

Kernel Treelets*

Hedi Xia[†] and Hector D. Ceniceros[†]

Abstract. A new method for hierarchical clustering of data points is presented. It combines treelets, a particular multiresolution decomposition of data, with a projection on a reproducing kernel Hilbert space. The proposed approach, called kernel treelets (KT), effectively substitutes the correlation coefficient matrix used in treelets with a symmetric, positive semi-definite matrix efficiently constructed from a symmetric, positive semi-definite kernel function. Unlike most clustering methods, which require datasets to be numeric, KT can be applied to more general data and yields a multi-resolution sequence of orthonormal bases on the data directly in feature space. The effectiveness and potential of KT in clustering analysis is illustrated with some examples.

1. Introduction. Treelets, introduced by Lee, Nadler, and Wasserman [1, 2], is a method to produce a multiscale, hierarchical representation of unordered data. The central premise of treelets and the treelet transform is to exploit sparsity and capture intrinsic localized structures with only a few features (attributes), represented in terms of an orthonormal basis. The treelet transform consists of a sequence of two-dimensional principal component analyses implemented efficiently via rotations. The resulting multiresolution representation of the data can be used for dimensionality reduction and for feature selection prior to regression/classification [1, 2].

Cluster analysis, also called clustering [3], is one of the basic tasks of unsupervised learning and is concerned with finding a partition of a set so that elements in one cluster are more similar to one another than they are to elements in another cluster, i.e. the corresponding equivalence class captures similarity of its elements. The clustering can be flat, where the partition is a collection of disjoint sets, or hierarchical [4], where a nested tree of partitions is produced. The treelet transform produces a hierarchical clustering *over attributes*. In this work, we propose to combine the kernel method [10, 11] with the treelet transform to obtain an efficient tool for hierarchical clustering analysis *over data points*. We call this method kernel treelet (KT). The central idea is to project the data onto a reproducing kernel Hilbert space (RKHS). This effectively substitutes the correlation coefficient matrix, used by the original treelet method as a measure of similarity among attributes, with a symmetric, positive semi-definite matrix that measures similarity among data points. The intuition behind this approach is that inner products provide a measure of data similarity and a projection onto a RKHS, done via the so-called kernel trick [10, 11], is a natural and efficient way to construct appropriate (dis)similarity matrices for a wide variety of datasets. We present some examples that demonstrate the potential of KT as an effective tool for data clustering analysis.

The typical complexity of hierarchical clustering methods is $O(n^3)$ (n denotes the number of data points in the dataset) but KT, like single linkage clustering [5], and complete linkage clustering [6] can be done in $O(n^2)$ operations. Most clustering methods are only directly applicable to numerical datasets. However, many modern datasets do not have clear represen-

*Submitted to the editors February 27, 2019

[†]Department of Mathematics, University of California Santa Barbara, 93106.

40 tations in \mathbb{R}^d due for instance to missing data, length difference, and non-numeric attributes.
 41 A typical solution to this problem usually involves finding a projection from each observation
 42 to \mathbb{R}^d as is the case for example in text vectorization [7], array alignment [8], and missing-data
 43 imputation [9]. However, these particular projections pose considerable challenges and might
 44 raise the bias of the model if false assumptions are made. KT does not have this limitation
 45 and can be applied to an ample range of datasets, including those mentioned above.

46 The rest of the paper is as follows. In Section 2 gives some basic background for the
 47 treelet transform and the kernel method. This is followed by the introduction of the KT
 48 model in Section 3. Section 4 presents some theory to help explain the success of the KT
 49 approach for clustering. Three examples of clustering analysis are given in section 5. In
 50 Section 6 an approach that combines KT with supervised learning is proposed to accelerate
 51 data hierarchical clustering. Finally, concluding remarks are given in Section 7.

52 **2. Background Information.** We give in this section a brief description of the treelet
 53 algorithm [1, 2] and the Kernel method [12] as background for the introduction of the KT
 54 model. Treelets are based on the repeated application of two dimensional rotations to a matrix
 55 measuring the similarity of attributes. So we start by reviewing Jacobi (also called Givens)
 56 rotations first.

A Jacobi rotation matrix J is an orthogonal matrix with at most 4 entries different from
 the identity. For a given symmetric matrix M and entry pq , the Jacobi matrix J is constructed
 so that

$$(J^T M J)_{pq} = (J^T M J)_{qp} = 0.$$

57 The construction of J is equivalent to finding the cosine (c) and sine (s) of the angle of rotation,
 58 which satisfy

$$59 \begin{bmatrix} c & -s \\ s & c \end{bmatrix} \begin{bmatrix} M_{pp} & M_{pq} \\ M_{qp} & M_{qq} \end{bmatrix} \begin{bmatrix} c & s \\ -s & c \end{bmatrix} = \begin{bmatrix} d_1 & 0 \\ 0 & d_2 \end{bmatrix}$$

60 subject to the constraint $c^2 + s^2 = 1$. The matrix J is then given $J_{pp} = J_{qq} = c$, $J_{pq} = -J_{qp} = s$,
 61 and $J_{ij} = \delta_{ij}$ for all the other entries. The new attributes

$$62 (2.1) \quad d_1 = M_{pp}c^2 + M_{qq}s^2 - 2M_{pq}cs,$$

$$63 (2.2) \quad d_2 = M_{pp}c^2 + M_{qq}s^2 + 2M_{pq}cs,$$

65 are referred to as the *sum and difference* variables or as the *scaling and detailing* variables.

A numerical stable way of computing J is as follows. Assuming $M_{pq} \neq 0$, compute

$$b = \frac{M_{pp} - M_{qq}}{2M_{pq}}$$

and define

$$t = \frac{\text{sgn}(b)}{|b| + \sqrt{b^2 + 1}}.$$

66 Then $c = \frac{1}{\sqrt{t^2 + 1}}$ and $s = ct$. A Jacobi rotation over a $n \times n$ matrix uses $O(1)$ space with $O(n)$
 67 work.

