# Recovering Trees with Convex Clustering*

Eric C. Chi† and Stefan Steinerberger‡

**Abstract.** Hierarchical clustering is a fundamental unsupervised learning task, whose aim is to organize a collection of points into a tree of nested clusters. Convex clustering has been proposed recently as a new way to construct tree organizations of data that are more robust to perturbations in the input data than standard hierarchical clustering algorithms. In this paper, we present conditions that guarantee when the convex clustering solution path recovers a tree and also make explicit how affinity parameters in the convex clustering formulation modulate the structure of the recovered tree. The proof of our main result relies on establishing a novel property of point clouds in a Hilbert space, which is of potentially independent interest.

**Key words.** Convex optimization, Fused lasso, Hierarchical clustering, Penalized regression, Sparsity

**AMS subject classifications.** 46C05, 49J99, 52C35

**1. Introduction.** Hierarchical clustering is a fundamental unsupervised learning task, whose aim is to organize a collection of points into a tree of nested clusters. To reinforce the idea that we seek a collection of nested clusters, we will often also refer to clusters as folders in this paper.

As an illustration, Figure 1 shows a collection of points in $\mathbb{R}^2$, labeled 1 to 18, that we seek to organize. Based on the Euclidean distances between the points, an intuitive organization is the following hierarchy of nested clusters. At the finest and first level of clustering, we partition the set $\{1, \ldots, 18\}$ into five subsets or folders:

$$F_{1,1} = \{1, 2, 3, 4, 5\}, \quad F_{1,2} = \{6, 7, 8\}, \quad F_{1,3} = \{9, 10, 11, 12, 13\},$$
$$F_{1,4} = \{14, 15, 16\}, \text{ and } F_{1,5} = \{17, 18\}.$$

At the second level of clustering, we merge the folders from the first level into a partition of two folders: $F_{2,1} = F_{1,1} \cup F_{1,2}$ and $F_{2,2} = F_{1,3} \cup F_{1,4} \cup F_{1,5}$.

Finally, at the third level of clustering, we merge the folders from the second level into a single folder: $F_{3,1} = F_{2,1} \cup F_{2,2}$. Figure 2 illustrates the described tree organization. Since each level of the tree consists of a partition of the data points, we refer to such hierarchical organizations as "partition trees."

There are many existing algorithms for automatically constructing partition trees, but perhaps the most often used algorithms in practice are collectively known as agglomerative hierarchical clustering methods [18, 21, 23, 30, 47]. Given a collection of points in $\mathbb{R}^p$, agglomerative hierarchical clustering methods recursively merge the points which are closest together

---

†Department of Statistics, North Carolina State University, Raleigh, NC (eric_chi@ncsu.edu).

‡Department of Mathematics, Yale University, New Haven, CT (stefan.steinerberger@yale.edu).

1

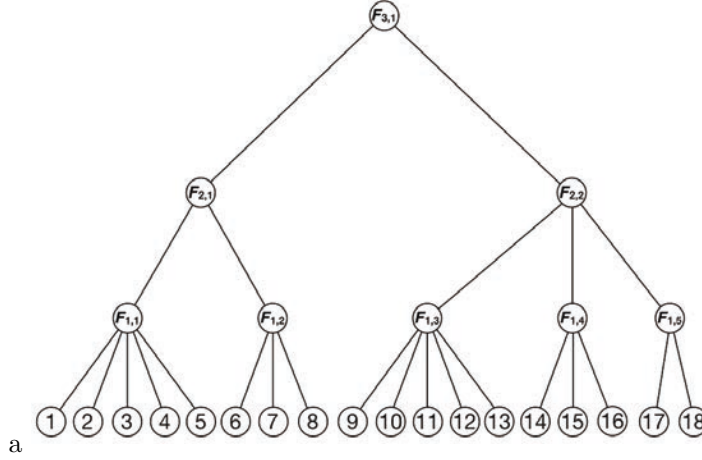Figure 1: Eighteen points in $\mathbb{R}^2$ to organize.



Figure 2: Partition Tree.

until all points are joined. Different choices in the definition of closeness lead to the different variants. Figure 3 shows two trees computed by two variants of the agglomerative hierarchical clustering. For each tree, the eighteen points reside in the "leaves" which are organized into a hierarchy of nested clusters that captures an increasingly coarser grouping structure as one progresses from the leaves to the root of the tree. The branch lengths in the tree quantify the similarity between pairs of points, or clusters at higher levels. We see that both trees recover binary partition trees that are similar to the ideal partition tree shown in Figure 2.
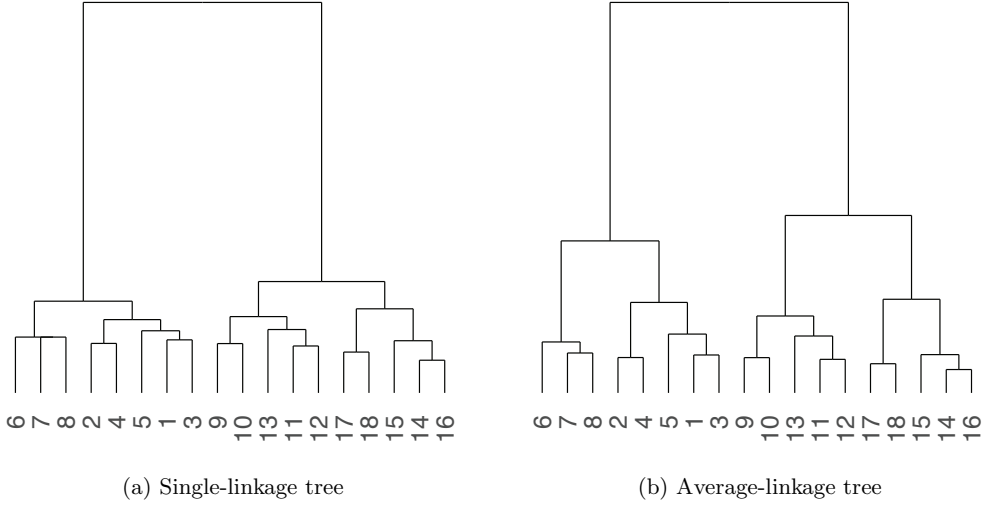
(a) Single-linkage tree                    (b) Average-linkage tree

Figure 3: Hierarchical clustering of data in Figure 1 under two different agglomeration methods.

**1.1. Convex Hierarchical Clustering?.** Although agglomerative hierarchical methods are widely used in practice, the greedy manner in which trees are constructed often results in an unstable mapping between input data and output tree. Indeed, agglomerative hierarchical clustering methods have been shown to be highly sensitive to perturbations in the input data, namely the resulting output trees can vary drastically with the addition of a little Gaussian noise to the data [10].

One promising alternative strategy for constructing trees stably relies on formulating the clustering problem as a continuous optimization problem. Following up on the initial proposal by [33], several recent works have shown that solving a sequence of convex optimization problems can recover tree organizations [9, 12, 19, 25, 32, 41]. Given $n$ points $x_1, \ldots, x_n$ in $\mathbb{R}^p$, we seek cluster centers (centroids) $u_i$ in $\mathbb{R}^p$ attached to point $x_i$ that minimize the convex criterion

$$(1.1) \qquad E_\gamma(u) = \frac{1}{2} \sum_{i=1}^{n} \|x_i - u_i\|^2 + \gamma \sum_{i<j} w_{ij} \|u_i - u_j\|,$$

where $\gamma$ is a nonnegative tuning parameter, $w_{ij}$ is a nonnegative affinity that quantifies the similarity between $x_i$ and $x_j$, and $u$ is the vector in $\mathbb{R}^{np}$ obtained by stacking the vectors $u_1, \ldots, u_n$ on top of each other. For now, we assume all norms are Euclidean norms; we will later consider arbitrary norms. The sum of squares data-fidelity term in (1.1) quantifies how well the centroids $u_i$ approximate the data $x_i$, while the sum of norms regularization term penalizes the differences between pairs of centroids $u_i$ and $u_j$. To expand on the latter, the regularization term is a composition of the group lasso [51] and the fused lasso [44] and incen-

(a) Gaussian Kernel Affinities                                    (b) Unit Affinities
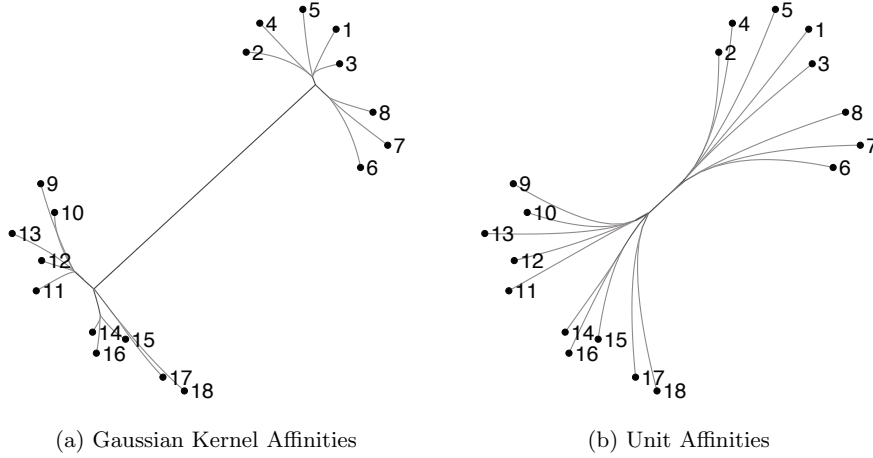
Figure 4: Solution paths of convex clustering using different affinities $w_{ij}$.

61   tivizes sparsity in the pairwise differences of centroid pairs. Overall, $E_\gamma(u)$ can be interpreted
62   as the energy of a configuration of centroids $u$ for a given relative weighting $\gamma$ between data-
63   fidelity and model complexity as quantified by the regularization term. We next elaborate
64   how $u(\gamma)$ varies as the tuning parameter $\gamma$ varies.

65      Because the objective function $E_\gamma(u)$ in (1.1) is strongly convex, for each value of $\gamma$ it
66   possesses a unique minimizer $u(\gamma)$, whose $n$ subvectors in $\mathbb{R}^p$ we denote by $u_i(\gamma)$. The tuning
67   parameter $\gamma$ trades off the relative emphasis between data fit and differences between pairs
68   of centroids. When $\gamma = 0$, the minimum is attained when $u_i = x_i$, namely when each point
69   occupies a unique cluster. As $\gamma$ increases, the regularization term encourages cluster centroids
70   to fuse together. Two points $x_i$ and $x_j$ with $u_i = u_j$ are said to belong to the same cluster.
71   For sufficiently large $\gamma$, the $u_i$ fuse into a single cluster, namely $u_i = \overline{x}$, where $\overline{x}$ is the average
72   of the data $x_i$ [12, 42]. Moreover, the unique global minimizer $u(\gamma)$ is a continuous function
73   of the tuning parameter $\gamma$ [10]; we refer to the continuous paths $u_i(\gamma)$, traced out from each
74   $x_i$ to $\overline{x}$ as $\gamma$ varies, collectively as the solution path. Thus, by computing $u_i(\gamma)$ for a sequence
75   of $\gamma$ over an appropriately sampled range of values, we hope to recover a partition tree.

