



LAWRENCE
LIVERMORE
NATIONAL
LABORATORY

Robust Local Scaling using Conditional Quantiles of Graph Similarities

J. J. Thiagarajan, P. Sattigeri, K. N. Ramamurthy,
B. Kailkhura

September 22, 2016

International Conference on Data Mining
Barcelona, Spain
December 12, 2016 through December 15, 2016

Disclaimer

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.

Robust Local Scaling using Conditional Quantiles of Graph Similarities

Jayaraman J. Thiagarajan
Lawrence Livermore National Labs
7000 E Avenue Livermore, CA 94550
Email: jjayaram@llnl.gov

Prasanna Sattigeri, Karthikeyan
Natesan Ramamurthy
IBM Research, Yorktown Heights, NY 10568
Email: {psattig, knatesa}@us.ibm.com

Bhavya Kailkhura
Syracuse University
Syracuse, NY 13244
Email: bkailkhu@syr.edu

Abstract—Spectral analysis of neighborhood graphs is one of the most widely used techniques for exploratory data analysis, with applications ranging from machine learning to social sciences. In such applications, it is typical to first encode relationships between the data samples using an appropriate similarity function. Popular neighborhood construction techniques such as k -nearest neighbor (k -NN) graphs are known to be very sensitive to the choice of parameters, and more importantly susceptible to noise and varying densities. In this paper, we propose the use of quantile analysis to obtain local scale estimates for neighborhood graph construction. To this end, we build an auto-encoding neural network approach for inferring conditional quantiles of a similarity function, which are subsequently used to obtain robust estimates of the local scales. In addition to being highly resilient to noise or outlying data, the proposed approach does not require extensive parameter tuning unlike several existing methods. Using applications in spectral clustering and single-example label propagation, we show that the proposed neighborhood graphs outperform existing locally scaled graph construction approaches.

1. Introduction

Neighborhood graphs are central to techniques that involve analysis and exploration of high-dimensional data. Constructing a graph involves encoding the relationships between the data samples using an appropriate similarity function. Spectral analysis of such graphs is the modus operandi in a variety of applications, including dimensionality reduction, image segmentation, text mining, and data analysis in general. These methods often involve the eigen-decomposition of the similarity (also referred to as the adjacency) matrix that reveal strong connections to graph properties such as connected components, the diameter of a graph and the degree of randomness.

Defining the notion of an appropriate neighborhood and adapting the analysis to the local scale or density of the data have been long-standing research problems. k -nearest neighborhood or ϵ -neighborhood graphs are the most commonly adopted approaches in practice. However, the instabilities arising due to noise or outlying data have plagued the performance of nearest neighbor graphs. Furthermore,

clusters with varying densities commonly occur in high-dimensions, which make the global neighborhood parameter choices highly unreliable. Consequently, a broad class of techniques that attempt to estimate the local scale to improve spectral analysis of data with varying densities, shapes and noise have been developed [1]. Another class of approaches translate the density variations into edge probabilities representing their significance in recovering the underlying structure [2]. Several of these methods still rely heavily on heuristics and parameter tuning to perform consistently across different domains.

In this paper, we explore the use of quantile analysis to obtain local scale estimates for building robust similarities. The proposed approach falls under the class of methods that construct stochastic graphs, which is carried out based on a novel, unsupervised quantile analysis framework. Quantile regression has been primarily used in analysis of datasets with heterogeneous properties, wherein traditional loss functions such as the ℓ_2 fail to account for biases in different parts of the data. In our context, spectral decomposition of graphs using squared ℓ_2 loss is widely used with a stable numerical solution (eigen decomposition). However, when we use large graphs and wish to discover its non-linear spectral structure at various quantiles, it becomes imperative to use a general optimization procedure. To this end, we build an auto-encoding neural network, which imposes a quantile loss between the input and reconstructed similarity matrices, to infer the conditional quantiles of similarity functions on graphs. The conditional quantiles are subsequently used to obtain robust estimates of the local scale. We show that the resulting neighborhood graphs outperform existing locally scaled graph construction approaches. Furthermore, we demonstrate through our experiments that the proposed method is highly resilient to noise and does not require extensive parameter tuning. Our contributions can be summarized as follows:

- We generalize the notion of quantile analysis to the similarity function defined on unsupervised graphs.
- We build a neural network architecture for efficient inference of conditional quantiles of neighborhood similarities.
- We relate the rate of edge decay across quantiles to the edge probabilities while constructing stochastic

neighborhood graphs.

- Using the proposed stochastic graphs, we develop a robust local scale estimation algorithm.
- We demonstrate that the proposed approach is resilient to noise and variations in local densities, and consistently outperforms existing approaches for spectral clustering.
- We propose a greedy algorithm for using the inferred graph similarities in label propagation with limited training examples.
- In an extreme setting, we evaluate the robustness of the graphs in label propagation with a single example (per class) and demonstrate substantial improvements.

2. Related Work

Constructing neighborhood graphs and performing spectral analysis of pairwise affinities are common to a wide-range of applications dealing with complex, high-dimensional data. Simple neighborhood techniques such as k-nearest neighbor (k-NN) graphs are known to be very sensitive to the choice of parameters. Alternatively, one can simply connect all points in a fully connected graph and rely on a scaling parameter σ to define the affinity between two points. In either case, the technique relies on setting a global parameter that does not take into account the variations in local densities. Consequently, several improvements have been proposed in the literature, through analysis of the underlying graph structure [1], [3], [4], and the stability [5], [6] of spectral clustering. In particular, techniques that obtain estimates of the local scale to handle data with varying densities and levels of noise have gained significant interest. One of the earliest strategies to estimate the local scales was developed by Zelnik-Manor and Perona [1]. Though this approach is effective even in high dimensions, it often leads to sub-optimal performances in presence of outliers/noise and in data with clusters of different densities. To alleviate this, Nadler *et al.* employed a coherence measure of a set of points for belonging to the same cluster [7]. Furthermore, Li *et al.* [8] proposed a warping model that maps the data into a new feature space for reliable clustering. Another interesting approach for local scale estimation stems from the use of proximity graphs (e.g. beta skeletons) to infer the neighborhood parameters [9], which was found to be resilient to noise.