68 **2.1. Treelets.** The treelets algorithm [1, 2] constructs a multiresolution basis and a corre-
69 sponding hierarchical clustering over the attributes of some datasets in \mathbb{R}^d , with which sparsity
70 sparsity can be exploited. In its most efficient implementation [2] it is an $O(\min(nd^2, n^2d) +$
71 $Ld)$ algorithm (L is the number of levels in the multiresolution).

72 The algorithm starts by constructing a $d \times d$ empirical covariance matrix A_0 and an
73 attribute similarity matrix M_0 given by

$$74 \quad (2.3) \quad (M_0)_{ij} = \sqrt{\frac{(A_0)_{ij}^2}{(A_0)_{ii}(A_0)_{jj}}} + \lambda |(A_0)_{ij}|.$$

75 where λ is a regularization, hyper-parameter.

76 The initial set S_0 of scaling indices is that of all the variables, i.e. $S_0 = \{1, 2, \dots, d\}$. Starting
77 with A_0 and S_0 , at each tree level $k = 1, \dots, L$ ($L < d$), A_k and S_k are constructed as follows:

78 1. Compute the $d \times d$ matrix M_k whose entries are given by

$$80 \quad (2.4) \quad (M_k)_{ij} = \sqrt{\frac{(A_{k-1})_{ij}^2}{(A_{k-1})_{ii}(A_{k-1})_{jj}}} + \lambda |(A_{k-1})_{ij}|.$$

81 2. Find the two indices α_k, β_k such that

$$82 \quad (2.5) \quad \alpha_k, \beta_k = \operatorname{argmax}_{\alpha, \beta \in S_{k-1}} (M_k)_{\alpha\beta}.$$

83 3. Calculate the Jacobi matrix J_k for α_k, β_k entry of A_{k-1} and set $A_k = J_k^T A_{k-1} J_k$.

84 4. Set aside the difference variable. Assuming, without loss of generality, that $(A_k)_{\alpha_k\alpha_k} \leq$
85 $(A_k)_{\beta_k\beta_k}$, set $S_k = S_{k-1} - \{\alpha_k\}$.

86 The Jacobi rotations in the treelet algorithm produce an orthogonal basis to represent the
87 data for each $k \in \{1, 2, 3, \dots, L\}$ ($L = d - 1$ being the maximum level of the tree). Defining

$$88 \quad (2.6) \quad B_k = J_k^T J_{k-1}^T \cdots J_2^T J_1^T,$$

89 then

$$90 \quad (2.7) \quad A_k = B_k A_0 B_k^T.$$

91 Consequently, every vector $v \in \mathbb{R}^d$ has a k -th basis representation $B_k v$. Furthermore, there
92 is a compressed k th basis representation obtained by dropping insignificant ($< \epsilon$) detailing
93 (non-scaling) indices of $B_k v$. That is, if we define e_i to be the i th column of the identity
94 matrix, the compressed k th basis representation is given by

$$95 \quad \tau_k(v) = B_k v - \sum_{\substack{i \notin S_k \\ |B_k v \cdot e_i| < \epsilon}} (B_k v \cdot e_i) e_i.$$

96 Treelets can also be viewed as a hierarchical clustering method *over attributes*. The
97 hierarchical clustering structure is stored in α_k, β_k . We start with trivial clustering where
98 each attribute is in its own cluster and labeled by itself. For each k , we merge clusters labeled
99 α_k and β_k and label it as β_k . This is feasible because each step k the set of all cluster labels
100 is exactly S_{k-1} . This operation gives a hierarchical tree for attribute clustering.

105 **2.2. Kernel Method.** The kernel method [12] allow us to map variables into a new feature
 106 (Hilbert) space via a kernel function. We now review briefly the basic concepts and ideas of
 107 this approach (see for example [10, 11] for a more comprehensive review).

108 A kernel K over some set X is a function $K : X \times X \rightarrow \mathbb{R}$. A symmetric and positive
 109 semi-definite (SPSD) kernel K has the properties:

$$110 \quad (2.8) \quad K(x_1, x_2) = K(x_2, x_1), \text{ for all } x_1, x_2 \in X.$$

$$111 \quad (2.9) \quad \sum_{i=1}^s \sum_{j=1}^s c_i c_j K(x_i, x_j) \geq 0, \text{ for all } \{x_1, \dots, x_s\} \in X \text{ and all } \{c_1, \dots, c_s\} \in \mathbb{R}$$

112
 113 If X is finite, then K is SPSPD if and only if $K(X, X)$ is a SPSPD matrix. If $X \subseteq \mathbb{R}^d$, there is a
 114 unique Hilbert space \mathbb{H} and a feature map $\Phi_K : \mathbb{R}^d \rightarrow \mathbb{H}$ associated to a SPSPD kernel K such
 115 that for all $x, y \in X$,

$$116 \quad (2.10) \quad K(x, y) = \langle \Phi_K(x), \Phi_K(y) \rangle_{\mathbb{H}},$$

117 where $\langle \cdot, \cdot \rangle_{\mathbb{H}}$ stands for the inner product in \mathbb{H} .

118 The space \mathbb{H} here is called a reproducing kernel Hilbert space (RKHS). The following are
 119 two common examples of SPSPD kernels:

120 1. Radial basis function (RBF) kernel

$$121 \quad (2.11) \quad K(x, y) = \exp\left\{-\frac{\|x - y\|^2}{2\sigma^2}\right\}.$$

122 2. Polynomial kernel

$$123 \quad (2.12) \quad K(x, y) = (\alpha \langle x, y \rangle + c_0)^r.$$

124 A kernel K for a set X can be restricted to a subset $Y \subseteq X$, and the SPSPD property is
 125 preserved under such restriction. If the task under consideration is clustering over a finite set,
 126 the selected kernel needs only be SPSPD on the (finite) set of all samples. Thus, we only need
 127 to check that the kernel matrix is SPSPD. If we need to extend the clustering outcome to other
 128 data, e.g. for clustering boosted classification, then X has to include the whole data space as
 129 a subset.