76      Figure 4 plots the $u_i$ as a function of $\gamma$ for two different sets of affinities $w_{ij}$. We will discuss
77   the differences in the recovered trees shortly, but for now we point out that computing $u(\gamma)$ for
78   a range of $\gamma$ indeed appears to recover trees that bear similarity to the desired partition tree
79   in Figure 2. Moreover, the $u_i(\gamma)$ are 1-Lipschitz functions of the data $x_i$ [11]. Consequently,
80   small perturbations to the input data $x_i$, are guaranteed to *not* result in disproportionately
81   large variations in the output $u_i(\gamma)$.

82      At this point, the solution path of convex clustering appears to stably recover partition
83   trees as desired. Nonetheless, questions remain as to whether convex clustering is a form
84   of convex hierarchical clustering. Specifically, (i) when is the solution path guaranteed to
85   produce a tree, and (ii) how do the affinities modulate the branch formation in the recovered

86 tree?

87     Hocking et al. provide a partial answer to the first question [19]. They prove that if unit
88 affinities are used, namely $w_{ij} = 1$ for all $i$ and $j$, and if 1-norms are used in the regularization
89 term in (1.1), then the solution path must be a tree. On the other hand, in the same paper,
90 they also provide an example, using the Euclidean norm in the regularization term, where
91 the solution path can fail to be a tree. Specifically, as the tuning parameter $\gamma$ increases, it is
92 possible for centroids to initially fuse and then "unfuse" before eventually fusing again. We
93 provide an example of this phenomenon in Appendix A.

94     The differences in the two recovered trees shown in Figure 4 motivate the second question.
95 Figure 4a shows the solution path when using Gaussian kernel affinities, namely for all $i$ and
96 $j$

$$ w_{ij} = \exp\left(-\frac{\|x_i - x_j\|_2^2}{\sigma}\right), $$

98 where $\sigma$ is a positive scale parameter. Gaussian kernel affinities have been empirically shown
99 to provide more aggressive fusion of folders closer to the leaves, and consequently more infor-
100 mative, hierarchical clustering results [10, 12, 19]. Figure 4b shows the solution path when
101 using unit affinities. We see that Gaussian kernel affinities can generate a solution path that
102 recovers the partition tree in Figure 2, while unit affinities can generate a solution path that
103 recovers a less "nested" approximation to the partition tree in Figure 2. The same sets of
104 points and folders are getting shrunk together in Figure 4a and Figure 4b, but less aggres-
105 sively in the latter as $\gamma$ increases. In Appendix B, we provide an additional real data example
106 highlighting how different the recovered trees can be under the two sets of affinities. Our
107 main result will complement these empirical observations with a theoretical argument for why
108 certain data-driven affinities, including but not limited to Gaussian kernel affinities, should
109 be preferred over others.

110     **1.2. Contributions.** In this paper, we answer the open questions of (i) why the solution
111 path of convex clustering can recover a tree and (ii) how affinities can be chosen to guarantee
112 recovery of a given partition tree on the data. We first answer these questions in the case
113 when Euclidean norms are employed in (1.1) and then later describe how our results can be
114 extended to more general data-fidelity terms and arbitrary norms in the regularization term.

115     We clarify how the theoretical contributions in this paper differ from existing theoretical
116 results in the convex clustering literature. Radchenko and Mukherjee in [34] present a pop-
117 ulation model for the convex clustering procedure and provide an analysis of the asymptotic
118 properties of the sample convex clustering procedure. We note that their analysis is specific
119 to using 1-norms in the regularization term, while we consider first the Euclidean norm before
120 generalizing to arbitrary ones. Zhu et al. in [54] provide conditions under which two true un-
121 derlying clusters can be identified by solving the convex clustering problem with appropriately
122 chosen affinities. Similarly, She [39] and Sharpnack et al. [38] present results when the convex
123 clustering solution can consistently recover groupings. Others present finite sample prediction
124 error bounds for recovery of a latent set of clusters [42, 46].

125     Our contributions differ from these prior works in two ways. First, we provide conditions
126 on the affinities that ensure that the solution path reconstructs an *entire* hierarchical partition

127  tree and clarify how these affinities can be explicitly tuned to recover a specific target tree.
128  With the exception of the work by Radchenko and Mukherjee in [34], all of the other works
129  present theoretical guarantees for recovering a *single* partition level rather than a nested
130  hierarchy of partitions. Second, in contrast to all of the previous work, we do not make any
131  distributional assumptions on the data. Instead, we focus in this paper on understanding the
132  behavior of the solution path as a function of the affinities used in the regularization term.
133  By understanding this dependency, we gain insight into why a commonly used data-driven
134  affinities choice, namely the Gaussian kernel, works so well in practice.

135      **1.3. Outline.** The rest of this paper proceeds as follows. In Section 2, we define structures
136  needed to construct affinities that will enable us to recover a desired partition tree and once
137  equipped with the necessary building blocks, give an overview of our main result. In Section 3,
138  we introduce a geometric lemma that is key to proving our main result. In Section 4, we
139  give proofs of the geometric lemma and our main theorem. In Section 5, we show how our
140  main result can be generalized to other data-fidelity terms and regularization term norms. In
141  Section 7, we conclude with a discussion on our results within the broader context of penalized
142  regression methods for clustering.

143      **2. Setup and Overview of Main Result.** Our main result shows that if the affinities
144  $w_{ij}$ arise from an underlying partition tree, then that tree can be reconstructed from the
145  solution path of the convex clustering problem. To proceed, we will need a formal definition
146  of a partition tree and then a judicious assignment of weights to the edges in the tree graph
147  corresponding to the partition tree.

148      **2.1. Partition Tree.** Let $\Omega = \{x_1, \ldots, x_n\} \subset \mathbb{R}^p$ be an arbitrary collection of points and
149  let $[n]$ denote the set of indices $\{1, \ldots, n\}$. Following the notation and language employed in
150  [2] and [29, 28], we say that $\mathcal{T}$ is a partition tree on the collection of points $\Omega$ consisting of
151  $\mathcal{P}_0, \ldots, \mathcal{P}_L$ partitions of $\Omega$ if it has the following properties:
152      1. The partition $\mathcal{P}_l = \{F_{l,1}, \ldots, F_{l,n_l}\}$ at level $l$ consists of $n_l$ disjoint non-empty subsets
153          of indices in $\{1, \ldots, n\}$, termed folders and denoted by $F_{l,i}, i \in [n_l]$.
154      2. The finest partition $\mathcal{P}_0$ contains $n_0 = n$ singleton "leaf" folders, namely $F_{0,i} = \{i\}$.
155      3. The coarsest partition $\mathcal{P}_L$ contains a single "root" folder, namely $F_{L,1} = [n]$.
156      4. Partitions are nested; if $F \in \mathcal{P}_l$, then $F \subseteq F'$ for some $F' \in \mathcal{P}_{l+1}$, namely each folder
157          at level $l - 1$ is a subset of a folder from level $l$. Note that we allow for $F = F'$.
158  A partition tree $\mathcal{T}$ on $\Omega$ can be seen as the collection of all folders at all levels, namely
159  $\mathcal{T} = \{F_{l,i} : 0 \leq l \leq L, i \in [n_l]\}$.

160      **2.2. Weighted Tree Graph.** We next assign every folder $F_{l,i} \in \mathcal{T}$ to a node and draw an
161  edge between nested folders in adjacent levels. Thus, if $F \in \mathcal{P}_l, F' \in \mathcal{P}_{l+1}$, and $F \subset F'$, then
162  we draw an edge $(F, F')$ between $F$ and $F'$. If we let $\mathcal{E}$ denote the set of all edges between
163  nested folders in adjacent levels, then the resulting graph $\mathcal{G} = (\mathcal{E}, \mathcal{T})$ is a tree.
164      We next assign weights on the edges in $\mathcal{E}$ as follows. Let $\varepsilon > 0$ be a fixed parameter,
165  whose value we will elaborate on shortly. Edges between level 0 folders and level 1 folders
166  receive a weight of 1. Edges between level 1 folders and level 2 folders receive a weight of
167  $\varepsilon$. Edges between level 2 folders and level 3 folders receive a weight of $\varepsilon^2$ and so on. Thus,
168  edges between level $l$ folders and level $l + 1$ folders receive a weight of $\varepsilon^l$. Figure 5a shows the

169    weighted tree graph $\mathcal{G}$ derived from the partition tree given in Figure 2.



(a) Weighted Tree Graph

(b) The path $p_{15}$ from 1 to 5 produces $w_{15} = 1$.

(c) The path $p_{17}$ from 1 to 7 produces $w_{17} = \varepsilon$.

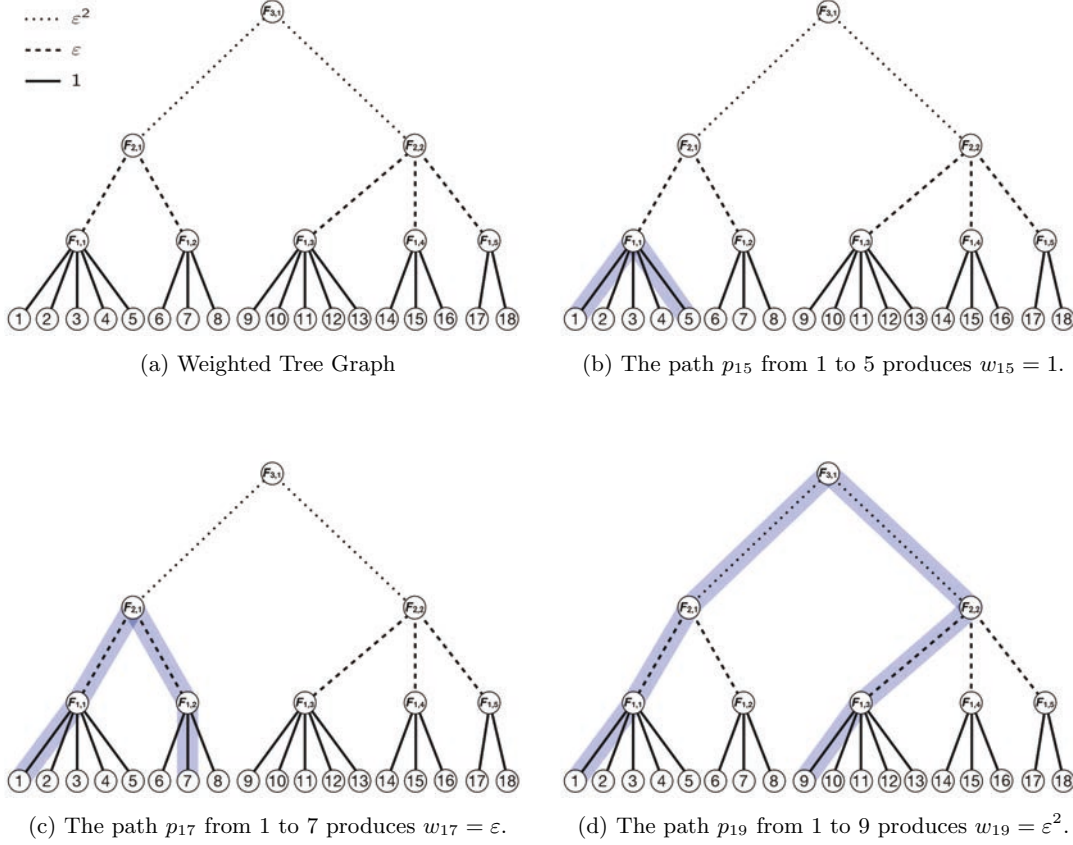(d) The path $p_{19}$ from 1 to 9 produces $w_{19} = \varepsilon^2$.