Since pairwise similarities are solely based on the Euclidean distances between samples in the input space, they reveal no information about the inherent class structure. An effective approach to capture that information is to exploit the underlying manifold structure so that samples belonging to the same manifold have consistently higher similarity while samples belonging to different manifolds do not. The path-based graphs [10] define a similarity measure that implicitly infers the underlying structure and produces a robust neighborhood graph. Another important class of approaches for graph construction attempt to build probabilistic graphs that reveal the relative significance of the different edges

[11], [12]. In particular, the consensus clustering algorithm in [2] circumvents the problem of parameter selection by creating an ensemble of clustering with different parameter choices and exploiting the theory of nearly uncoupled Markov chains to construct a probabilistic graph, which is finally used with spectral clustering for robust analysis. This iterative procedure can be practically infeasible in high-dimensions, and their performance can be affected by varying densities. In this paper, we propose to adopt ideas from quantile analysis to construct robust graph similarities and thereby alleviate the inherent challenges with local scale estimation algorithms.

3. Inferring Conditional Quantiles of Graph Similarity

Supervised regression is a common statistical approach employed to analyze the relationships between the predictor variables and a response variable. Regression with squared ℓ_2 loss, *aka* the method of least squares, estimates the conditional mean of the response variable for the given predictors. This is sufficient when the data is homogeneous; however, when the data is heterogeneous, merely estimating the conditional mean is insufficient, as estimates of the standard errors are often biased. To comprehensively analyze such heterogeneous datasets, quantile regression is a better alternative. Quantile regression aims at estimating either the conditional median or other quantiles of the response variable [13], [14].

3.1. Definition: Quantile Loss

Quantile losses as fidelity measures for function approximation has applications in computational biology [15], survival analysis [16], workforce analytics [17], economics [18] and data analysis [19], [20] to name a few. For a scalar residual r , the quantile loss is a *check function* defined as

$$q_\tau(r) = (-\tau + 1[r \geq 0])r. \quad (1)$$

The quantile loss is piecewise linear and shares the robustness properties of the ℓ_1 loss by not penalizing the outliers as harshly as the squared ℓ_2 loss. Note that in (1) the quantile loss is equivalent to the ℓ_1 loss when $\tau = 0.5$. Consequently, similar to ℓ_1 , the quantile loss is non-differentiable at the origin and forces the residuals close to the origin to be exactly zero which may not be preferred in some applications.

A smoothed ‘huberized’ version of the quantile loss has been recently proposed [21]. This is defined as

$$\rho_\tau(r) = \begin{cases} \tau|r| - \frac{\kappa\tau^2}{2} & \text{if } r < -\tau\kappa, \\ \frac{1}{2\kappa}r^2 & \text{if } r \in [-\kappa\tau, (1-\tau)\kappa], \\ (1-\tau)|r| - \frac{\kappa(1-\tau)^2}{2}, & \text{if } r > (1-\tau)\kappa. \end{cases} \quad (2)$$

The quantile Huber loss permits small residuals by behaving like a squared ℓ_2 loss near the origin and hence may be preferred in regression settings. Figure 1 shows the quantile ($\tau = 0.3$) and quantile Huber ($\tau = 0.3$) loss functions. Furthermore, the quantile loss itself is a special case of

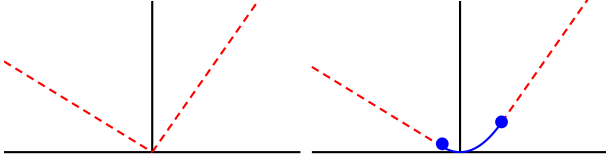


Figure 1. Quantile ($\tau = 0.3$) loss (left) and quantile Huber ($\tau = 0.3$) loss (right). The quantile Huber loss is obtained by smoothing the quantile loss at the origin.

quantile Huber as $\kappa \rightarrow 0$. Both quantile and quantile Huber losses are additive along the co-ordinates of the residual and hence the loss for multi-dimensional residuals is easily defined as $\rho_\tau(\mathbf{r}) = \sum_{i=1}^N \rho_\tau(r_i)$. In the rest of the paper, we will use the Huberized version of quantile loss.

3.2. Quantile Huber Loss in Spectral Graph Decomposition

Although quantile losses are traditionally used only in regression settings, it can be beneficial to employ them in unsupervised learning where the goal is to explore the structure of the data in lieu of fitting a function. Alternately, we could consider the response variable to be as the input data itself and thereby infer the underlying structure using the loss. However, the meaning of this formulation is problem-specific and needs careful consideration. We will focus on the problem of spectral decomposition of graphs in this paper.