130 **3. The KT Model.** The objective KT is produce a hierarchical dataclustering for some
 131 set $\mathcal{D} = \{d_1, \dots, d_n\}$ given a SPSPD kernel $K : \mathcal{D} \times \mathcal{D} \rightarrow \mathbb{R}$ measuring the data similarity.
 132 Instead of using the $d \times d$ empirical covariance matrix A_0 in the initial step (2.3) of the treelet
 133 transform, we replace A_0 with the $n \times n$ kernel matrix and apply the rest of the steps of
 134 treelets algorithm. In more detail:

- 135 1. First calculate the $n \times n$ kernel matrix $(A_0)_{ij} = K(d_i, d_j)$. A_0 is a SPSPD matrix
 136 because K is SPSPD kernel, and each column (or row) corresponds to a data point in
 137 \mathcal{D} .
- 138 2. Apply the treelet algorithm with hyper-parameter λ and $L = n - 1$ using the A_0 of
 139 step 1. In our experiments, λ is set to 0 but it can also be tuned as in treelets.
- 140 3. The hierarchical clustering produced by the treelet transform can now be viewed as
 141 a clustering of columns (or rows) of A_0 and consequently as a clustering of the data
 142 points of the set \mathcal{D} .

146 **3.1. Illustration.** To illustrate how the KT method work we use the following 5 point
 147 two-dimensional dataset:

$$148 \quad (3.1) \quad \mathcal{D} = \left\{ \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 2 \\ 0 \end{bmatrix}, \begin{bmatrix} 2 \\ -1 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ 1 \end{bmatrix} \right\}.$$

149 The data points are displayed in Figure 1. If we choose the RBF kernel (2.11) with $\sigma = 0.5$,

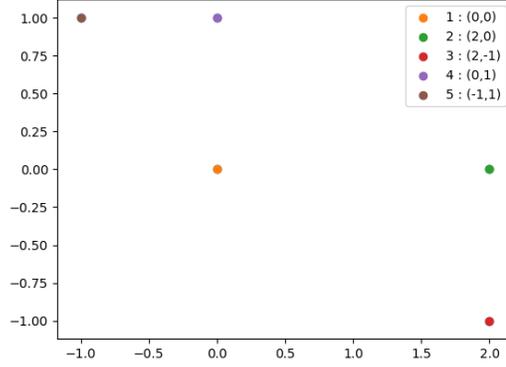


Figure 1. Distribution of dataset $\mathcal{D} = \{(0, 0), (2, 0), (2, -1), (0, 1), (-1, 1)\}$

then the kernel matrix becomes

$$A_0 = \begin{bmatrix} 1 & e^{-8} & e^{-10} & e^{-2} & e^{-4} \\ e^{-8} & 1 & e^{-2} & e^{-10} & e^{-20} \\ e^{-10} & e^{-2} & 1 & e^{-16} & e^{-24} \\ e^{-2} & e^{-10} & e^{-16} & 1 & e^{-2} \\ e^{-4} & e^{-20} & e^{-24} & e^{-2} & 1 \end{bmatrix} \approx \begin{bmatrix} 1.000 & 0.000 & 0.000 & 0.135 & 0.018 \\ 0.000 & 1.000 & 0.135 & 0.000 & 0.000 \\ 0.000 & 0.135 & 1.000 & 0.000 & 0.000 \\ 0.135 & 0.000 & 0.000 & 1.000 & 0.135 \\ 0.018 & 0.000 & 0.000 & 0.135 & 1.000 \end{bmatrix}.$$

150 Applying treelets to it gives a hierarchical tree for columns (or rows) of A_0 and as remarked
 151 above, for the data points themselves as Figure 2 illustrates.

152 **4. Theory.** We now prove that the kernel projection is equivalent to working with an
 153 $n \times n$ symmetric positive semi-definite matrix, regardless of the dimension of the KRHS \mathbb{H} ,
 154 and that this matrix can be efficiently evaluated through the kernel K . We also suggest a
 155 definition of a *clustering frame* and *clustering equivalence* for similarity-based clustering that
 156 allows us to connect the results of the clustering analysis for the original dataset with those
 157 of the transformed, projected set and thus explain the usefulness of KT approach.

158 Hereto \mathcal{D} denotes the dataset and D the corresponding matrix whose columns are the
 159 data points in \mathcal{D} and similarly for functions of \mathcal{D} .

160 **Lemma 4.1.** *Let $X \subseteq \mathbb{R}^d$. For every finite dataset $\mathcal{D} = \{d_1, \dots, d_n\} \subseteq X$ and an SPSD*
 161 *kernel K , there exists an orthonormal basis B of the RKHS \mathbb{H} such that*

$$162 \quad (4.1) \quad [\Phi_K(d_i)]_B = \begin{bmatrix} \delta_i \\ 0 \end{bmatrix} \quad \text{for } i = 1, \dots, n,$$

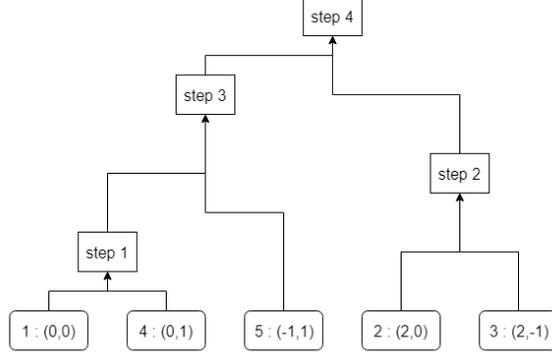


Figure 2. Clustering tree for $\mathcal{D} = \{(0,0), (2,0), (2,-1), (0,1), (-1,1)\}$.