Figure 5: Weighted Tree: Edges that are solid lines have weight 1. Edges that are dashed lines have weight $\varepsilon$. Edges that are dotted lines have weight $\varepsilon^2$.

170      We are finally ready to construct $w_{ij}$ from the weighted tree graph. Let $F_{0,i}$ and $F_{0,j}$ be
171    leaf nodes in the graph $\mathcal{G}$ and let $p_{ij}$ be the sequence of edges in $\mathcal{E}$ that form the path between
172    $F_{0,i}$ and $F_{0,j}$. Then we set $w_{ij}$ to be the smallest weight of edges contained in $p_{ij}$. In other
173    words, $w_{ij}$ is the smallest edge weight one sees in traveling from $i$ to $j$. Figure 5b shows that
174    the path $p_{15}$ from 1 to 5 in the weighted graph $\mathcal{G}$ leads to the affinity assignment $w_{15} = 1$.
175    Figure 5c and Figure 5d show additional examples of how affinities are derived from the edge
176    weights in $\mathcal{G}$.

177      **2.3. Main Result.** We now state our main result.

178      Theorem 2.1. *There exists $\varepsilon_0 > 0$, depending on the data and the tree structure (which we*
179    *assume defines the $w_{ij}$ as outlined above in Section 2.2), so that for all $\varepsilon \in (0, \varepsilon_0)$ the solution*

180   *path*

$$u(\gamma) = \underset{u_1,\ldots,u_n}{\arg\min} \sum_{i=1}^{n} \|x_i - u_i\|^2 + \gamma \sum_{i,j=1}^{n} w_{ij}\|u_i - u_j\|,$$

182   *as parametrized by $\gamma \in (0, \gamma_0)$ traces out exactly the partition tree structure underlying the*
183   *affinities $w_{ij}$ before collapsing into a point for some large, but finite, $\gamma_0$.*

184      Informally speaking, this means that as $\gamma$ increases, elements from the same folder collapse
185   into a single point, these folders (now single points) move themselves (or rather, the fused
186   points move in a coordinated manner) and then collapse again in a way predicted by the tree
187   (i.e. folders sharing a parent folder collapse). This evolution continues on until all points have
188   collapsed into a single point (which happens for a finite value $\gamma_0$). We have no precise bound
189   on the times $\gamma$ at which these collapses happen but by making $\varepsilon_0$ sufficiently small, there is
190   an arbitrary long time between stages of collapsing. The proof of Theorem 2.1 also gives a
191   bound on $\gamma_0$ as a byproduct.
192
193      **Remarks** Several additional remarks are in order.
194      1. At first blush, it appears that the data $x_i$ plays no role in the recovered partition tree
195         as the affinities $w_{ij}$ dictate the trajectory of the solution path. In practice, however,
196         one would *never* use $w_{ij}$ that did not depend on the data. We study the convex
197         clustering solution path separate of any particular data-driven choice of the affinities,
198         but intuitively the affinity $w_{ij}$ should be inversely proportional to the distance between
199         $x_i$ and $x_j$. Theorem 2.1 further clarifies a sufficient condition on how *rapidly* (i.e.
200         geometrically fast) the affinity $w_{ij}$ should decrease as the distance between $x_i$ and $x_j$
201         increases for all pairs of data points, to ensure the solution path is a tree. To further
202         clarify the importance of using $w_{ij}$ that respect the geometry of the data, we give an
203         example of a solution path that is *not* a tree as a consequence of using $w_{ij}$ that do not
204         respect the geometry of the data in Appendix A.
205      2. The affinities do not need to have exactly the structure described in Section 2.2.
206         A more precise statement would be that there exists an $\varepsilon_0$ such that whenever we
207         associate weight $\varepsilon_1 \in (0, \varepsilon_0)$ to the first level, then there exists an $\varepsilon$ (depending on
208         everything and $\varepsilon_0, \varepsilon_1$) such that if we associate weight $\varepsilon_2 \in (0, \varepsilon)$ to the second level
209         there exists an $\varepsilon_3$ (depending on everything and $\varepsilon_0, \varepsilon_1, \varepsilon_2$ etc.). Simply put, it suffices
210         to have a sufficiently clear separation of scales encoded in the affinities.
211         Indeed, Figure 6 shows the Gaussian kernel affinities $w_{1j}$ between $x_1$ and the remain-
212         ing $x_j$ for $j = 2, \ldots, 18$ from the example in Figure 1. We observe clear separation
213         of scales encoded in the Gaussian kernel affinities that align with the partition tree
214         and corresponding weighted graph $\mathcal{G}$ in Figure 5a. Similar plots of the set of affinities
215         associated with each data point reveal alignment with the partition tree and corre-
216         sponding weighted graph $\mathcal{G}$. The key quality of the Gaussian kernel should be readily
217         apparent, namely the Gaussian kernel naturally encodes, in a data-driven way, a ge-
218         ometric decay in weights that is sufficient to reconstruct a partition tree embedded
219         in Euclidean space. We emphasize, however, that there is nothing special about the
220         Gaussian kernel, and its rapid decay in weights is not even necessary. Any data-driven
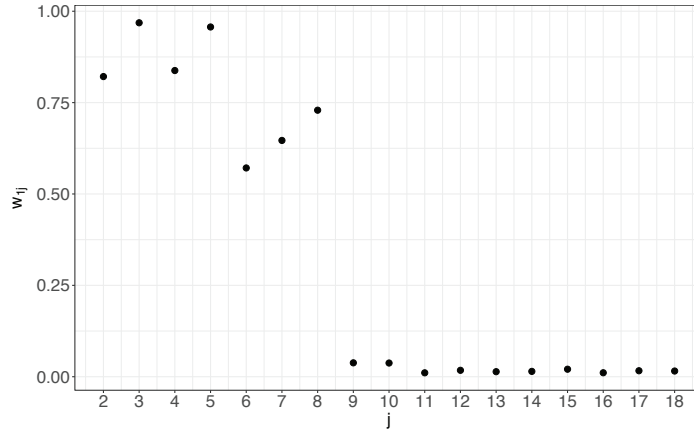
Figure 6: Gaussian kernel affinities $w_{1j}$ between $x_1$ and the other $x_j$ from the example in Figure 1.

affinities possessing a sufficient separation of scales will produce similar trees.

3. The result is completely independent of where the $\{x_1, \ldots, x_n\} \in \mathbb{R}^p$ are located in space. Their location, however, affects the critical scale $\varepsilon_0$.

4. The statement guarantees that points $u_i$ fuse together with respect to the folder structure before moving to fuse with other points and their respective folder structure, however, we do not have clear control over whether they intersect (in the sense of two $u_i, u_j$ belonging to different folders occupying the same point in space for some value of $\gamma$) in between or not. Generically, this will not happen but, for a non-generic set of $x_i$, it is possible to arrange for the $u_i$ to indeed intersect, then move apart again before finally fusing for a larger value of $\gamma$. This is a consequence of our lack of conditions on the position of the points $x_i$. If the $x_i$ are located in space in a way that actually reflects the tree structure, then they will fuse upon intersecting for the first time.

**3. A Geometric Lemma.** We establish a geometric Lemma that is of intrinsic interest: it states that for any set of distinct points $\{u_1, \ldots, u_n\} \in \mathbb{R}^p$, one of these points $u$ (indeed, one on the boundary of the convex hull of all the points) has the property that for a suitable "viewing direction" $v \in \mathbb{R}^p$ most points are clearly visible when standing in the point $u$ and looking towards the viewing direction (in the sense of having a large inner product). We now phrase this more precisely below. Recall that the convex hull of a set $S$, denoted by $\operatorname{conv} S$ is the smallest convex set containing the set $S$.

**Lemma 3.1.** *For every set $S = \{u_1, \ldots, u_n\} \subset \mathbb{R}^p$ of $n \geq 3$ distinct points, there exists*

$$u \in S \cap \partial \operatorname{conv} S \quad and \quad v \in \mathbb{R}^p \quad satisfying \ \|v\| = 1$$

240    *such that*

241    (3.1)
$$\frac{1}{n} \sum_{\substack{i=1 \\ u_i \neq u}}^{n} \left\langle \frac{u_i - u}{\|u_i - u\|}, v \right\rangle \geq \frac{1}{2}.$$

242    The statement can be summarized as follows: for a suitable point $u \in S \cap \partial \operatorname{conv} S$, if we
243    map the direction to all other points onto the unit sphere $\mathbb{S}^p$, then convexity implies that there
244    is a great circle on $\mathbb{S}^p$ such that all these directions are on one side of the great circle or on it.
245    This can be interpreted as the dualization of the fact that there is a supporting hyperplane
246    touching the boundary of the convex hull in such a way that all of $\operatorname{conv} S$ is on one side. The
247    statement claims the existence of a boundary point $u$ such that the average projection point
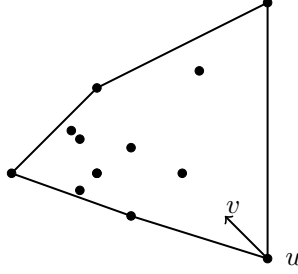248    is bounded away from that great circle by a universal constant.



Figure 7: A set of points in $\mathbb{R}^2$: there exists a point $u$ on the boundary of the convex hull and
a direction $v$ such that the average inner product of $(u_i - u)/\|u_i - u\|$ and $v$ is bounded away
from 0 by a universal constant.

