A rigorous framework for data clustering

Thomas Strohmer Department of Mathematics University of California, Davis

Statistics Seminar Series UC Davis, Nov. 7, 2019



Acknowledgements

Joint work with Shuyang Ling (NYU)



This work is sponsored by the NSF-DMS and the National Geospatial-Intelligence Agency.







Trying to detect Alzheimer's from MRIs





Trying to detect Alzheimer's from MRIs



- 75 year old Control
- 75 year old MCI

75 year old AD



Unsupervised Learning via diffusion maps



Diffusion component representation of two normal brains (dark blue and light blue) and two Alzheimer brains (red, yellow)

Sometimes our diffusion component analysis approach worked very well, and sometimes it failed. Why?

Alzheimer MRI data are too complicated to develop a thorough theoretical analysis.

Look at existing theory for simpler situations.



Sometimes our diffusion component analysis approach worked very well, and sometimes it failed. Why?

Alzheimer MRI data are too complicated to develop a thorough theoretical analysis.

Look at existing theory for simpler situations.

Surprise:

No useful theory about performance of diffusion maps.



Sometimes our diffusion component analysis approach worked very well, and sometimes it failed. Why?

Alzheimer MRI data are too complicated to develop a thorough theoretical analysis.

Look at existing theory for simpler situations.

Surprise:

No useful theory about performance of diffusion maps.

Bigger surprise:

There is no useful theory about performance of spectral clustering, despite the fact that it has been around for decades!

Spectral clustering: learning the "shape" of data



Existing theory for the performance of spectral clustering either assumes that the clusters are infinitely far apart or that the associated graph has disconnected components as expressed in a "perfect" binary weight matrix.

Spectral clustering, graph cuts, community detection



Three NP-hard problems ...



Data clustering and unsupervised learning

Question: Given a set of *N* data points in \mathbb{R}^d , how to partition them into *k* clusters based on their similarity?



K-means: minimize the total within-cluster sum of squared error to estimate the partition.

$$\min_{\{\Gamma_{I}\}_{l=1}^{k}} \sum_{l=1}^{k} \underbrace{\sum_{i \in \Gamma_{I}} \left\| x_{i} - \underbrace{\frac{1}{|\Gamma_{I}|} \sum_{i \in \Gamma_{I}} x_{i}}_{\text{within-cluster sum of squares}} \right\|^{2}}_{\text{within-cluster sum of squares}}$$

where $\{\Gamma_l\}_{l=1}^k$ is a partition of $\{1, \dots, N\}$.



K-means only works for datasets with individual clusters:

isotropic and within convex boundaries, well-separated





Kernel k-means and nonlinear embedding

Goal: map the data into a feature space so that they are well-separated and *k*-means would work.



How: locally-linear embedding, isomap, multidimensional scaling, Laplacian eigenmaps, diffusion maps, etc.

Focus: We will focus on Laplacian eigenmaps. Spectral clustering consists of Laplacian eigenmaps followed by *k*-means clustering.



Spectral clustering

Spectral clustering¹ consists of two steps:

- Laplacian eigenmaps
- "Rounding" procedure (e.g. k-means)



Laplacian eigenmap

¹[Luxburg, 07], [Shi, Malik, 02], [Belkin, Niyogi, 03]



The Graph Laplacian

One key ingredient of spectral clustering is the graph Laplacian. Given $\{\mathbf{x}_i\}_{i=1}^N \in \mathbb{R}^d$, we construct a similarity (weight) matrix \boldsymbol{W} .



The graph Laplacian is

L = D - W

where $\boldsymbol{D} = \text{diag}(\boldsymbol{W} \mathbf{1}_N)$ is the degree matrix.

The spectra of *L* are related to the graph's connectivity, especially $\lambda_2(L)$.

The *incidence matrix* ∇ of a graph is

- $\nabla := \begin{cases} \nabla = -1 & \text{if } v \text{ is the initial vertex of the edge} \\ \nabla := \begin{cases} \nabla = 1 & \text{if } v \text{ is the terminal vertex of the edge} \end{cases}$

$$\left(\nabla = \mathbf{0} \quad \text{if } \mathbf{v} \text{ is not in the edge}, \right.$$

where we assume each edge has an arbitrary (but fixed) orientation. If we have an undirected graph, we can obtain the incidence matrix by choosing a (fixed) orientation of the edges. If f is a function acting on the vertices, then

$$\nabla f = \{f(\mathbf{v}_j) - f(\mathbf{v}_i)\}_{i,j}$$

Hence ∇ is a kind of difference operator.

The unweighted graph Laplacian *L* can be written as $L = \nabla^T \nabla$.



Denote the eigenvectors of \boldsymbol{L} by $\boldsymbol{u}_1, \cdots, \boldsymbol{u}_N$ with associated real, non-negative eigenvalus $\lambda_1, \cdots, \lambda_N$.