164 where the left hand side stands for the representation of $\Phi_K(d_i)$ in the basis B and $\delta_i \in \mathbb{R}^n$.
 165 Moreover, the matrix $[\delta_1 \ \delta_2 \ \dots \ \delta_n]$ is symmetric, positive semi-definite.

166 *Proof.* We apply the Gram-Schmidt orthogonalization process to the maximal linearly in-
 167 dependent subset of $\{\Phi_K(d_1), \dots, \Phi_K(d_n)\}$ and get a set of orthonormal vectors $\{\hat{\beta}_1, \dots, \hat{\beta}_m\}$,
 168 where

169 (4.2)
$$m = \dim(\text{span}\{\Phi_K(d_1), \dots, \Phi_K(d_n)\}) \leq n.$$

171 We extend this set to a orthonormal basis $\hat{B} = \{\hat{\beta}_1, \dots, \hat{\beta}_m, \dots\}$ of \mathbb{H} . Then, for all $i \in$
 172 $\{1, 2, \dots, n\}$, $[\Phi_K(d_i)]_{\hat{B}}$ is 0 after the m -th entry, and consequently after n -th entry, so there
 173 exists $\hat{d}_i \in \mathbb{R}^n$ such that

174 (4.3)
$$[\Phi_K(d_i)]_{\hat{B}} = \begin{bmatrix} \hat{d}_i \\ 0 \end{bmatrix}.$$

 175

176 The $n \times n$ matrix $[\hat{d}_1 \ \hat{d}_2 \ \dots \ \hat{d}_n]$ can be written in its singular value decomposition

177 (4.4)
$$[\hat{d}_1 \ \hat{d}_2 \ \dots \ \hat{d}_n] = U\Sigma V^T,$$

 178

179 where U and V are orthogonal matrices and Σ is a diagonal matrix with non-negative entries.

180 We can now define a new orthonormal basis $B = \{\beta_1, \dots, \beta_m, \dots\}$ through the change of
 181 basis matrix $\begin{bmatrix} VU^T & 0 \\ 0 & I \end{bmatrix}$. Let $\delta_i = VU^T \hat{d}_i$ for all $i \in \{1, 2, \dots, n\}$, then

182 (4.5)
$$[\Phi_K(d_i)]_B = \begin{bmatrix} VU^T & 0 \\ 0 & I \end{bmatrix} [\Phi_K(d_i)]_{\hat{B}} = \begin{bmatrix} VU^T \hat{d}_i \\ 0 \end{bmatrix} = \begin{bmatrix} \delta_i \\ 0 \end{bmatrix}.$$

 183

184 The projected data $\Phi_K(d_i)$ in the basis B is $[\delta_i \ 0]^T$ and the matrix

185 (4.6)
$$[\delta_1 \ \delta_2 \ \dots \ \delta_n] = VU^T [\hat{d}_1 \ \hat{d}_2 \ \dots \ \hat{d}_n] = V\Sigma V^T$$

 186

187 is symmetric and positive semi-definite. ■

188 **Corollary 4.2.** Let $X \subseteq \mathbb{R}^d$ and $\mathcal{D} = \{d_1, \dots, d_n\} \subseteq X$ a dataset. For every $x \in X$, define
 189 $\Psi(x)$ as the first n components of $\Phi_K(x)$ in the basis B of Lemma 4.1, i.e.

$$190 \quad (4.7) \quad \begin{bmatrix} \Psi(x) \\ * \end{bmatrix} = [\Phi_K(x)]_B.$$

192 Then,

$$193 \quad (4.8) \quad \langle \Psi(D), \Psi(D) \rangle = \langle \Phi_K(D), \Phi_K(D) \rangle_{\mathbb{H}} = K(D, D).$$

195 Here, $\langle \Psi(D), \Psi(D) \rangle$ denotes the matrix with entries $\langle \Psi(d_i), \Psi(d_j) \rangle$ and similarly for
 196 $\langle \Phi_K(D), \Phi_K(D) \rangle_{\mathbb{H}}$ and $K(D, D)$.

Proof. For all $d_i \in \mathbb{H}$,

$$\begin{bmatrix} \Psi(d_i) \\ 0 \end{bmatrix} = [\Phi_K(d_i)]_B,$$

that is $\Psi(d_i) = \delta_i$. From Lemma 4.1, we have that $\Psi(D) = [\delta_1 \ \delta_2 \ \dots \ \delta_n]$ is symmetric, positive semi-definite and

$$\langle \Psi(D), \Psi(D) \rangle = [\delta_1 \ \delta_2 \ \dots \ \delta_n]^2 = \langle \Phi_K(D), \Phi_K(D) \rangle_{\mathbb{H}} = K(D, D).$$

197 **4.1. Clustering Equivalences.**

198 **Definition 4.3.** A clustering frame is a pair (\mathcal{D}, f) where \mathcal{D} is a finite, ordered dataset and
 199 $f : \mathcal{D} \times \mathcal{D} \rightarrow \mathbb{R}$ is a mapping that measures similarity for the points in the dataset \mathcal{D} .

200 **Definition 4.4.** Two clustering frames (\mathcal{D}_1, f_1) and (\mathcal{D}_2, f_2) are equivalent, and we write
 201 $(\mathcal{D}_1, f_1) = (\mathcal{D}_2, f_2)$, if and only if $f_1(\mathcal{D}_1, \mathcal{D}_1) = f_2(\mathcal{D}_2, \mathcal{D}_2)$.

202 For any similarity-based clustering method, using two equivalent clustering frames gives the
 203 same clustering outcome.