249    We will use Lemma 3.1 to study the regularization term in (1.1), namely the functional

250
$$J(u) = \sum_{i,j=1}^{m} \|u_i - u_j\| \qquad \text{for a given set of distinct points } \{u_1, u_2, \ldots, u_m\} \subset \mathbb{R}^p.$$

251    The functional $J$ is clearly minimized for any collection of $u_i$ that are all identical. Con-
252    sequently, any collection of distinct $u_i$ represents a suboptimal configuration of centroids
253    and therefore admits a descent direction that leads to a decrease in energy. The power of
254    Lemma 3.1 is that it identifies a direction that guarantees a large amount of decrease in $J$.
255    To see this, we write down the directional derivative of $J$ explicitly.
256    The directional derivative of moving $u_j$ in direction $v \in \mathbb{R}^p$, normalized to $\|v\| = 1$ is

257     computed as

$$
\begin{aligned}
\left\langle \frac{\partial J}{\partial u_j}, v \right\rangle
&= \lim_{t \to 0} \frac{1}{t} \sum_{i \neq j} \| u_i - (u_j + tv) \| - \| u_i - u_j \| \\
&= \lim_{t \to 0} \frac{1}{t} \sum_{i \neq j} \sqrt{\langle u_i - (u_j + tv), u_i - (u_j + tv) \rangle} - \| u_i - u_j \| \\
&= \sum_{i \neq j} \lim_{t \to 0} \frac{1}{t} \left( \sqrt{\| u_i - u_j \|^2 - 2t \langle u_i - u_j, v \rangle + t^2} - \| u_j - u_i \| \right) \\
&= -\sum_{i \neq j} \left\langle \frac{u_i - u_j}{\| u_i - u_j \|}, v \right\rangle.
\end{aligned}
$$

258     (3.2)

259     The expression for the directional derivative given in (3.2), in conjunction with Lemma 3.1,
260     shows that it is always possible to find one point such that moving it $\delta$ in a certain direction
261     decreases the entire functional by at least $(n/2)\delta$. The existence of a direction of guaranteed
262     minimum decrease in $J$ will be essential in proving Theorem 2.1.

263

264         The following variant of Lemma 3.1 will also be useful in applications.

        **Lemma 3.2.** *For every set $S = \{u_1, \ldots, u_n\} \subset \mathbb{R}^p$ of $n \geq 3$ points such that not all of them
are in the same place, there exists*

$$
u \in S \cap \partial \operatorname{conv} S \quad and \quad v \in \mathbb{R}^p \quad satisfying \ \|v\| = 1
$$

265     *such that*

266     (3.3)
$$
\frac{1}{n} \sum_{\substack{i=1 \\ u_i \neq u}}^{n} \left\langle \frac{u_i - u}{\| u_i - u \|}, v \right\rangle \geq \frac{1}{4}.
$$

267         Before proceeding to proofs of the geometric lemmata and main result, we also note the
268     following consequence because of its intrinsic interest. We give a proof of Corollary 3.3 in
269     Appendix C.

270         **Corollary 3.3.** *Let $S = \{u_1, \ldots, u_n\} \subset \mathbb{R}^p$ be a set of distinct points. Then there exist at
271     least $n/6$ points $u \in S$ having the property that for some $\|v\| = 1$*

272
$$
\frac{1}{n} \sum_{\substack{i=1 \\ u_i \neq u}}^{n} \left\langle \frac{u_i - u}{\| u_i - u \|}, v \right\rangle \geq \frac{1}{4}.
$$

273         This simple statement has non-trivial implications: Lemma 3.1 may seem like these van-
274     tage points from which to observe the entirety of the set without having too many small inner
275     products are rare. To the contrary, Corollary 3.3 declares that the property is surprisingly
276     common and enjoyed by a universal fraction of all points. While we do not use Corollary 3.3
277     in the proof of our main result, we believe this result to be of substantial independent interest
278     since it can be interpreted as a basic statement (with universal constants) in a general Hilbert
279     space. It could be of interest to further pursue this line of investigation.

**4. Proofs.** We now prove Lemma 3.1, Lemma 3.2, and Theorem 2.1.

**4.1. Geometric Lemmata.**

*Proof of Lemma 3.1.* Let $S = \{u_1, u_2, \ldots, u_n\}$. Select an arbitrary $u \in \partial S \cap \operatorname{conv} S$, and let $y \in S$ be a point in the set furthest from $u$ (there may be more than one such point), formally

$$\|u - y\| = \max_{1 \leq i \leq n} \|u - u_i\| \tag{4.1}$$

It is easy to see that $y$ resides on the boundary of the convex hull; $y$ is in fact an extreme point. We now show that $u$, equipped with the viewing direction vector $v_1 = (y - u)/\|y - u\|$, or $y$, equipped with the viewing direction vector $v_2 = -v_1$, has the desired property. We first show that for every $u_i \notin \{u, y\}$

$$\left\langle \frac{u_i - u}{\|u_i - u\|}, v_1 \right\rangle + \left\langle \frac{u_i - y}{\|u_i - y\|}, v_2 \right\rangle \geq 1. \tag{4.2}$$

Since we are only dealing with three points $u, y$, and $u_i$, all angles are determined by the corresponding triangle, which we can assume without loss of generality to reside in $\mathbb{R}^2$. Moreover, the invariance under dilation, translation and rotation enables us to assume that $u = (0, 0)$ and $y = (1, 0)$. If we write $u_i = (a, b)$, then the expression on the left hand side of (4.2) simplifies to

$$\left\langle \frac{u_i - u}{\|u_i - u\|}, v_1 \right\rangle + \left\langle \frac{u_i - y}{\|u_i - y\|}, v_2 \right\rangle = \frac{a}{\sqrt{a^2 + b^2}} + \frac{1 - a}{\sqrt{(1 - a)^2 + b^2}}, \tag{4.3}$$

and the condition on the distances $\|u - u_i\|$ and $\|y - u_i\|$ required by (4.1) implies that

$$\max\left\{a^2 + b^2, (1 - a)^2 + b^2\right\} \leq 1. \tag{4.4}$$

Minimizing the expression in (4.3) subject to the constraint in (4.4) gives us the desired inequality in (4.2); almost equality is attained for $u_i$ very close to either $u$ or $y$ and equality is attained for $(a, b) = (1/2, \sqrt{3}/2)$. We then sum the left and right hand sides of (4.2) over $i = 1, \ldots, n$ to arrive at the inequality

$$\sum_{\substack{i=1 \\ u_i \neq u}}^{n} \left\langle \frac{u_i - u}{\|u_i - u\|}, v_1 \right\rangle + \sum_{\substack{i=1 \\ u_i \neq y}}^{n} \left\langle \frac{u_i - y}{\|u_i - y\|}, v_2 \right\rangle \geq n, \tag{4.5}$$

which follows from realizing that each of the sums contains one term that is equal to 1 and that the remaining sum runs over all $u_i \notin \{u, y\}$ yielding at least a total of $n - 2$. Thus at least one of the two terms is size $n/2$ and we obtain the desired result.  ∎

*Proof of Lemma 3.2.* Let $S = \{u_1, u_2, \ldots, u_n\}$ be a set of points not all of which are in the same place. Then the diameter of the set is not 0 and there exist two points, that we call w.l.o.g. $u_1$ and $u_2$ such that $\|u_1 - u_2\| = \operatorname{diam}(S)$. Let us suppose the number of points that

are co-located with $u_1$ is $n_1$, the number of points that are co-located with $u_2$ is $n_2$ and the number of points everywhere else is $n_3$. Clearly,

$$n_1 + n_2 + n_3 = n.$$

307  The main idea is now to derive two independent lower bounds. One of them will be tighter
308  when $n_1 + n_2$ is large (compared to $n$) and one will be tighter when $n_1 + n_2$ is small (compared
309  to $n$). We can then always apply the stronger of the two bounds and that will end up in
310  resulting a lower bound of $n/4$ regardless of what the values of $n_1$ and $n_2$ are.
311

**Bound 1.** We could pick $u$ to be $u_1$ and its viewing direction vector $v_1 = (u_2 - u_1)/\|u_2 - u_1\|$ or, conversely, the point $u_2$ and the vector $v_2 = (u_2 - u_1)/\|u_2 - u_1\|$ to be $u$ and $v$ respectively. We note that, since we chose the points to be of maximal distance, all arising inner products are nonnegative. Therefore

$$\sum_{\substack{i=1 \\ u_i \neq u_1}}^{n} \left\langle \frac{u_i - u_1}{\|u_i - u_1\|}, v_1 \right\rangle \geq n_2$$

and

$$\sum_{\substack{i=1 \\ u_i \neq u_2}}^{n} \left\langle \frac{u_i - u_2}{\|u_i - u_2\|}, v_2 \right\rangle \geq n_1.$$

312  Altogether, there is a pair of vectors $u$ and $v$ that achieves a sum of inner products of at least
313  $\max\{n_1, n_2\}$, which is a good bound when either of those two numbers is large (but true in
314  all cases). On the other hand, since we are only considering that small subset of points, the
315  bounds naturally become quite loose when $n_1 + n_2$ is small.
316

**Bound 2.** On the other hand, we can remove all the points co-located with either $u_1$ or $u_2$ except for one in each set, leaving us with $n - n_1 - n_2 + 2$ points. We can now apply the previous argument which guarantees the existence of a vector $u$ and a vector $v$ with

$$\sum_{\substack{i=1 \\ u_i \neq u}}^{n} \left\langle \frac{u_i - u}{\|u_i - u\|}, v \right\rangle \geq \frac{n - n_1 - n_2 + 2}{2}.$$

317  We see that this bound is quite good when $n_1$ and $n_2$ are small, in particular we recover the
318  original bound for distinct points whenever $n_1 = n_2 = 1$.
319

320  **Conclusion.** Having both bounds at our disposal, we can always guarantee the existence
321  of a pair $u$ and $v$ such that the lower bound is at least

322  
323  $$\max\left\{ \frac{n - n_1 - n_2 + 2}{2}, n_1, n_2 \right\} \geq \frac{1}{2}\left( \frac{n - n_1 - n_2 + 2}{2} + \frac{n_1 + n_2}{2} \right) \geq \frac{n}{4}$$

where the last line makes use of the inequality

$$\max\{x, y, z\} \geq \frac{x}{2} + \frac{y}{4} + \frac{z}{4} \qquad \text{for all } x, y, z \geq 0$$

324  since the maximum has to exceed every weighted average. ∎

**4.2. Main Theorem. Outline:** The proof is based on the self-similarity of the statement. We essentially show that points at the lowest level fuse in the right way with points in the same leaves (those who have mutual affinity 1). Once they are fused, we show that they stay fused for all subsequent values of $\gamma$. The newly emerging problem turns out to be exactly of the same type as the original one: we re-interpret fused points as single points with a mutual interaction now at scale $\sim \varepsilon$ (which becomes the dominant scale since points with $w_{ij} = 1$ are already fused). This makes crucial use of the geometry of the 1-norm. At every step, the arguments will go through provided $\varepsilon$ is sufficiently small (but positive) and since the tree is of finite height, the result follows. To be more precise, the argument will proceed as follows.

1. We assume that the $x_i$ are fixed and that the $u_i$ are solutions of the minimization problem

$$\inf_{u_1,\ldots,u_n} \left[ \sum_{i=1}^{n} \|x_i - u_i\|^2 + \gamma \sum_{i,j=1}^{n} w_{ij} \|u_i - u_j\| \right].$$

   Plugging in an example shows that the minimal energy is uniformly bounded in $\gamma$. This has some basic implications: the $u_i$ cannot be too far away from the $x_i$ and not too far away from each other.