Let us define a undirected graph with N nodes denoted by its similarity matrix $\mathbf{W} \in \mathbb{R}^{N \times N}$. The (i, j) th element w_{ij} corresponds to the similarity of the edge between the two nodes. \mathbf{W} is symmetric positive semi-definite, has all positive entries, and the maximum value of the entries is 1. The low-rank decomposition of this graph can be used for applications such as spectral clustering. For simplicity, we will consider an L-R decomposition of \mathbf{W} , where $\mathbf{L} \in \mathbb{R}^{N \times P}$ and $\mathbf{R} \in \mathbb{R}^{N \times P}$. In this case, our goal is to measure the fidelity of the approximation $\hat{\mathbf{W}} = \mathbf{L}\mathbf{R}^T$. The corresponding optimization with quantile loss is posed as

$$\min_{\mathbf{L}, \mathbf{R}} \rho_\tau(\mathbf{W} - \mathbf{L}\mathbf{R}^T). \quad (3)$$

Setting $P = 1$ for simplicity, this becomes

$$\min_{\mathbf{l}, \mathbf{r}} \sum_{i,j} \rho_\tau(w_{ij} - l_i r_j). \quad (4)$$

When τ is high, the positive residuals will be penalized less compared to the negative residuals and hence $l_i r_j$ will underestimate w_{ij} for most i and j . Setting the negative values of $l_i r_j$ to zero will lead to a sparse similarity matrix. On the other hand, lower values of τ will correspond to lesser number of negative values in $\{l_i r_j | \forall i, j \in 1, \dots, N\}$, thereby resulting in a dense graph. Hence, we can consider the spectral decomposition of graphs with quantile penalties as a way of robustly sparsifying the graph while preserving the essential spectral structure. Note, optimizing for $\hat{\mathbf{W}}$ without imposing any spectral structural constraints ($\hat{\mathbf{W}} = \mathbf{L}\mathbf{R}^T$) will result in all the elements of $\hat{\mathbf{W}}$ being

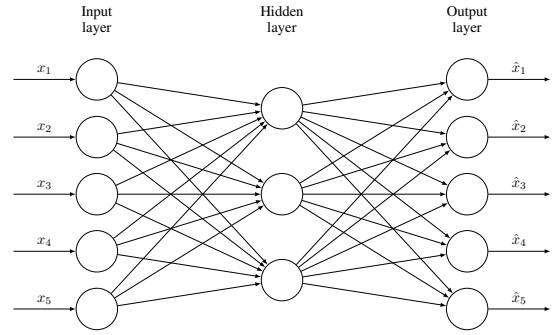


Figure 2. An example architecture with 5-dimensional inputs and 3-dimensional representations in the hidden layer.

equal to the τ^{th} quantile of the elements of \mathbf{W} . Hence it is important to preserve the spectral structure when obtaining the conditional quantile estimates.

3.3. Non-linear Spectral Decomposition using Auto-Encoders

We propose a non-linear extension to the spectral decomposition approach described in the previous section. To infer conditional quantiles of graph similarity, we solve the problem $\rho_\tau(\mathbf{W} - \hat{\mathbf{W}})$ using an auto-encoding neural network. An auto-encoder learns to reconstruct/predict its input by learning successive non-linear encodings and subsequent non-linear decodings using *hidden layers*. Even though the network is learning an identity function, by placing constraints such as low-dimensionality, sparsity, etc. on the hidden layers, useful structure learning can be performed. The hidden unit activations can then be used as feature representations. In the simplest form, we can express the reconstructed similarity as the composition, $\hat{\mathbf{W}} = g(f(\mathbf{W}))$, where f is a non-linear transformation and g is its inverse. The auto-encoder attempts to learn both f and g by restricting them to specific forms of non-linearity. In particular the following forms are widely used and well-understood: $\mathbf{H} = f(\mathbf{W}) = \phi(\Psi^T \mathbf{W})$ and $\hat{\mathbf{W}} = g(\mathbf{H}) = \phi(\Psi \mathbf{H})$ where $\Psi \in \mathbb{R}^{N \times P}$ with $P < N$ and ϕ is an elementwise non-linearity (*aka* the activation function), such as the logistic sigmoid function. Comparing this formulation to the linear decomposition $\hat{\mathbf{W}} = \mathbf{L}\mathbf{R}^T$, one can readily draw parallels: \mathbf{R}^T is the forward transformation f and \mathbf{L} is the inverse g . We can now denote the non-linear spectral learning with squared ℓ_2 loss using the objective

$$\hat{\Psi} = \underset{\Psi}{\operatorname{argmin}} \|\mathbf{W} - \phi(\Psi \phi(\Psi^T \mathbf{W}))\|_2^2. \quad (5)$$

This can be optimized using an auto-encoder where \mathbf{W} input data and Ψ are the weights between the input and hidden layers. A similar framework has been considered by Tian *et al.* [22], where the authors show the equivalence between spectral clustering and the objective of an auto-encoder. To

learn conditional quantile estimates of the input similarity matrix, the optimization can be re-posed using the quantile Huber loss ρ_τ as,

$$\hat{\Psi} = \underset{\Psi}{\operatorname{argmin}} \rho_\tau (\mathbf{W} - \phi(\Psi \phi(\Psi^T \mathbf{W}))).$$

An example architecture is provided in Figure 2 with one hidden layer. This can be generalized easily to multiple hidden layers resulting in deep and more complex representations. Interestingly, the linear L-R decomposition can be recovered by considering a single hidden layer with P units ($P < N$), and ϕ set to the identity function. As a result, $\Psi^T \mathbf{W} = \mathbf{R}^T$ and $\Psi = \mathbf{L}$. The ease of back-propagation based optimization enables the incorporation of additional constraints on the hidden layer weights and activation. In fact, this is equivalent to adding regularization terms on \mathbf{R} and \mathbf{L} in the linear case.

We implement the above approach using stochastic gradient descent with mini-batch operations and automatic differentiation [23]. The derivative of the quantile Huber loss with respect to the residual r can be obtained as

$$\rho'_\tau(r) = \begin{cases} -\tau & \text{if } r < -\tau\kappa, \\ \frac{r}{\kappa} & \text{if } r \in [-\kappa\tau, (1-\tau)\kappa], \\ (1-\tau) & \text{if } r > (1-\tau)\kappa. \end{cases} \quad (6)$$

Based on this, we create a custom automatic differentiation procedure for quantile Huber so that the gradients can be back-propagated through the network.