204 A pertinent example of clustering frame equivalence is the following. If an SPSD kernel
 205 K corresponds to the projection Φ_K , then there is a RKHS \mathbb{H} such that

$$206 \quad (4.9) \quad K(D, D) = \langle \Phi_K(D), \Phi_K(D) \rangle_{\mathbb{H}}.$$

208 Therefore, the clustering on \mathcal{D} using similarity mapping K corresponds to the clustering on
 209 $\Phi_K(\mathcal{D})$ using the inner product in \mathbb{H} , i.e.

$$210 \quad (4.10) \quad (\mathcal{D}, K) = (\Phi_K(\mathcal{D}), \langle \cdot, \cdot \rangle_{\mathbb{H}}).$$

4.2. Kernel Treelets Clustering Equivalences. Recall that the treelet transform applied to some dataset \hat{D} produces a hierarchical clustering *over its attributes*, or in other words, a clustering *over the data* of \hat{D}^T . If we apply the treelet transform to $\hat{D} = \Psi^T(D)$, then the outcome would be a hierarchical clustering over the columns of $\hat{D}^T = (\Psi^T(D))^T = \Psi(D)$. Therefore, the clustering frame for the KT is $(\Psi(D), \langle \cdot, \cdot \rangle)$ and because

$$\langle \Phi_K(D), \Phi_K(D) \rangle_{\mathbb{H}} = \langle \Psi(D), \Psi(D) \rangle,$$

212 we have $(\Psi(\mathcal{D}), \langle \cdot, \cdot \rangle) = (\Phi_K(\mathcal{D}), \langle \cdot, \cdot \rangle_{\mathbb{H}})$. Finally, from (4.10) we get the clustering frame
 213 equivalence

$$214 \quad (4.11) \quad (\mathcal{D}, K) = (\Phi_K(\mathcal{D}), \langle \cdot, \cdot \rangle_{\mathbb{H}}) = (\Psi(\mathcal{D}), \langle \cdot, \cdot \rangle),$$

216 which implies that applying the treelet transform to the matrix $\hat{D} = \Psi^T(D)$ provides the
 217 same hierarchical clustering as that produced by the clustering frame (\mathcal{D}, K) . Also, rather
 218 than computing A_0 directly through \hat{D} , one computes this matrix efficiently with the *kernel*
 219 *trick* (4.9) : $A_0 = K(D, D)$.

220 **4.3. Complexity.** The complexity of computing kernel matrix is $O(\xi n^2)$, where ξ is the
 221 complexity of applying kernel function to a pair of data and $\xi = d$ if the data is numeric.
 222 Computing the rotation steps can be seen as applying treelets to a $d = n$ matrix, and thus it
 223 complexity is $O(Ld) = O((n-1)n) = O(n^2)$ if properly optimized as in treelets. So the total
 224 complexity of KT is $O(\xi n^2)$.

225 **5. Examples.** We implemented KT and the following examples in Python with the pack-
 226 ages Numpy [15] and Scikit-learn [16]. Plots were generated with Matplotlib [17]. The treelets
 227 part of our implementation is not optimized, so its cost is $O(n^3)$ operations (an $O(n^2)$ imple-
 228 mentations is also possible [1]). The hyperparameter λ is set to 0 in all of the experiments
 229 below.

230 **5.1. Clustering for Six Datasets.** To illustrate how KT works as a hierarchical clustering
 231 method over data, we use first an example from Scikit-learn [16] which consists of 6 datasets,
 232 each of which has 1500 two-dimensional data points (i.e. $n = 1500$ and $d = 2$). We can
 233 visualize each dataset and each cluster by plotting each observation as a point in the plane.
 234 Each of the first five datasets consists of data drawn from multiple shapes with an error in
 235 distance. The sixth dataset is a uniform random sample from $[0, 1]^2$ to show how clustering
 236 methods work for uniform distributed data and specially how smooth the boundaries of their
 237 partitions are.

238 **Figure 3** compares the performance of KT with different kernels with that of some other
 239 clustering methods for the six aforementioned datasets. The number of clusters and hyper-
 240 parameters are tuned for each method and the sample sizes are set to 1000 for each instance
 241 of the KT method. Each row in **Figure 3** represents a dataset and each column represents
 242 a clustering method. In this experiment, KT with RBF kernel is the method that performs
 243 clustering closest to human intuition for all first five datasets. Only spectral clustering (column
 244 5) has a similar performance. The sixth dataset shows that KT is affected by the relative
 245 density deficiency in some area due to sampling and shows porous boundaries. The excellent
 246 performance of the RBF KT on the first five datasets can be traced to the fact that these
 247 datasets are to some extent Euclidean distance-based, which corresponds to the assumptions
 248 for RBF kernel.

249 **5.2. Clustering for a Social Network Dataset.** We now consider an example of network
 250 analysis from the Stanford Network Analysis Project [18]. This is a dataset consisting of
 251 ‘circles of friends’ (or ‘friends lists’) from Facebook. It has $n_V = 4039$ surveyed individuals,
 252 which can be viewed as vertices of an undirected graph. Each two vertices are connected with