2. We then study a subset of points $\{x_1, \ldots, x_n\}$ contained in a leaf of the tree. This means that their mutual affinity satisfies $w_{ij} = 1$ and the affinity between any of these points to any other point not in the leaf of the partition is at most $\varepsilon$.

3. We then focus exclusively on these point sets and prove that for $\gamma$ sufficiently large, these sets are necessarily fused in a point. This is where Lemma 3.2 will be applied.

4. Once we establish that for $\gamma$ sufficiently large, the point sets in the leaf are fused into exactly one point as desired, the full statement essentially follows by induction since these fused points interact exactly as individual points used to do; having common parents in the tree becomes the next-level analogue of being associated to the same leaf. The result then follows.

*Proof.* We introduce the energy of the minimal energy configuration for $\gamma > 0$ as

$$E(\gamma) = \inf_u E_\gamma(u) = \inf_u \left[ \sum_{i=1}^{n} \|x_i - u_i\|^2 + \gamma \sum_{i<j} w_{ij} \|u_i - u_j\| \right].$$

By setting $u_1 = u_2 = \cdots = u_n$ and putting these points in the center of mass of $\{x_1, \ldots, x_n\}$, we observe that this energy is uniformly bounded for all $\gamma$

$$E_{\sup} = \sup_{\gamma > 0} E(\gamma) \leq \sum_{i=1}^{n} \left\| x_i - \frac{1}{n} \sum_{i=1}^{n} x_i \right\|^2 < \infty.$$

We decompose the energy functional $E(\gamma)$ as

(4.6) $$E(\gamma) = E_1(\gamma) + E_2(\gamma),$$

where

$$E_1(\gamma) = \sum_{i=1}^{n} \|x_i - u_i\|^2 + \gamma \sum_{(i,j)\in\mathcal{E}_1} \|u_i - u_j\|,$$

356    where $\mathcal{E}_1 = \{(i,j) : w_{ij} = 1\}$ and

357
$$E_2(\gamma) = \gamma \sum_{(i,j)\in\mathcal{E}_2} w_{ij}\|u_i - u_j\|,$$

358    where $\mathcal{E}_2 = \{(i,j) : w_{ij} \le \varepsilon < 1\}$. The decomposition (4.6) makes explicit that, for $\varepsilon$ suffi-
359    ciently small, the functional $E_2(\gamma)$ can be interpreted as an error term, while the dominant
360    dynamics are determined by $E_1(\gamma)$. We now claim that for $\gamma$ sufficiently large (where suffi-
361    ciently large depends on everything except the parameter $\varepsilon$) any subset of the points $u_i$ whose
362    mutual affinities are 1 (i.e. all the members of one of the leaves in the tree) are fused in a
363    point. The argument can be made quantitative and we will give an explicit bound on $\gamma$ that
364    will be sufficient.
365
366        We will now ensure that we can assume that all points are distinct. The energy $E$ is a
367    continuous functional. This means that we can move any potentially clumped points apart
368    by accepting an arbitrarily small increase of energy; the remainder of the argument works
369    as follows: if points happen to be clumped together – but not in exactly one point but in
370    several – then we may move all of them an arbitrarily small bit. We can accept an arbitrarily
371    small increase of energy as long as we are able to then deduce a definite decrease in energy
372    afterwards (that will depend on the diameter of the $u_i$); this contradiction shows that the
373    clumping has to occur in exactly one point. The next step in the argument is dynamical: we
374    compute the effect of moving one of the points an infinitesimal amount (this is already using
375    the assumption that all $u_i$ are distinct). Reusing the computation in (3.2), we see that

376    (4.7)
$$\left\langle \frac{\partial E}{\partial u_j}, v \right\rangle = 2\left\langle u_j - x_j, v \right\rangle - \gamma \sum_{\substack{i=1 \\ i\neq j,(i,j)\in\mathcal{E}_1}}^{n} \left\langle \frac{u_i - u_j}{\|u_i - u_j\|}, v \right\rangle$$

377
$$+ \left\langle \frac{\partial}{\partial u_j}\gamma \sum_{(i,j)\in\mathcal{E}_2} w_{ij}\|u_i - u_j\|, v \right\rangle.$$

378        The first term on the right hand side of (4.7) is bounded above by

379    (4.8)
$$2\left|\langle u_j - x_j, v\rangle\right| \le 2\|x_j - u_j\| \le 2\sqrt{E_{\sup}},$$

380    and the third term on the right hand side of (4.7) is bounded above by

381    (4.9)
$$\left\|\frac{\partial}{\partial u_j}\gamma \sum_{(i,j)\in\mathcal{E}_2} w_{ij}\|u_i - u_j\|\right\| = \gamma\left\|\sum_{i:(i,j)\in\mathcal{E}_2,i\neq j} w_{ij}\frac{u_i - u_j}{\|u_i - u_j\|}\right\| \le \gamma\varepsilon n.$$

        Lemma 3.2 guarantees that there exists $u_j$ for which the second term on the right hand
    side of (4.7) is
$$-\gamma \sum_{\substack{i=1 \\ i\neq j,(i,j)\in\mathcal{E}_1}}^{n} \left\langle \frac{u_i - u_j}{\|u_i - u_j\|}, v \right\rangle \le -\frac{\gamma}{4}\#\left\{1 \le i \le n : (i,j) \in \mathcal{E}_1\right\}.$$

382   The proof of Lemma 3.1 is even stronger and guarantees that if $\|u_i - u_j\| = \operatorname{diam}\{u_1, \ldots, u_n\}$,
383   then either $u_i$ or $u_j$ has the desired property and can be moved in a suitable direction $v$.
384   Plugging the $u_j$ and $v$ from Lemma 3.1 into both sides of (4.7) and applying inequalities (4.8)
385   and (4.9), we arrive at the following inequality.

386   (4.10)     $$\left\langle \frac{\partial E}{\partial u_j}, v \right\rangle \le D(\gamma) = 2\sqrt{E_{\sup}} + \gamma \varepsilon n - \frac{\gamma}{4} \# \{1 \le i \le n : (i,j) \in \mathcal{E}_1\}.$$

A crucial observation is that for

$$\varepsilon < \frac{1}{4n} \# \{1 \le i \le n : (i,j) \in \mathcal{E}_1\}$$

we can conclude the existence of $\gamma$ sufficiently large (depending on all the other parameters)
so that $D(\gamma) < 0$. This, however, means the point configuration $\{u_1, \ldots, u_n\}$ cannot be a
minimizer of the functional since we found a point $u_j$ and a direction $v$ such that moving
$u_j$ into direction $v$ decreases the functional. This is a contradiction unless we are somehow
forbidden to apply Lemma 3.2: the only assumption in Lemma 3.2 is that not all points $u_i$
are in the same place. Thus we see that, for $\gamma$ sufficiently large, all points in $\mathcal{E}_1$ are fused. A
simple computation shows that these points have to be fused for all

$$\gamma \ge \frac{4\sqrt{E_{\sup}}}{\# \{1 \le i \le n : (i,j) \in \mathcal{E}_1\} - 4\varepsilon n}.$$

387   (This lower bound is not sharp; in practice, points will already be fused for smaller values of
388   $\gamma$.) A careful inspection of the proof shows that we do not require $w_{ij} = 1$ for points in the
389   same partition: it suffices if $1 \le w_{ij} \le c$ for some constant $c$ if subsequent parameter choices
390   of $\gamma$ are allowed to depend on that. The full statement now follows by induction: points in
391   leaves become a single point, their parent structure determines the next collection of leaves
392   and the product of their affinities determines the new affinities. Since there are only finitely
393   many levels to the tree, the process eventually terminates. ∎

394      **5. Extensions of the Main Theorem.** The proof of Theorem 2.1 relies on rather ele-
395   mentary analysis and consequently is quite flexible. Indeed, the proof can be immediately
396   extended to more general notions of energy of the type

397     $$E_\gamma(u) = \phi(x_1, \ldots, x_n, u_1, \ldots, u_n) + \gamma \sum_{i<j} w_{ij} \|u_i - u_j\|_X,$$

398   where $X$ is an arbitrary norm on $\mathbb{R}^p$ and $\phi$ is assumed to satisfy the following properties:
399       1. The function $\phi : \mathbb{R}^{p \times n} \to \mathbb{R}_{\ge 0}$ is differentiable and enforces some degree of data-fidelity
400          and compactness. More precisely, at one extreme $\phi$ should be minimized when $u_i = x_i$,
401          for example $\phi$ is nonnegative for all $u$ and $\phi(x_1, \ldots, x_n, x_1, \ldots, x_n) = 0$. At the other
402          extreme, $\phi$ should diverge whenever $\|u\|$ diverges. We want $\phi$ to have the property
403          of ensuring that minimizing the energy implies that all $u_i$ are trapped in a universal
404          convex set (determined by the $x_i$ but independent of $\gamma$). This amounts to a type of
405          growth condition on $\phi$ and many of the functions one would canonically choose will
406          have that property.

2. For all $u$ for which

$$\phi(x_1,\ldots,x_n,u_1,\ldots,u_n) + \gamma \sum_{i,j=1}^{n} w_{ij} \|u_i - u_j\|_X \le \inf_{x \in \mathbb{R}^p} \phi(x_1,\ldots,x_n,x,\ldots,x),$$

we have

$$\left\| \frac{\partial}{\partial u_i} \phi(x_1,\ldots,x_n,u_1,\ldots,u_n) \right\| \le c$$

where $c$ only depends on $\gamma$ and $\{x_1,\ldots,x_n\}$.

The argument proceeds in exactly the same way and makes crucial use of the fact that any two norms in a finite-dimensional Euclidean space are equivalent up to constants, namely

$$c_5\|x\|_{\ell^2} \le \|x\|_X \le c_6\|x\|_{\ell^2}.$$

Since constants can always be absorbed in $\gamma$, this reduces to our case, namely $X = \ell^2$.

*Proof.* (Sketch of the argument) Setting all $u_i = x$ and minimizing over $x$ implies that the energy is uniformly bounded in $\gamma$ (with a bound depending only on $\{x_1,\ldots,x_n\}$). Since the norm $X$ is comparable to the Euclidean norm, this implies that any minimizing configuration $\{u_1,\ldots,u_n\}$ has to have a bounded diameter (with a bound depending only on $\{x_1,\ldots,x_n\}$). Then, for $\gamma$ sufficiently large (depending on $c$), Lemma 3.1 implies a direction of decay and thus points are eventually fused. We leave the precise details to the interested reader. ∎

We close this section by noting that the generality of our result opens the door to intriguing applications. For example, one potential application of our extension is to construct partition trees of regression coefficients in clustered regression [5, 22, 39, 48]. We leave these investigations as future work.