4. Proposed Local Scale Estimation

We propose to use the conditional quantiles of the graph similarity function to obtain estimates of the local scale for each sample in the dataset. To this end, we first construct a stochastic graph by studying the persistence of edges in reconstructed graphs at different quantiles, and then create random realizations of the neighborhood to obtain a robust estimate of the scale parameter.

4.1. Constructing Stochastic Graphs

The accuracy of spectral clustering depends, among other factors, on the appropriate choice of the scale parameter (and k in k -NN graphs). Since global neighborhood parameters are insufficient for modeling disparate sampling densities across different clusters, it is common to define a more general, dense affinity matrix that incorporates local scaling. A popular alternative to choosing a single parameter is to define the affinity between two samples \mathbf{x}_i and \mathbf{x}_j as follows:

$$w_{ij} = \exp\left(-\frac{d(\mathbf{x}_i, \mathbf{x}_j)}{\sigma_i \sigma_j}\right), \quad (7)$$

where $d(\cdot, \cdot)$ denotes an appropriate distance function and σ_i, σ_j are the local scales corresponding to the samples \mathbf{x}_i and \mathbf{x}_j respectively. For example, in [1], this parameter is defined as $\sigma_i = d(\mathbf{x}_i, \mathbf{x}_k^i)$, where \mathbf{x}_k^i is the k^{th} nearest neighbor of \mathbf{x}_i . While this graph similarity construction

Algorithm 1 Estimate the local scale for all samples in the input data \mathbf{X}

- 1) Compute the locally scaled affinity matrix, \mathbf{W}_ℓ , using (7)

Construct stochastic graph:

- 2) Train autoencoders with the quantile huber loss for \mathbf{W}_ℓ (Section 3.3), at different values of the quantile parameter τ
- 3) Compute the number of edges, β_τ , at each quantile
- 4) Construct a stochastic graph with edge probabilities given in (8)

Estimate local scale:

- 5) Draw R independent realizations of the neighborhood graph from the edge probabilities
 - 6) For each neighborhood graph, estimate the local scale of a sample as the average distance to all its neighbors
 - 7) Obtain the final local scale estimates as the median of the R realizations
-

tends to produce improved results in practice, it still relies heavily on the choice of the parameter k (set to 7 in [1]). Furthermore, a single value of k may not cluster data effectively in the presence of noise or under non-linear geometric transformations [9]. In this paper, we argue that exploring the conditional quantiles of graph similarities will enable us to obtain good estimates of local scale thereby leading to robust performances.

Our proposed approach begins by constructing a locally scaled affinity matrix using (7) and inferring the conditional quantiles using the approach described in Section 3.3. More specifically, we train a set of auto-encoding neural networks following the architecture in Figure 2 to recover the affinity matrix at different values of the quantile parameter, $\tau = \{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9\}$. In each network, we use the locally scaled affinity matrix (\mathbf{W}_ℓ) as the input and obtain the output from the decoder layer, \mathbf{W}_ℓ^τ , representing the affinity matrix at the τ^{th} quantile. An interesting observation is that this approach produces a unique set of neighborhood graphs monotonously parameterized by τ . In other words, the graph monotonously loses edges as τ increases and depending on the data distribution, becomes disconnected at a certain quantile. While it might be interesting to analyze the clusters revealed at different quantiles of the similarity function, we believe that the rate at which the graph loses its edges directly reveals the significance of the edges in recovering the underlying structure. For example, let us consider the synthetic data in Figure 3(a) that contains samples from two spirals. While the original affinity matrix is significantly dense, by ignoring edges with trivial affinities, the recovered graph at the 0.1 quantile (Figure 3(b)) contains only 4939 edges. Using only a fraction of the edges, it still recovers the underlying structure effectively. As shown in Figure 3(e), the