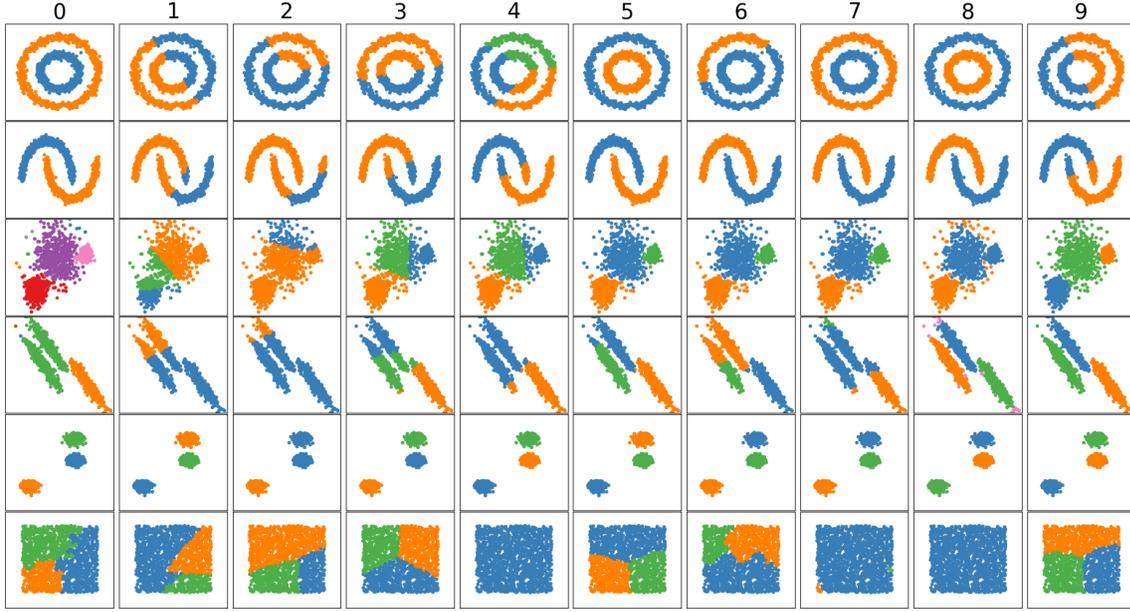


Figure 3. Comparison of different clustering algorithms on 6 datasets. 0 - KT (RBF), 1 - KT ($r = 1$), 2 - KT ($r = 2$), 3 - Mini batch Kmeans, 4 - Mean shift, 5 - Spectral clustering, 6 - Ward, 7 - Agglomerative clustering, 8 - DBSCAN, 9 - Gaussian mixture.

253 an edge if the corresponding individuals are friends and not otherwise. The edges are not
 254 weighted and the total number is 88234. We use KT to obtain a hierarchical clustering on
 255 this dataset.

256 Denote the set of vertices on the graph as V and define a kernel function $K : V \times V \rightarrow \mathbb{R}$
 257 such that for every $v_1, v_2 \in V$

$$(5.1) \quad K(v_1, v_2) = \begin{cases} 1045 & \text{if } v_1 = v_2. \\ 1 & \text{if } v_1, v_2 \text{ are connected.} \\ 0 & \text{otherwise.} \end{cases}$$

260 The number 1045 is the largest degree of all vertices.

261 **Lemma 5.1.** *The matrix $K(V, V)$ defined by (5.1) is SPSD.*

262 *Proof.* Clearly, $K(V, V)$ symmetric. Moreover it is diagonally dominant with a positive
 263 diagonal:

$$264 \quad \sum_{j \neq i} |K(V, V)_{i,j}| = \text{degree}(v_i) \leq \max_j \text{degree}(v_j) = 1045 = K(V, V)_{i,i}, \quad i = 1, 2, \dots, n_V.$$

266 Now, since $K(V, V)$ is symmetric all its eigenvalues are real. Suppose there is a negative
 267 eigenvalue $-\lambda^2$. Then $K(V, V) + \lambda^2 I$ is singular but this is impossible because $K(V, V) + \lambda^2 I$
 268 is strictly diagonally dominant. Therefore $K(V, V)$ is SPSD. ■

269 To estimate the performance of KT as a multiresolution, hierarchical clustering method
 270 on this dataset, we do the following evaluation. For each cluster partition in the hierarchy, i.e.
 271 for each tree level $k = 1, \dots, L = n - 1$, we compute its confusion matrix and the corresponding
 272 true positive rate and false positive rate. The confusion matrix, is a $2 \times x$ array recording
 273 the number of true positives (true predicted connections), true negatives (true predicted no-
 274 connections), false positives (false predicted connections), and false negatives (false predicted
 275 no-connections) for pairwise associations. The true positive rate (TPR) measures the propor-
 276 tion of two nodes being in the same cluster given that the two nodes are connected. The false
 277 positive rate (FPR) measures the proportion of two nodes being in the same cluster given that
 278 the two nodes are not connected. For each tree level, $k = 1, \dots, n - 1$, there corresponds a point
 279 in the plane $(\text{TPR}(k), \text{FPR}(k))$ and interpolating these points we obtain the so-called receiver
 280 operating characteristic (ROC) curve, which is displayed in [Figure 4](#). The performance of the
 281 KT for clustering on this dataset is excellent. As a reference, the line $y = x$ corresponds to a
 282 prediction accuracy of random guessing. For the KT clustering, at about 20% FPR we obtain
 283 almost 100% TPR, i.e. if we take 20% of nodes that are not connected and place them in
 284 the clusters obtained by the KT method then they will be connected to other vertices with
 285 almost probability 1. Another measure of the effectiveness of the hierarchical clustering is the
 286 so-called area under the curve, which is the numerical integral of the ROC over $[0, 1]$. For the
 287 KT the AUC is 0.958, very close to the optimal value 1.

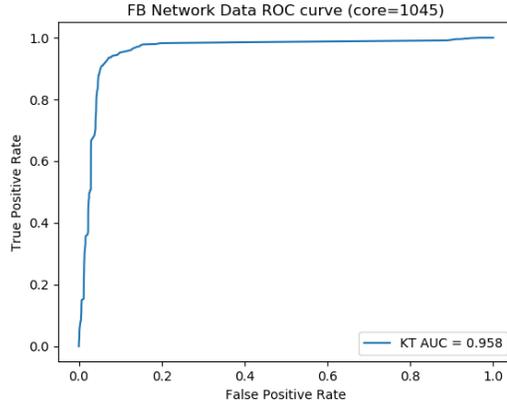


Figure 4. Receiver operating curve (ROC), as the tree level increases, for the KT clustering for the social network example.

288 **5.3. Clustering for a Dataset with Missing Information.** Our last example is a dataset
 289 with missing information. We use the mice protein expression (MPE) dataset [19] from the
 290 UCI Machine Learning Repository. This is a dataset consisting of 1080 observations for 8
 291 classes of mice, each of which containing 77 expression levels of different proteins with some
 292 of the entries missing.