**6. Convex Clustering in High-dimensional Spaces.** We now briefly provide some practical guidance in using convex clustering in high-dimensional spaces. Beyer et al. showed in [4] that over a broad class of data distributions, as the ambient dimensional increases, distances from a point to its nearest neighbors become indistinguishable from distances to its farthest neighbors. Thus, at first glance, it is unclear whether tree organizations can be recovered from high-dimensional data using convex clustering, a method in which distance metrics play a central role. Fortunately, many high-dimensional data sets encountered in engineering and science can be approximated reliably by a lower dimensional representation or embedding. In some cases, high-dimensional data consist of many features that contain little to no information about the clustering structure and should be dropped. In this case, one may consider computing a sparse convex clustering solution path [46]. In other cases, where there are more nuanced relationships among most or even all the features, we may turn to nonlinear dimension reduction methods. Indeed, manifold learning [3, 13, 15, 43, 35] has proven to be effective as a nonlinear dimension reduction technique in many scientific domains where very high-dimensional measurements are recorded such as in bioinformatics [17, 20, 27, 50] and neuroscience [7, 6, 8, 36, 40, 45]. Upon some reflection, this is not surprising, as these studies

443   collect high-dimensional data that are generated from natural processes that are subject to
444   physical constraints and are thus intrinsically low-dimensional.
445       In light of these observations, we recommend the following simple strategy. First, embed
446   high-dimensional data into a low-dimensional space, and then compute a convex clustering
447   solution path using the low-dimensional representation of the data. This strategy is especially
448   natural if one uses diffusion maps, since the diffusion distance between two points in high-
449   dimensions can be approximated by the Euclidean distance in the lower dimensional diffusion
450   maps space [13]. Once points are embedded in the diffusion maps space, one can use Gaussian
451   kernel affinities and compute the convex clustering solution path using the Euclidean norm in
452   the regularization term.

453       **7. Discussion.** In this paper, we answered the question of when the convex clustering
454   solution path can recover a tree. The key to ensuring the recovery of a well nested partition
455   tree is the use of affinities that encourage the fusions within a folder before fusions with
456   higher level folders and so on as the tuning parameter $\gamma$ increases. By choosing the edge
457   weight parameter $\varepsilon$ sufficiently small, different folders have very little incentive to interact,
458   and the optimization problem is essentially decoupled. As $\gamma$ increases, the same procedure
459   repeats itself.
460       We end with a discussion on the relationship between convex and non-convex formulations
461   of penalized regression based clustering. Although we focus in this paper on the ability of
462   convex clustering to recover a potentially deep hierarchy of nested folders, our result also sheds
463   light on a gap in theory and practice that convex clustering's performance can be significantly
464   improved when using non-uniform data-driven affinities when seeking a shallow or single level
465   of nested folders. In practice, Gaussian kernel affinities have been observed to work well, but
466   these affinity choices have until now lacked formal justification.
467       Indeed, non-uniform affinities provide the link between convex clustering and other penal-
468   ized regression-based clustering methods that use folded concave penalties. It is well known
469   that 1-norm penalties lead to parameter estimates that are shrunk towards zero. This shrink-
470   age toward zero is the price for simultaneously estimating the support, or locations of the
471   nonzero entries, in a sparse vector as well the values of the nonzero entries. In the context
472   of convex clustering, the centroid estimates $u_i$ are shrunk towards the grand mean $\overline{x}$. Con-
473   sequently, others have proposed employing a folded concave penalty instead of a norm in the
474   regularization terms [31, 26, 49]. Folded concave penalties induce milder shrinkage in exchange
475   for giving up convexity in the optimization problem, which means that iterative algorithms
476   can typically at best converge only to a KKT point.
477       Suppose we were to employ a folded concave penalty, such as the smoothly clipped absolute
478   deviation [16] or minimax concave penalty [53], and seek to minimize the following alternative
479   objective to (1.1)

480   (7.1)
$$\tilde{E}_\gamma(u) = \frac{1}{2}\sum_{i=1}^{n}\|x_i - u_i\|^2 + \gamma\sum_{i<j}\varphi\left(\|u_i - u_j\|\right),$$

481   where each $\varphi : [0, \infty) \mapsto [0, \infty)$ has the following properties: (i) $\varphi$ is concave and differentiable
482   on $(0, \infty)$, (ii) $\varphi$ vanishes at the origin, and (iii) the directional derivative of $\varphi$ exists and is
483   positive at the origin.

Since $\varphi$ is concave and differentiable, for all positive $z$ and $\tilde{z}$

$$\varphi(z) \le \varphi(\tilde{z}) + \varphi'(\tilde{z})(z - \tilde{z}).$$

In other words, the first order Taylor expansion of a differentiable concave function $\varphi$ provides a tight global upper bound at the expansion point $\tilde{z}$. Thus, we can construct a function that is a tight upper bound of the function $\tilde{E}_\gamma(u)$

(7.2)
$$g_\gamma(u \mid \tilde{u}) = \frac{1}{2} \sum_{i=1}^{n} \|x_i - u_i\|^2 + \gamma \sum_{i<j} w_{ij} \|u_i - u_j\| + c_7,$$

where $c_7$ is a constant that does not depend on $u$, and $w_{ij}$ are affinities that depend on $\tilde{u}$, namely

$$w_{ij} = \varphi'\left(\|\tilde{u}_i - \tilde{u}_j\|\right).$$

Note that if we take $\tilde{u}_i$ to be the data $x_i$, and $\varphi(z)$ to be the following variation on the error function

$$\varphi(z) = \int_0^z e^{-\frac{\alpha^2}{\sigma}} \, d\alpha,$$

then the bounding function given in (7.2) coincides, up to an irrelevant shift and scaling, with the convex clustering objective using Gaussian kernel affinities.

The function $g_\gamma(u \mid \tilde{u})$ is said to majorize the function $\tilde{E}_\gamma(u)$ at the point $\tilde{u}$ [24] and minimizing it corresponds to performing one step of the local linear-approximation algorithm [37, 55], which is a special case of the majorization-minimization (MM) algorithm [24]. Thus, we can see that employing Gaussian kernel affinities corresponds to taking one step of a local linear-approximation algorithm applied to a penalized regression based clustering with an appropriately chosen folded concave penalty.

In practice, variants that employ folded concave penalties take multiple steps of the local linear approximation. So at the $k$th step,

$$u^{(k)} = \arg\min_u \frac{1}{2} \sum_{i=1}^{n} \|x_i - u_i\|^2 + \gamma \sum_{i<j} \varphi'\left(\|u_i^{(k-1)} - u_j^{(k-1)}\|\right) \|u_i - u_j\|.$$

As affinities represent a data-driven way to approximate the partition tree, one can see that employing folded concave penalties corresponds to implicitly recomputing the affinities, which corresponds to refining our estimate of the partition tree based on the data.

In light of this current work, this last observation raises two interesting questions: (i) what partition tree is being recovered by a solution path of a penalized regression-based clustering method that uses a folded concave penalty and (ii) when is the recovered partition tree substantially different than the tree corresponding to a one-step local linear approximation? We leave these questions to future work.

**Appendix A. Example of Non-Tree Solution Path.**

We recreate a configuration of points in $\mathbb{R}^2$ and affinities similar to those used in [19], which yield a solution path that is not a tree. Consider the following four points, $x_1 = (-0.25, 3), x_2 = (0.25, 3), x_3 = (2, 0)$, and $x_4 = (-2, 0)$, and employ affinities $w_{12} = 9, w_{13} = w_{24} = 30$, and $w_{ij} = 1$ for all remaining $i$ and $j$ pairs. Figure 8 shows snapshots of the evolution of the solution paths for $u_1(\gamma)$ (red), $u_2(\gamma)$ (blue), $u_3(\gamma)$ (green), and $u_4(\gamma)$ (purple) as $\gamma$ increases. We see that $u_1(\gamma) = u_2(\gamma)$ for a continuous range of $\gamma$ greater than $10^{-2.05}$ and strictly less than $10^{-1.64}$ (Figure 8d and Figure 8e) but that $u_1(\gamma) \neq u_2(\gamma)$ for a continuous range of $\gamma$ greater than $10^{-1.64}$ and less than $10^{-0.85}$ (Figure 8e, Figure 8f, and Figure 8g).



(a) $\gamma = 0$    (b) $\gamma = 10^{-2.63}$    (c) $\gamma = 10^{-2.19}$    (d) $\gamma = 10^{-2.05}$

(e) $\gamma = 10^{-1.64}$    (f) $\gamma = 10^{-1.39}$    (g) $\gamma = 10^{-1.12}$    (h) $\gamma = 10^{-0.85}$
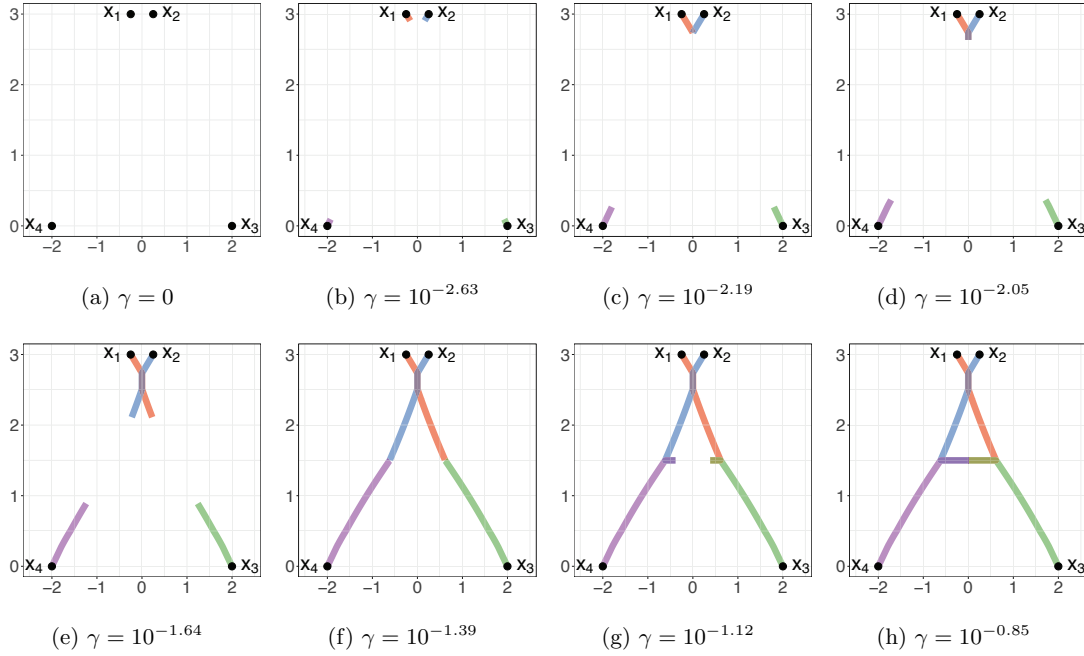
Figure 8: Snapshots of the solution path as the parameter $\gamma$ increases.

We emphasize that in order to generate this degenerate solution path, we needed to use affinities that *do not* reflect the geometry of the data. The largest affinities, $w_{13}$ and $w_{24}$, are between the two pairs of points that are furthest apart from each other.