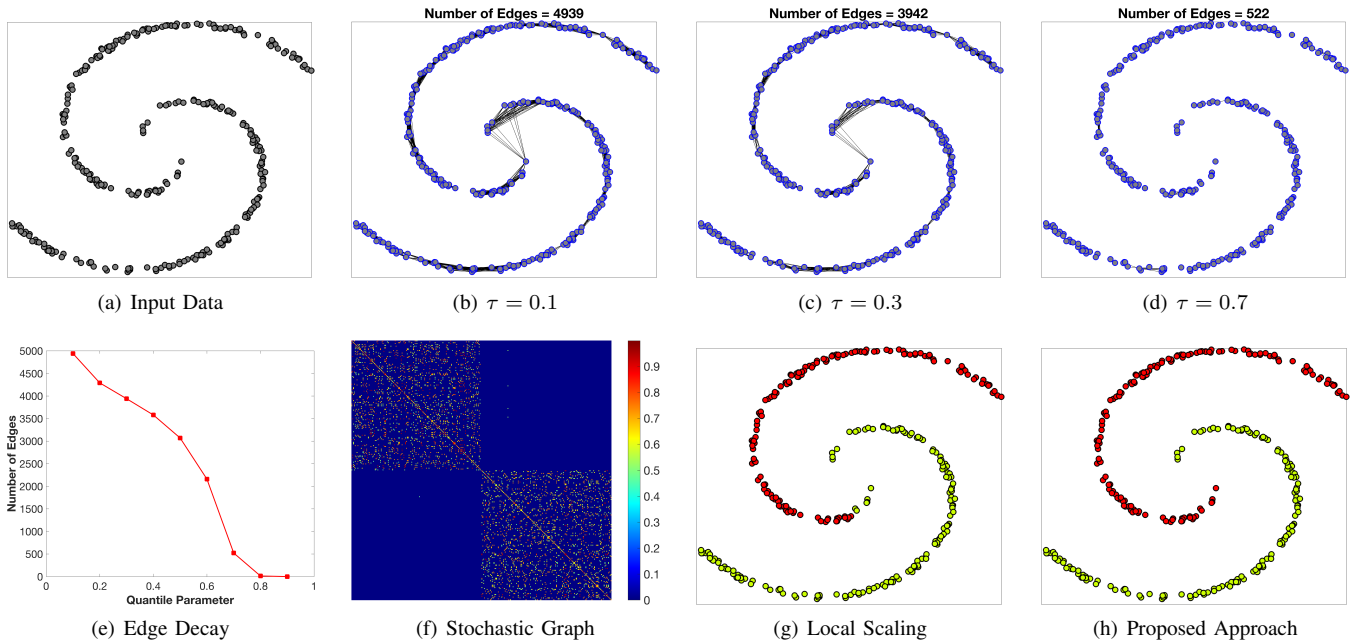


Figure 3. Two spirals dataset. We demonstrate the proposed local scale estimation from the conditional quantiles of the similarity function. The reconstructed affinities (b-d) are used to measure the rate at which edges decay as a function of τ (e) and subsequently used to construct the stochastic graph (f). Finally, the results of spectral clustering reveal the effectiveness of the proposed approach.

graph becomes increasingly sparse at higher quantiles and creates disconnected components. While intuition suggests that edges which persist at higher quantiles are crucial to recovering the underlying clusters, it is also important to note that the significance of an edge depends on the level of sparsity in the recovered graph, regardless of the quantile used. The latter observation will ensure that edges from regions with disparate densities are treated equally. Hence, we propose to construct a stochastic graph, wherein the probability of an edge is proportional to the sparsity of the graph from the highest τ at which the edge persists before disappearing.

Denoting the number of edges in the affinity matrix recovered at quantile τ as β_τ , the probability of an edge, e_{ij} is measured as

$$p(e_{ij}) = \max\left(\delta, 1 - \frac{\beta_{\hat{\tau}}}{\beta_{0.1}}\right), \quad (8)$$

where $\hat{\tau}$ corresponds to the highest quantile at which e_{ij} persists and δ is the minimal probability assigned to all edges in the graph at $\tau = 0.1$ (in all our experiments, $\delta = 0.4$). Figure 3(f) illustrates the estimated stochastic graph for the two spirals dataset.

4.2. Algorithm

In this section, we describe the proposed algorithm for obtaining robust local scale estimates using the stochastic graphs. Existing approaches such as consensus clustering [2] directly use stochastic graphs as inputs to spectral clustering. An important downside of such approaches is that they

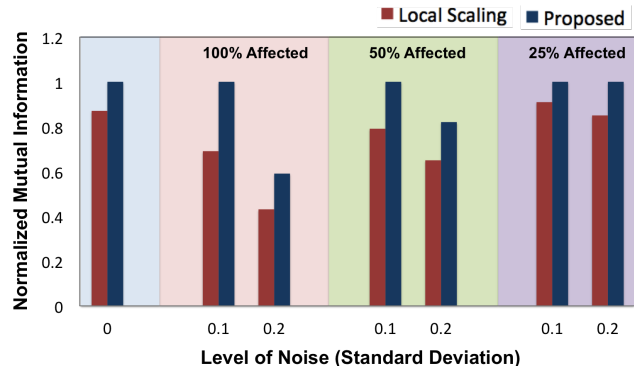


Figure 4. Impact of noise on clustering performance of local scaling in 7 and the proposed approach with the twospirals dataset in Figure 3(a).

require multiple iterations of refinement to construct robust graph affinities. In contrast, we propose to generate multiple random realizations of the neighborhood using the inferred edge probabilities and estimate the local scales in each of the realizations independently. Denoting a random realization of the stochastic graph as $\tilde{\mathbf{W}}^{(1)}$, we estimate the local scale for a sample \mathbf{x}_i as

$$\sigma_i^{(1)} = \frac{1}{|\mathcal{N}_i|} \sum_{j \in \mathcal{N}_i} d(\mathbf{x}_i, \mathbf{x}_j), \quad (9)$$

where \mathcal{N}_i denotes the set of neighbors of \mathbf{x}_i identified using $\tilde{\mathbf{W}}^{(1)}$. The estimate of the local scale is obtained as the median of the local scales at each of the realizations, $\sigma_i = \text{median}\left([\sigma_i^{(1)}, \sigma_i^{(2)}, \dots, \sigma_i^{(R)}]\right)$, where R is the total

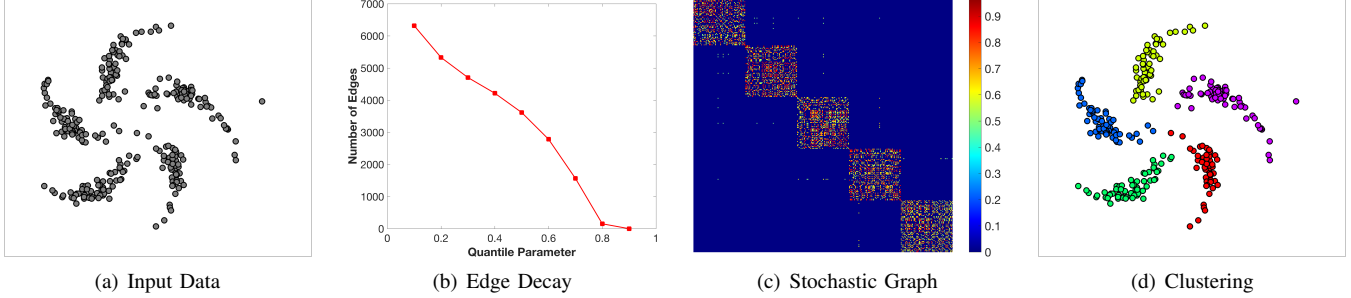


Figure 5. Pin wheel dataset. Robust scale estimation leads to accurate recovery of the underlying clusters.