293 We employ KT to obtain a hierarchical clustering on this dataset. First, we standardize
 294 the data so that each attribute has empirical mean 0 and standard deviation 1. Then, we

295 define the following RBF kernel. For for all observations u, v in the dataset V ,

$$296 \quad (5.2) \quad K(u, v) = \exp \left\{ - \frac{32}{|E_{uv}|} \sum_{i \in E_{uv}} \|u_i - v_i\|^2 \right\},$$

297

298 where E_{uv} is the set of indices for which the data is available in both u and v . We check that
 299 $E_{uv} \neq \emptyset$ so that K is well-defined. The number 32 is a parameter empirically selected. We
 300 confirmed numerically that the kernel matrix is SPSP.

301 We compare the predicted clusters and the true labels according to pairwise scores. **Fig-**
 302 **ure 5** shows how KT performs compared to the popular KMeans clustering method. Here,
 303 we measure the true positive rate as the proportion of two records being in the same cluster
 304 given that they are from mice of the same type, and the false positive rate as the proportion
 305 of two records being in the same cluster given that they are from mice of different type. As
 306 in the previous example of the network dataset, we draw the ROC and calculate its AUC.
 307 The AUC of KT is much greater than that of KMeans ($0.726 > 0.579$), demonstrating KT's
 308 superiority for this dataset.

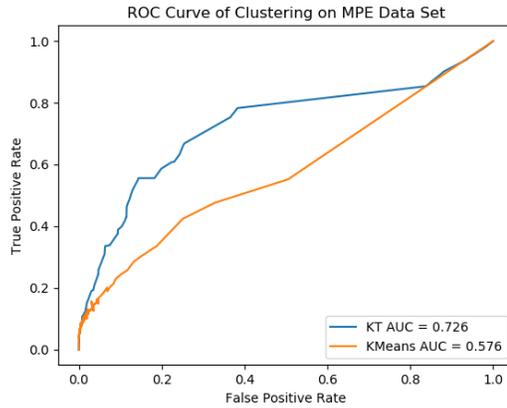


Figure 5. Comparison of receiver operating curve (ROC) for KT and KMeans on the MPE dataset.

309 **6. Accelerating Kernel Treelets.** The treelet approach is meant for small to moderate
 310 size high-dimensional datasets (small n , large d) because of its $O(n^2)$ complexity. In the
 311 context of spectral clustering, Yan, Huang and Jordan [20] proposed a preprocessing method,
 312 using classical k -means or random projections trees, to reduce the size of the input dataset
 313 and thus accelerate the spectral analysis of the method.

314 In a similar spirit, we propose here to apply KT to a moderate size sample $\mathcal{S} \subset \mathcal{D}$ for an
 315 initial clustering and then use a kernel support vector machine (SVM) to assign cluster labels to
 316 the rest of the data, i.e. the all $d_j \in \mathcal{D} \setminus \mathcal{S}$. To illustrate this approach, we consider again the six
 317 datasets of Example 1. **Figure 6** demonstrates how the number of sample points n_S affects the
 318 clustering result. Each column represents KT using the RBF kernel for different n_S , denoted
 319 KTn_s in Table 1. The run time is displayed in Table 1, columns 2-7. The hyper-parameter
 320 $\sigma = 0.1$ is tuned towards $n_S = 1000$ case and is used for all other sample sizes. Note that as

321 KT1500 is of full sample size, it does not trigger the kernel SVM whereas KT1499 does. From
 322 **Figure 6** and **Table 1**, we observe that the minimum optimal n_s for the first 5 datasets is 1000,
 323 100, 1000, 200, 50, respectively (n_s is dataset dependent, as expected) and thus the overall
 324 cost of the clustering analysis could be substantially reduced with this approach. Furthermore,
 325 the fourth dataset shows that optimal hyper-parameter σ is n_s -dependent. The RBF kernel
 326 can be considered as a weighted average of distance and connectivity, where a larger σ means
 327 a higher weight on distance. For the same $\sigma = 0.1$, as sample size n_s increases, the clustering
 328 result becomes more distance-based rather than connectivity-based, demonstrating that the
 329 optimal σ for those sample sizes is actually smaller.

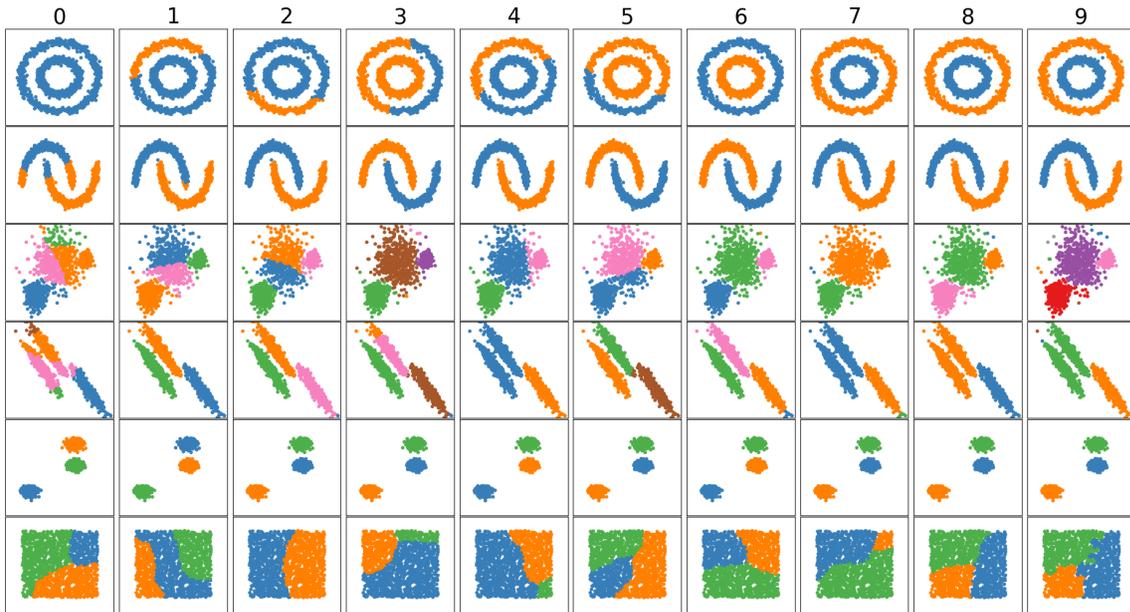


Figure 6. Comparison of different number-of-cluster estimate on 6 datasets.