**Appendix B. Comparison of Unit versus Gaussian Kernel Affinities on Vote Data.**

To illustrate the superiority of Gaussian kernel affinities over unit affinities often observed on real data, we compute the convex clustering solution paths under the two kinds of affinities on US senate voting data in 2001 [1, 14]. We removed duplicate voting records, restricting our attention to 29 senators – 15 Democrats, 13 Republicans, and 1 Independent (Jim Jeffords, who was a Republican prior to 2001) – and their votes on 13 issues ranging over domestic, foreign, economic, military, environmental, and social concerns. The raw data consisted of 29
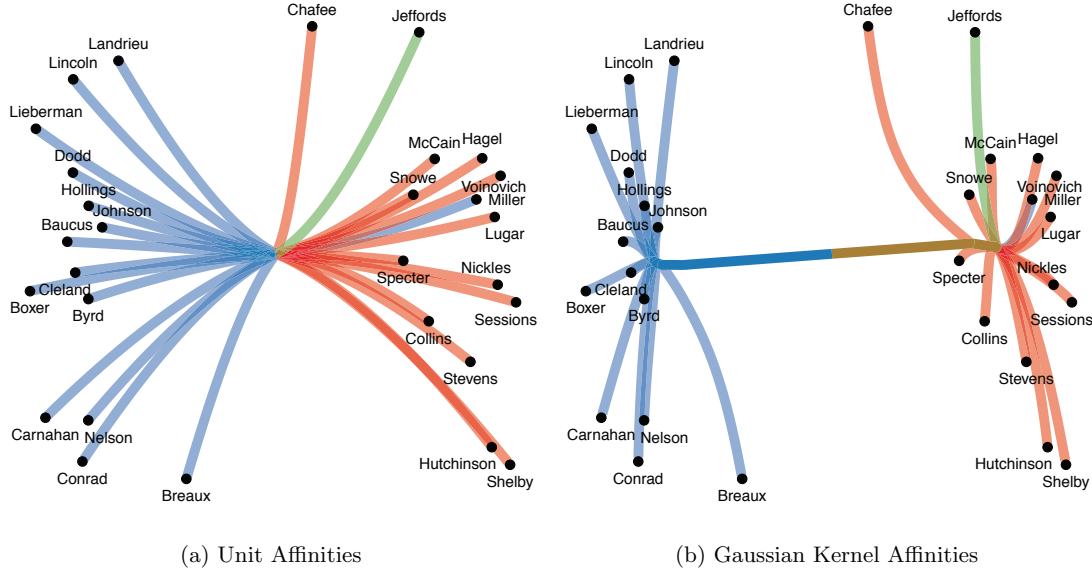
(a) Unit Affinities

(b) Gaussian Kernel Affinities

Figure 9: US Senate Vote Data: Solution path as the parameter $\gamma$ increases.

534  binary vectors of length 13, which we centered and scaled. Figure 9 shows the solution paths
535  under the two kinds of affinities; for visualization purposes we projected $u_i(\gamma) \in \mathbb{R}^{13}$ onto the
536  first two principal components of the centered and scaled data matrix. We color coded the
537  solution paths to reflect senator party affiliations: Democrats in blue, Republicans in red, and
538  the Independent in green. As an aside, we identify an outlying Democrat in Zell Miller, who
539  had a track record for supporting Republican policies during his tenure. He notably supported
540  Republican President George W. Bush against John Kerry, the Democratic nominee in the
541  2004 presidential election.
542      Figure 9a and Figure 9b show the resulting clustering paths under unit affinities, $w_{ij} = 1$
543  for all $i$ and $j$, and Gaussian kernel affinities respectively. In the latter case, we use a commonly
544  used data-driven strategy of choosing a local scale parameter $\sigma_{ij}$ that is pair dependent [52],
545  namely

546
$$w_{ij} = \exp\left(-\frac{\|x_i - x_j\|_2^2}{\sigma_{ij}}\right).$$

547  We first compute a local measure of scale $\sigma_i$, which is the median Euclidean distance between
548  the $i$th point $x_i$ and its 5-nearest neighbors. We then set $\sigma_{ij} = \sigma_i \sigma_j$.
549      The solution path in Figure 9a exhibits exactly *one* fusion event as $\gamma$ increases, namely at
550  the end of the solution path. In contrast, the solution path in Figure 9b exhibits fusions that
551  initially group together senators in their respective parties, before the two main groups fuse
552  at the end of the solution path. Figure 10a and Figure 10b show points along the solution
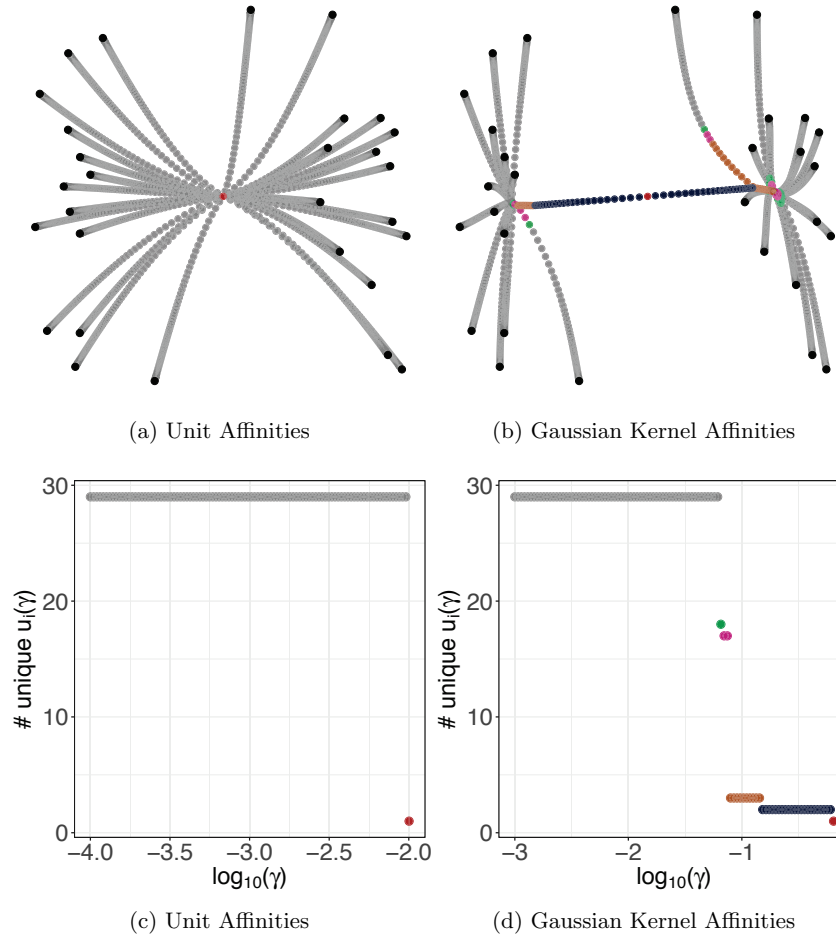553  paths obtained from unit and Gaussian kernel affinities respectively, color coded according to

(a) Unit Affinities            (b) Gaussian Kernel Affinities

(c) Unit Affinities            (d) Gaussian Kernel Affinities

Figure 10: US Senate Vote Data: The number of unique $u_i(\gamma)$ as a function of $\gamma$.

the number of unique $u_i(\gamma)$ as $\gamma$ varies. Figure 10c and Figure 10d plot the number of unique
$u_i(\gamma)$ as $\gamma$ varies under unit and Gaussian kernel affinities respectively. Indeed, we see that in
this real example, the unit affinities produce a rather useless tree, namely one with *no* nesting
at all. In contrast, the Gaussian kernel affinities produce a tree that organizes the senators
into partitions that respect party affiliations. Figure 10b also shows that John Chaffee, who
was one of the more liberal Republicans, fuses somewhat later to the Republican group and
also shows that John Breaux, whose centrist voting tendencies at times led Republicans to
seek his help in swaying a few critical Democratic votes, fuses somewhat later to the Democrat
group.

## Appendix C. Proof of Corollary 3.3.

*Proof.* Lemma 3.1 guarantees the existence of a point $u$, call it $\tilde{u}_1$, and viewing direction vector $v_1$ that satisfies inequality (3.1). Remove $\tilde{u}_1$ from the set $S = \{u_1, \ldots, u_n\}$ and apply Lemma 3.1 to the new set $S \backslash S_1$, where $S_1 = \{\tilde{u}_1\}$. Repeat this procedure $k$ times and let $S_k$ denote the set of $k$ points, $\{\tilde{u}_1, \ldots, \tilde{u}_k\}$, that satisfy inequality (3.1) for the sets $S, S \backslash S_1, \ldots, S \backslash S_{k-1}$ respectively. Lemma 3.1 guarantees the existence of a point $u \in S \backslash S_k$ and viewing direction vector $v$ such that

(C.1)
$$\frac{1}{n-k} \sum_{\substack{u_i \in S \backslash S_k \\ u_i \neq u}} \left\langle \frac{u_i - u}{\|u_i - u\|}, v \right\rangle \geq \frac{1}{2}.$$

The Cauchy-Bunyakovsky-Schwarz inequality tells us that

(C.2)
$$\left\langle \frac{u_i - u}{\|u_i - u\|}, v \right\rangle \geq -1,$$

for all $u_i \in S_k$. Inequalities (C.1) and (C.2) together imply that

(C.3)
$$\sum_{\substack{i=1 \\ u_i \neq u}}^{n} \left\langle \frac{u_i - u}{\|u_i - u\|}, v \right\rangle \geq \frac{n-k}{2} - k$$

Finally, for $k \leq n/6$, we see that the right hand side of (C.3) is bounded below by $n/4$ which implies the desired result. ∎

**Acknowledgments.** We thank Raphy Coifman for pointing out Corollary 3.3.

## REFERENCES

[1] AMERICANS FOR DEMOCRATIC ACTION, *2001 voting record: Shattered promise of liberal progress*, ADA Today, 57 (2002), pp. 1–17.

[2] J. I. ANKENMAN, *Geometry and Analysis of Dual Networks on Questionnaires*, PhD thesis, Yale University, 2014.

[3] M. BELKIN AND P. NIYOGI, *Laplacian Eigenmaps for Dimensionality Reduction and Data Representation*, Neural Computation, 15 (2003), pp. 1373–1396.

[4] K. S. BEYER, J. GOLDSTEIN, R. RAMAKRISHNAN, AND U. SHAFT, *When Is "Nearest Neighbor" Meaningful?*, in Proceedings of the 7th International Conference on Database Theory, ICDT '99, London, UK, UK, 1999, Springer-Verlag, pp. 217–235.

[5] H. D. BONDELL AND B. J. REICH, *Simultaneous Regression Shrinkage, Variable Selection, and Supervised Clustering of Predictors with OSCAR*, Biometrics, 64 (2008), pp. 115–123.

[6] B. M. BROOME, V. JAYARAMAN, AND G. LAURENT, *Encoding and Decoding of Overlapping Odor Sequences*, Neuron, 51 (2006), pp. 467–482.