number of independent realizations. Finally, the estimated scaled are used to construct the graph affinities as shown in (7). The steps of the algorithm are summarized in Algorithm 1. Figure 3(g) and (h) show the spectral clustering results obtained with the local scaling approach in (7) and the proposed approach respectively. As it can be observed, using the conditional quantiles leads to robust scale estimates and consequently the proposed approach accurately recovers the two spirals. We also analyze the sensitivity of the resulting graphs to noise in the data, by adding random Gaussian noise ($\sigma^{noise} = \{0.1, 0.2\}$) to different number of randomly chosen samples (25%, 50%, 100%). From the results in Figure 4, the proposed scaling approach is highly resilient to the noise in data in all cases. Figure 5 illustrates the estimated probabilities and clustering performance of the proposed method on the pinwheel dataset. Another interesting observation is that the choice of parameter k in the construction of the initial affinity \mathbf{W}_ℓ does not significantly affect the performance of the proposed approach. For example, we evaluated the clustering results on the twospirals data by varying k between 5 and 15 and found that our method consistently recovers the true clusters effectively.

5. Experiments

5.1. Spectral Clustering

In this section, we evaluate the performance of the proposed graph construction approach in spectral clustering. We have explored our approach using a number of datasets from the UCI machine learning repository [24]. These data sets are characterized for having clusters of varying density, scale and shape—where spectral algorithms using a global scale are known to perform poorly. We evaluate the clustering performance using the Normalized Mutual Information (NMI) score, defined as

$$NMI(Y; \hat{Y}) = \frac{2I(Y; \hat{Y})}{H(Y) + H(\hat{Y})}, \quad (10)$$

where $I(Y; \hat{Y})$ is the mutual information between sets Y and \hat{Y} , and $H(Y), H(\hat{Y})$ are the entropies of the two sets.

For comparison, we computed the clustering performance of the local scaling approach in [1] and the path-based similarities in [10]. In addition, we implemented a

Algorithm 2 Perform label propagation using greedy walk on the graph laplacian kernel

Input: Data \mathbf{X} with the first C samples labeled, graph affinity \mathbf{W}

- 1) Compute the graph laplacian using (11) and build the graph laplacian kernel $\mathbf{K} = \mathbf{L}^\dagger$
- 2) Measure Euclidean distances on the RKHS to construct the distance matrix \mathbf{S}

For each unlabeled example j , initialize $ind = j$, $iter = 0$,

- 3) **while** $iter < maxwalk$:
 - a) Store the distances $\gamma^{iter}(\mathbf{x}_{ind}, \mathbf{x}_i) = \min(S(ind, i), \gamma^{iter-1}(\mathbf{x}_{ind}, \mathbf{x}_i)), \forall i$
 - b) Determine the nearest neighbor for \mathbf{x}_{ind} , $t = \arg \min \mathbf{S}(ind, :)$
 - c) **if** $t \in \{1, \dots, C\}$, **break**
 - d) $ind = t, iter = iter + 1$
 - 4) $\text{Label}(\mathbf{x}_j) = \arg \min_i \gamma^{iter}$
-

consensus clustering approach, which used different values for the parameter k to obtain the local scale estimates and adopted the consensus inference approach in [2] to perform clustering. For the proposed approach, we fixed the number of random realizations in Algorithm 1, $R = 25$. Table 1 shows the NMI obtained using the different techniques and our approach outperforms the baseline techniques in all cases.

5.2. Label Propagation using a Single Example

Another interesting application of graph similarities is in propagating labels to a large set of test examples using a limited training set. In its extreme case, the problem of classification with only one labeled example per class can benefit significantly from robust graphs [25]. In this section, we develop a greedy label propagation strategy based on graph similarities and evaluate the effectiveness of the proposed graph construction in this application.

5.2.1. Greedy Walk on Graph Laplacian Kernels. Existing methods for label propagation often rely on graph-based

Dataset	Local Scaling [1]	Consensus Clustering [2]	Path-based Similarities [10]	Proposed Approach
Two spirals	0.87	1.0	0.92	1.0
Pinwheel	0.95	0.97	0.97	1.0
Glass	0.37	0.38	0.39	0.42
Breast Cancer	0.79	0.82	0.78	0.85
Wine	0.91	0.89	0.91	0.93
E-coli	0.57	0.63	0.58	0.63
Leaf	0.66	0.73	0.68	0.75
Parkinson	0.31	0.32	0.33	0.37

TABLE 1. PERFORMANCE COMPARISON OF SPECTRAL CLUSTERING. WE COMPARE THE PROPOSED ROBUST AFFINITIES WITH BASELINE GRAPH CONSTRUCTION TECHNIQUES. IN EACH THE METHOD THAT PRODUCES THE HIGHEST NMI IS MARKED IN BOLD.

methods with local and global consistency constraints [26]. In particular, constructing appropriate RKHS (Reproducing Kernel Hilbert Space) kernels by transforming the spectrum of the graph over labeled and unlabeled data together has been effective. Hence, we adopt the approach in [25] to construct a graph laplacian kernel from weighted neighborhood graphs for label propagation.