330 **7. Concluding remarks.** In the paper we describe a novel approach, kernel treelets (KT),
 331 for hierarchical clustering. The method relies on applying the treelet transform to an $n \times n$
 332 matrix measuring data similarities in a feature, reproducing kernel Hilbert space. We show
 333 with some examples that KT can be as useful as other hierarchical clustering methods and
 334 is especially competitive for datasets without numerical matrix representation and or then
 335 there is missing data. The KT approach also shows significant potential for semi-supervised
 336 learning tasks and as a pre-processing, post-processing step in deep-learning. Work in these
 337 directions is underway.

338 **8. Acknowledgements.** HDC gratefully acknowledges support from the National Science
 339 Foundation through grant DMS 1818821.

REFERENCES

KTn_s	1	2	3	4	5	6
0 - KT50	0.011	0.012	0.013	0.011	0.011	0.01
1 - KT100	0.035	0.044	0.039	0.045	0.033	0.028
2 - KT200	0.109	0.099	0.128	0.12	0.121	0.132
3 - KT300	0.225	0.217	0.242	0.269	0.259	0.235
4 - KT500	0.551	0.568	0.62	0.569	0.652	0.536
5 - KT800	1.315	1.513	1.534	1.378	1.699	1.295
6 - KT1000	2.016	2.055	2.336	2.098	2.782	1.941
7 - KT1200	2.88	2.94	3.242	3.004	4.146	2.77
8 - KT1499	4.438	4.532	5.4	4.713	6.788	4.341
9 - KT1500	4.472	4.69	5.398	4.807	6.782	4.274

Table 1

KTn_s denotes the use of KT to a sample of size n_s . The numbers on columns 2-7 are the run times for the clusters in Figure 7.

- 341 [1] Ann B Lee and Boaz Nadler. Treelets— a tool for dimensionality reduction and multi-scale analysis of
342 unstructured data. In *Artificial Intelligence and Statistics*, pages 259–266, 2007.
- 343 [2] Ann B. Lee, Boaz Nadler, and Larry Wasserman. Treelets: an adaptive multi-scale basis for sparse
344 unordered data. *The Annals of Applied Statistics*, 2(2):435–471, 2008.
- 345 [3] Robert C. Tryon. *Cluster analysis: Correlation profile and orthometric (factor) analysis for the isolation*
346 *of unities in mind and personality*. Edwards brother, Incorporated, lithoprinters and publishers, 1939.
- 347 [4] Stephen C Johnson. Hierarchical clustering schemes. *Psychometrika*, 32(3):241–254, 1967.
- 348 [5] Robin Sibson. Slink: an optimally efficient algorithm for the single-link cluster method. *The computer*
349 *journal*, 16(1):30–34, 1973.
- 350 [6] Daniel Defays. An efficient algorithm for a complete link method. *The Computer Journal*, 20(4):364–366,
351 1977.
- 352 [7] David S Doermann. An introduction to vectorization and segmentation. In *International Workshop on*
353 *Graphics Recognition*, pages 1–8. Springer, 1997.
- 354 [8] Yongjin Wang, Foteini Agrafioti, Dimitrios Hatzinakos, and Konstantinos N Plataniotis. Analysis of
355 human electrocardiogram for biometric recognition. *EURASIP journal on Advances in Signal Pro-*
356 *cessing*, 2008(1):148658, 2007.
- 357 [9] Fiona M Shrive, Heather Stuart, Hude Quan, and William A Ghali. Dealing with missing data in a multi-
358 question depression scale: a comparison of imputation methods. *BMC medical research methodology*,
359 6(1):57, 2006.
- 360 [10] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning. Data Mining, Inference,*
361 *and Prediction*. Springer, New York, second edition, 2016.
- 362 [11] S. Theodoridis. *Machine Learning. A Bayesian and Optimization Perspective*. Academic Press, London,
363 2015.
- 364 [12] Mark A Aizerman. Theoretical foundations of the potential function method in pattern recognition
365 learning. *Automation and remote control*, 25:821–837, 1964.
- 366 [13] Naomi S Altman. An introduction to kernel and nearest-neighbor nonparametric regression. *The American*
367 *Statistician*, 46(3):175–185, 1992.
- 368 [14] Bernhard E Boser, Isabelle M Guyon, and Vladimir N Vapnik. A training algorithm for optimal margin
369 classifiers. In *Proceedings of the fifth annual workshop on Computational learning theory*, pages 144–
370 152. ACM, 1992.
- 371 [15] Travis E Oliphant. *A guide to NumPy*, volume 1. Trelgol Publishing USA, 2006.
- 372 [16] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-
373 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and
374 E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*,

- 375 12:2825–2830, 2011.
- 376 [17] John D. Hunter. Matplotlib: A 2d graphics environment. *Computing In Science & Engineering*, 9(3):90–
377 95, 2007.
- 378 [18] Jure Leskovec and Julian J Mcauley. Learning to discover social circles in ego networks. In *Advances in*
379 *neural information processing systems*, pages 539–547, 2012.
- 380 [19] Clara Higuera, Katherine J Gardiner, and Krzysztof J Cios. Self-organizing feature maps identify proteins
381 critical to learning in a mouse model of Down syndrome. *PloS one*, 10(6):e0129126, 2015.
- 382 [20] Donghui Yan, Ling Huang, and Michael I. Jordan. Fast Approximate Spectral Clustering *Proceedings of*
383 *the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 907–
384 916, 2009.

385

REFERENCES