[7] S. L. BROWN, J. JOSEPH, AND M. STOPFER, *Encoding a Temporally Structured Stimulus with a Temporally Structured Neural Representation*, Nature Neuroscience, 8 (2005), pp. 1568–76.

[8] L. CARRILLO-REID, F. TECUAPETLA, D. TAPIA, A. HERNÁNDEZ-CRUZ, E. GALARRAGA, R. DRUCKER-COLIN, AND J. BARGAS, *Encoding Network States by Striatal Cell Assemblies*, Journal of Neurophysiology, 99 (2008), pp. 1435–1450.

[9] G. K. CHEN, E. C. CHI, J. M. RANOLA, AND K. LANGE, *Convex Clustering: An Attractive Alternative to Hierarchical Clustering*, PLoS Computational Biology, 11 (2015), p. e1004228.

[10] E. C. CHI, G. I. ALLEN, AND R. G. BARANIUK, *Convex Biclustering*, Biometrics, 73 (2017), pp. 10–19.

[11] E. C. CHI, B. R. GAINES, W. W. SUN, H. ZHOU, AND J. YANG, *Provable convex co-clustering of tensors.* arXiv:1803.06518 [stat.ME], 2018.

[12] E. C. CHI AND K. LANGE, *Splitting Methods for Convex Clustering*, Journal of Computational and Graphical Statistics, 24 (2015), pp. 994–1013.

[13] R. R. COIFMAN AND S. LAFON, *Diffusion Maps*, Applied and Computational Harmonic Analysis, 21 (2006), pp. 5–30.

[14] J. DE LEEUW AND P. MAIR, *Gifi methods for optimal scaling in R: The Package homals*, Journal of Statistical Software, 31 (2009), pp. 1–21.

[15] D. L. DONOHO AND C. GRIMES, *Hessian Eigenmaps: Locally Linear Embedding Techniques for High-dimensional Data*, Proceedings of the National Academy of Sciences, 100 (2003), pp. 5591–5596.

[16] J. FAN AND R. LI, *Variable selection via nonconcave penalized likelihood and its oracle properties*, Journal of the American Statistical Association, 96 (2001), pp. 1348–1360.

[17] J. M. GARCÍA-GÓMEZ, J. GÓMEZ-SANCHIS, P. ESCANDELL-MONTERO, E. FUSTER-GARCIA, AND E. SORIA-OLIVAS, *Sparse Manifold Clustering and Embedding to Discriminate Gene Expression Profiles of Glioblastoma and Meningioma Tumors*, Computers in Biology and Medicine, 43 (2013), pp. 1863–1869.

[18] J. C. GOWER AND G. J. S. ROSS, *Minimum spanning trees and single linkage cluster analysis*, Journal of the Royal Statistical Society. Series C (Applied Statistics), 18 (1969), pp. 54–64.

[19] T. D. HOCKING, A. JOULIN, F. BACH, AND J.-P. VERT, *Clusterpath an algorithm for clustering using convex fusion penalties*, in Proceedings of the 28th International Conference on Machine Learning (ICML-11), 2011, pp. 745–752.

[20] X. JIANG, X. HU, H. SHEN, AND T. HE, *Manifold Learning Reveals Nonlinear Structure in Metagenomic Profiles*, in 2012 IEEE International Conference on Bioinformatics and Biomedicine, 2012, pp. 1–6.

[21] S. C. JOHNSON, *Hierarchical clustering schemes*, Psychometrika, 32 (1967), pp. 241–254.

[22] Z. T. KE, J. FAN, AND Y. WU, *Homogeneity pursuit*, Journal of the American Statistical Association, 110 (2015), pp. 175–194.

[23] G. N. LANCE AND W. T. WILLIAMS, *A general theory of classificatory sorting strategies: 1. hierarchical systems*, The Computer Journal, 9 (1967), pp. 373–380.

[24] K. LANGE, D. R. HUNTER, AND I. YANG, *Optimization transfer using surrogate objective functions*, Journal of Computational and Graphical Statistics, 9 (2000), pp. 1–20.

[25] F. LINDSTEN, H. OHLSSON, AND L. LJUNG, *Just relax and come clustering! A convexification of k-means clustering*, tech. report, Linköpings Universitet, 2011.

[26] Y. MARCHETTI AND Q. ZHOU, *Solution path clustering with adaptive concave penalty*, Electronic Journal of Statistics, 8 (2014), pp. 1569–1603.

[27] E. MARRAS, A. TRAVAGLIONE, AND E. CAPOBIANCO, *Manifold Learning in Protein Interactomes*, Journal of Computational Biology, 18 (2010), pp. 81–96.

[28] G. MISHNE, R. TALMON, I. COHEN, R. R. COIFMAN, AND Y. KLUGER, *Data-driven tree transforms and metrics*, IEEE Transactions on Signal and Information Processing over Networks, 4 (2018), pp. 451–466.

[29] G. MISHNE, R. TALMON, R. MEIR, J. SCHILLER, M. LAVZIN, U. DUBIN, AND R. R. COIFMAN, *Hierarchical coupled-geometry analysis for neuronal structure and activity pattern discovery*, IEEE Journal of Selected Topics in Signal Processing, 10 (2016), pp. 1238–1253.

[30] F. MURTAGH, *A survey of recent advances in hierarchical clustering algorithms*, The Computer Journal, 26 (1983), pp. 354–359.

[31] W. PAN, X. SHEN, AND B. LIU, *Cluster analysis: Unsupervised learning via supervised learning with a non-convex penalty*, Journal of Machine Learning Research, 14 (2013), pp. 1865–1889.

[32] A. PANAHI, D. DUBHASHI, F. D. JOHANSSON, AND C. BHATTACHARYYA, *Clustering by Sum of Norms: Stochastic Incremental Algorithm, Convergence and Cluster Recovery*, in Proceedings of the 34th International Conference on Machine Learning, D. Precup and Y. W. Teh, eds., vol. 70 of Proceedings of Machine Learning Research, International Convention Centre, Sydney, Australia, 06–11 Aug 2017, PMLR, pp. 2769–2777.

[33] K. PELCKMANS, J. DE BRABANTER, J. SUYKENS, AND B. DE MOOR, *Convex clustering shrinkage*, in PASCAL Workshop on Statistics and Optimization of Clustering Workshop, 2005.

[34] P. RADCHENKO AND G. MUKHERJEE, *Convex clustering via l1 fusion penalization*, Journal of the Royal Statistical Society: Series B (Statistical Methodology), 79 (2017), pp. 1527–1546.

[35] S. T. ROWEIS AND L. K. SAUL, *Nonlinear Dimensionality Reduction by Locally Linear Embedding*, Science, 290 (2000), pp. 2323–2326.

[36] D. SAHA, K. LEONG, C. LI, S. PETERSON, G. SIEGEL, AND B. RAMAN, *A Spatiotemporal Coding Mechanism for Background-Invariant Odor Recognition*, Nature Neuroscience, 16 (2013), pp. 1830–1839.

[37] E. D. SCHIFANO, R. L. STRAWDERMAN, AND M. T. WELLS, *Majorization-minimization algorithms for nonsmoothly penalized objective functions*, Electronic Journal of Statistics, 4 (2010), pp. 1258–1299.

[38] J. SHARPNACK, A. SINGH, AND A. RINALDO, *Sparsistency of the edge lasso over graphs*, in AISTATS, 2012.

[39] Y. SHE, *Sparse regression with exact clustering*, Electronic Journal of Statistics, 4 (2010), pp. 1055–1096.

[40] M. STOPFER, V. JAYARAMAN, AND G. LAURENT, *Intensity versus Identity Coding in an Olfactory System*, Neuron, 39 (2003), pp. 991–1004.

[41] D. SUN, K.-C. TOH, AND Y. YUAN, *Convex Clustering: Model, Theoretical Guarantee and Efficient Algorithm.* arXiv:1810.02677 [cs.LG], 2018.

[42] K. M. TAN AND D. WITTEN, *Statistical properties of convex clustering*, Electronic Journal of Statistics, 9 (2015), pp. 2324–2347.

[43] J. B. TENENBAUM, V. D. SILVA, AND J. C. LANGFORD, *A Global Geometric Framework for Nonlinear Dimensionality Reduction*, Science, 290 (2000), pp. 2319–2323.

[44] R. TIBSHIRANI, M. SAUNDERS, S. ROSSET, J. ZHU, AND K. KNIGHT, *Sparsity and smoothness via the fused lasso*, Journal of the Royal Statistical Society: Series B (Statistical Methodology), 67 (2005), pp. 91–108.

[45] J. T. VOGELSTEIN, Y. PARK, T. OHYAMA, R. A. KERR, J. W. TRUMAN, C. E. PRIEBE, AND M. ZLATIC, *Discovery of Brainwide Neural-Behavioral Maps via Multiscale Unsupervised Structure Learning*, Science, 344 (2014), pp. 386–392.

[46] B. WANG, Y. ZHANG, W. W. SUN, AND Y. FANG, *Sparse convex clustering*, Journal of Computational and Graphical Statistics, 27 (2018), pp. 393–403.

[47] J. H. WARD, *Hierarchical grouping to optimize an objective function*, Journal of the American Statistical Association, 58 (1963), pp. 236–244.

[48] D. M. WITTEN, A. SHOJAIE, AND F. ZHANG, *The cluster elastic net for high-dimensional regression with unknown variable grouping*, Technometrics, 56 (2014), pp. 112–122.

[49] C. WU, S. KWON, X. SHEN, AND W. PAN, *A new algorithm and theory for penalized regression-based clustering*, Journal of Machine Learning Research, 17 (2016), pp. 1–25.

[50] Z.-H. YOU, Y.-K. LEI, J. GUI, D.-S. HUANG, AND X. ZHOU, *Using Manifold Embedding for Assessing and Predicting Protein Interactions from High-throughput Experimental Data*, Bioinformatics, 26 (2010), pp. 2744–2751.

[51] M. YUAN AND Y. LIN, *Model selection and estimation in regression with grouped variables*, Journal of the Royal Statistical Society: Series B (Statistical Methodology), 68 (2006), pp. 49–67.

[52] L. ZELNIK-MANOR AND P. PERONA, *Self-tuning spectral clustering*, in Advances in Neural Information Processing Systems 17, L. K. Saul, Y. Weiss, and L. Bottou, eds., MIT Press, 2005, pp. 1601–1608.

[53] C.-H. ZHANG, *Nearly unbiased variable selection under minimax concave penalty*, The Annals of Statistics, 38 (2010), pp. 894–942.

[54] C. ZHU, H. XU, C. LENG, AND S. YAN, *Convex optimization procedure for clustering: Theoretical revisit*, in Advances in Neural Information Processing Systems 27, Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, eds., 2014, pp. 1619–1627.

[55] H. ZOU AND R. LI, *One-step sparse estimates in nonconcave penalized likelihood models*, The Annals of Statistics, 36 (2008), pp. 1509–1533.