Given the graph affinity matrix $\mathbf{W} \in \mathbb{R}^{N \times N}$, we construct the normalized graph laplacian \mathbf{L} (not to be confused with the matrix \mathbf{L} in Section 3.2) as follows:

$$\mathbf{L} = \mathbf{D}^{-\frac{1}{2}}(\mathbf{D} - \mathbf{W})\mathbf{D}^{-\frac{1}{2}}, \quad (11)$$

where \mathbf{D} is the degree matrix whose diagonal entries are defined as $D_{ii} = \sum_j w_{ij}$. Let $\mathcal{F}(G)$ denote the linear space of real-valued functions defined on the graph G and $\{\lambda_i, \mathbf{u}_i\}_{i=1}^N$ denote the eigen spectrum of the corresponding laplacian \mathbf{L} . Now, we define a Hilbert space of functions on G , $\mathcal{H}(G) = \{\mathbf{g} | \mathbf{g}^T \mathbf{u}_i = 0, \forall i\}$, which is a linear subspace of $\mathcal{F}(G)$ orthogonal to the eigenvectors of \mathbf{L} with zero eigenvalues. Similar to the analysis in [27], we can show that the pseudo-inverse of \mathbf{L} is the reproducing kernel of $\mathcal{H}(G)$. The resulting matrix $\mathbf{K} = \mathbf{L}^\dagger$ is referred to as the *graph laplacian kernel*.

In our problem setup, the data contains C classes with one labeled example per class. Without loss of generality, we assume that the first C samples in the input dataset correspond to the labeled examples and the rest are unlabeled. Note that, the distance between two samples can be obtained by measuring their Euclidean distance in the RKHS,

$$d(\mathbf{x}_i, \mathbf{x}_j) = \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|_2 = K_{ii} + K_{jj} - 2K_{ij}. \quad (12)$$

Using this, we build the distance matrix $\mathbf{S} \in \mathbb{R}^{N \times N}$ and employ a greedy walk scheme, shown in Algorithm 2, for propagating the labels. In all our experiments, the maximum length of the greedy walk, *maxwalk*, was fixed at 5. The performance of this propagation strategy relies heavily on the robustness of the graph similarity \mathbf{W} to noise and outliers.

5.2.2. Results. We evaluate the performance of different locally scaled affinity matrices in label propagation using a variety of challenging binary classification datasets from

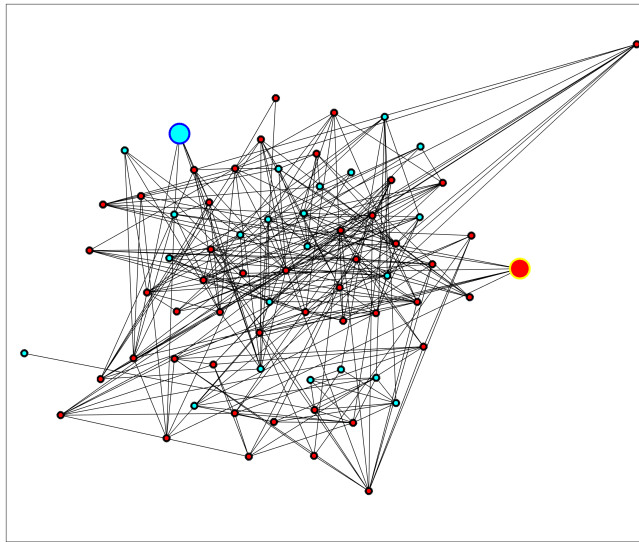
Dataset	Local Scaling [1]	Path-based Similarities [10]	Proposed Approach
Blood Transfusion	76.05	77.1	84.9
Breast Cancer	82.7	84.9	90.55
Echocardiogram	70.06	73.2	88.49
Kidney Disease	66.8	67.5	71.9
SPECT Heart	70	68.5	86
Thoracic Surgery	68.3	66.2	75.8
Arcene	59	61.4	71.3

TABLE 2. PERFORMANCE OF DIFFERENT GRAPH SIMILARITY CONSTRUCTION APPROACHES IN CLASSIFICATION USING A SINGLE EXAMPLE.

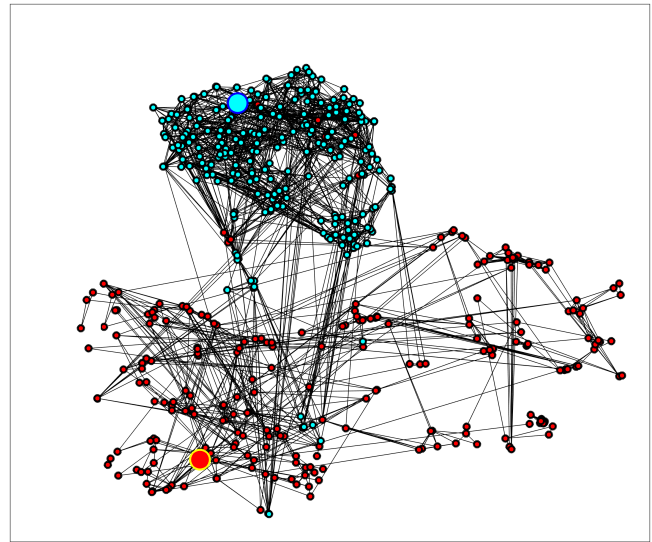
the UCI repository. In each dataset, we randomly chose one training example from each class and computed the accuracy of the label propagation scheme in Algorithm 2. We repeated the experiment for 10 independent trials and the average classification accuracies are shown in Table 2. As in the spectral clustering experiments, we compared the proposed approach to (7) and path-based similarities [10]. The effectiveness of the proposed neighborhood graphs is apparent from the improvements in the classification performance (as high as 18%) over the baseline methods. Figure 6 illustrates the graphs obtained using the proposed approach for the *echocardiogram* and *breast cancer* datasets.

