- with noise. In *Kdd*, volume 96, pages 226–231, 1996.
- [30] Keinosuke Fukunaga and Larry Hostetler. The estimation of the gradient of a density function, with applications in pattern recognition. *Information Theory, IEEE Transactions on*, 21(1):32–40, 1975.
- [31] Yizong Cheng. Mean shift, mode seeking, and clustering. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 17(8):790–799, 1995.
- [32] Dorin Comaniciu and Peter Meer. Mean shift: A robust approach toward feature space analysis. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 24(5):603–619, 2002.
- [33] Jia Li, Surajit Ray, and Bruce G Lindsay. A non-parametric statistical approach to clustering via mode identification. *Journal of Machine Learning Research*, 8(8), 2007.
- [34] Frédéric Chazal, Leonidas J Guibas, Steve Y Oudot, and Primoz Skraba. Persistence-based clustering in riemannian manifolds. *Journal of the ACM (JACM)*, 60(6):41, 2013.
- [35] F. Chazal, B.T. Fasy, F. Lecci, B. Michel, A. Rinaldo, and L. Wasserman. Robust topological inference: Distance to a measure and kernel distance. *arXiv preprint arXiv:1412.7197.*, 2014.
- [36] Y.C. Chen, C.R. Genovese, and L. Wasserman. Statistical inference using the morse-smale complex. *arXiv preprint arXiv:1506.08826*, 2015.
- [37] Y.C. Chen, C.R. Genovese, and L. Wasserman. A comprehensive approach to mode clustering. *Electronic Journal of Statistics*, 10(1):210–241, 2016.
- [38] Anil Chaturvedi, Paul E Green, and J Douglas Carroll. K-modes clustering. *Journal of Classification*, 18(1):35–55, 2001.
- [39] Larry Wasserman, Martin Azizyan, and Aarti Singh. Feature selection for high-dimensional clustering. *arXiv preprint arXiv:1406.2240*, 2014.
- [40] M-cores code release. <http://github.com/hhjiang/mcores>.
- [41] Diabetic retinopathy detection. <https://www.kaggle.com/c/diabetic-retinopathy-detection>.
- [42] Lawrence Hubert and Phipps Arabie. Comparing partitions. *Journal of classification*, 2(1):193–218, 1985.
- [43] Nguyen Xuan Vinh, Julien Epps, and James Bailey. Information theoretic measures for clusterings comparison: Variants, properties, normalization and correction for chance. *The Journal of Machine Learning Research*, 11:2837–2854, 2010.
- [44] Trevor Hastie, Robert Tibshirani, Jerome Friedman, and James Franklin. The elements of statistical learning: data mining, inference and prediction. *The Mathematical Intelligencer*, 27(2):83–85, 2005.