References

- [1] L. Zelnik-manor and P. Perona, “Self-tuning spectral clustering,” in *Advances in Neural Information Processing Systems 17*. MIT Press, 2005, pp. 1601–1608.
- [2] S. Race and C. Meyer, “A flexible iterative framework for consensus clustering,” *arXiv preprint arXiv:1408.0972*, 2014.
- [3] M. Maier, U. von Luxburg, and M. Hein, “How the result of graph clustering methods depends on the construction of the graph,” *ESAIM: Probability and Statistics*, vol. 17, pp. 370–418, Jan. 2013.
- [4] E. Biçici and D. Yuret, “Locally scaled density based clustering,” in *Adaptive and Natural Computing Algorithms: 8th International Conference, ICANNGA 2007, Warsaw, Poland, April 11-14, 2007, Proceedings, Part I*, 2007, pp. 739–748.



(a) Echocardiogram



(b) Breast Cancer

Figure 6. Locally scaled affinities for the label propagation experiments. The two training examples are shown as the bigger circles. The 2-D embeddings of the samples are created, using the t-SNE algorithm, for visualization.

- [5] L. Huang, D. Yan, N. Taft, and M. I. Jordan, "Spectral clustering with perturbed data," in *Advances in Neural Information Processing Systems 21*, 2009, pp. 705–712.
- [6] U. von Luxburg, M. Belkin, and O. Bousquet, "Consistency of spectral clustering," *Ann. Statist.*, vol. 36, no. 2, pp. 555–586, 04 2008.
- [7] B. Nadler and M. Galun, "Fundamental limitations of spectral clustering methods," in *Advances in Neural Information Processing Systems 19*, B. Schölkopf, J. Platt, and T. Hoffman, Eds. MIT Press, 2007.
- [8] Z. Li, J. Liu, S. Chen, and X. Tang, "Noise robust spectral clustering," in *2007 IEEE 11th International Conference on Computer Vision*, Oct 2007, pp. 1–8.
- [9] C. D. Correa and P. Lindstrom, "Locally-scaled spectral clustering using empty region graphs," in *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2012, pp. 1330–1338.
- [10] H. Chang and D.-Y. Yeung, "Robust path-based spectral clustering with application to image segmentation," in *Tenth IEEE International Conference on Computer Vision (ICCV'05) Volume 1*, vol. 1, Oct 2005, pp. 278–285 Vol. 1.
- [11] M. Potamias, F. Bonchi, A. Gionis, and G. Kollios, "K-nearest neighbors in uncertain graphs," *Proc. VLDB Endow.*, vol. 3, no. 1-2, pp. 997–1008, Sep. 2010.
- [12] C. D. Meyer and C. D. Wessell, "Stochastic data clustering," *SIAM Journal on Matrix Analysis and Applications*, vol. 33, no. 4, pp. 1214–1236, 2012.
- [13] R. Koenker and G. Bassett Jr, "Regression quantiles," *Econometrica: Journal of the Econometric Society*, pp. 33–50, 1978.
- [14] I. Takeuchi, Q. V. Le, T. D. Sears, and A. J. Smola, "Nonparametric quantile estimation," *The Journal of Machine Learning Research*, vol. 7, pp. 1231–1264, 2006.
- [15] H. Zou and M. Yuan, "Regularized simultaneous model selection in multiple quantiles regression," *Computational Statistics & Data Analysis*, vol. 52, no. 12, pp. 5296–5304, 2008.
- [16] R. Koenker and O. Geling, "Reappraising medfly longevity: A quantile regression survival analysis," *Journal of the American Statistical Association*, vol. 96, pp. 458–468, 2001.
- [17] K. N. Ramamurthy, K. R. Varshney, and M. Singh, "Quantile regression for workforce analytics," in *Proc. IEEE GlobalSIP*, 2013.
- [18] R. Koenker and K. F. Hallock, "Quantile regression," *Journal of Economic Perspectives, American Economic Association*, pp. 143–156, 2001.
- [19] A. Bhatlele, A. R. Titus, J. J. Thiagarajan, N. Jain, T. Gamblin, P. T. Bremer, M. Schulz, and L. V. Kale, "Identifying the culprits behind network congestion," in *Parallel and Distributed Processing Symposium (IPDPS), 2015 IEEE International*, May 2015, pp. 113–122.
- [20] K. N. Ramamurthy, A. Y. Aravkin, and J. J. Thiagarajan, "Beyond l2-loss functions for learning sparse models," in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2016, pp. 4692–4696.
- [21] A.Y. Aravkin, P. Kambadur, A. Lozano, and R. Luss, "Orthogonal matching pursuit for sparse quantile regression," in *Data Mining (ICDM), International Conference on*. IEEE, 2014, pp. 11–19.
- [22] F. Tian, B. Gao, Q. Cui, E. Chen, and T.-Y. Liu, "Learning deep representations for graph clustering," in *AAAI*, 2014, pp. 1293–1299.
- [23] D. Maclaurin, D. Duvenaud, and R. P. Adam, "Autograd: Effortless gradients in pure Numpy," in *ICML AutoML Workshop*, 2015.
- [24] M. Lichman, "UCI machine learning repository," 2013. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [25] H. Chang and D.-Y. Yeung, "Graph laplacian kernels for object classification from a single example," in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, vol. 2, June 2006, pp. 2011–2016.
- [26] D. Zhou, O. Bousquet, T. N. Lal, J. Weston, and B. Schölkopf, "Learning with local and global consistency," in *Advances in Neural Information Processing Systems 16*, 2004, pp. 321–328.
- [27] M. Herbster, M. Pontil, and L. Wainer, "Online learning over graphs," in *Proceedings of the 22nd International Conference on Machine Learning*, ser. ICML '05. ACM, pp. 305–312.