The Laplacian eigenmap is defined as

$$\begin{bmatrix} \varphi(\boldsymbol{x}_1) \\ \vdots \\ \varphi(\boldsymbol{x}_N) \end{bmatrix} := \underbrace{[\boldsymbol{u}_1, \cdots, \boldsymbol{u}_n]}_{\boldsymbol{U}} \in \mathbb{R}^{N \times n}$$

where $\{u_l\}_{l=1}^n$ are the eigenvectors w.r.t. the *n* smallest eigenvalues.



Laplacian eigenmaps, k-means, spectral clustering



 φ maps data in \mathbb{R}^d to \mathbb{R}^n ; coordinates in terms of eigenvectors:

$$\varphi:\underbrace{\mathbf{X}_{i}}_{\mathbb{R}^{d}}\longrightarrow\underbrace{\varphi(\mathbf{X}_{i})}_{\mathbb{R}^{n}}.$$

Then we apply *k*-means to $\{\varphi(\mathbf{x}_i)\}_{i=1}^N$ to perform clustering.



Pros and Cons of spectral clustering

Pros:

- Spectral clustering enjoys high popularity and conveniently applies to various settings.
- Supposedly works better than vanilla *k*-means

Cons:

- Rigorous justification of spectral clustering is still lacking.
- Two-step procedure complicates the theoretic analysis.

Our goal:

- A different route: **convex relaxation** of spectral clustering.
- Establish a theoretical framework for spectral clustering.

A graph cut perspective

The matrix W is viewed as a weight matrix of a graph with N vertices.



Partitioning the dataset into k clusters is equivalent to finding a k-way graph cut such that any pair of induced subgraphs is not well-connected.

Graph cut

The cut is defined as the weight sum of edges whose two ends are in different subsets,

$$\operatorname{cut}(\Gamma,\Gamma^c):=\sum_{i\in\Gamma,j\in\Gamma^c}w_{ij}$$

where Γ is a subset of vertices and Γ^c is its complement.

Warning: unbalanced cuts!





RatioCut

The ratio cut induced by the partition $\{\Gamma_a\}_{a=1}^k$ is equal to

$$\mathsf{RatioCut}(\{\Gamma_a\}_{a=1}^k) = \sum_{a=1}^k \frac{\mathsf{cut}(\Gamma_a, \Gamma_a^c)}{|\Gamma_a|}.$$

In particular, if k = 2,

$$\mathsf{RatioCut}(\Gamma,\Gamma^c) = \frac{\mathsf{cut}(\Gamma,\Gamma^c)}{|\Gamma|} + \frac{\mathsf{cut}(\Gamma,\Gamma^c)}{|\Gamma^c|}.$$

Minimizing RatioCut is NP-hard in general!



RatioCut and the graph Laplacian

Let $\mathbf{1}_{\Gamma_a}(\cdot)$ be an indicator vector which maps a vertex to a vector in \mathbb{R}^N via

$$\mathsf{I}_{\Gamma_a}(I) = egin{cases} \mathsf{1}, & I \in \Gamma_a, \ \mathsf{0}, & I \notin \Gamma_a. \end{cases}$$

Relating RatioCut to the graph Laplacian

$$\operatorname{cut}(\Gamma_{a},\Gamma_{a}^{c}) = \left\langle \boldsymbol{L}, \ \boldsymbol{1}_{\Gamma_{a}}\boldsymbol{1}_{\Gamma_{a}}^{\top} \right\rangle = \boldsymbol{1}_{\Gamma_{a}}^{\top}\boldsymbol{L}\boldsymbol{1}_{\Gamma_{a}}$$

$$\operatorname{RatioCut}(\{\Gamma_{a}\}_{a=1}^{k}) = \sum_{a=1}^{k} \frac{1}{|\Gamma_{a}|} \left\langle \boldsymbol{L}, \ \boldsymbol{1}_{\Gamma_{a}}\boldsymbol{1}_{\Gamma_{a}}^{\top} \right\rangle = \left\langle \boldsymbol{L}, \ \boldsymbol{X} \right\rangle,$$

$$\underbrace{\boldsymbol{X} := \sum_{a=1}^{k} \frac{1}{|\Gamma_{a}|} \boldsymbol{1}_{\Gamma_{a}}\boldsymbol{1}_{\Gamma_{a}}^{\top}}_{\text{a block-diagonal projection matrix}}$$

Convex relaxation of RatioCut

Question: What constraints does X satisfy for any given $\{\Gamma_a\}_{a=1}^k$? $\begin{bmatrix} \frac{1}{n_1} J_{n_1} & \mathbf{0} & \cdots & \mathbf{0} \end{bmatrix}$

$$\boldsymbol{X} = \boldsymbol{\Pi} \begin{bmatrix} \overline{n_1} & \overline{n_1} & \overline{n_1} & \overline{n_1} & \overline{n_1} \\ \mathbf{0} & \frac{1}{n_2} & \overline{n_2} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \frac{1}{n_k} & \mathbf{J}_{n_k} \end{bmatrix} \boldsymbol{\Pi}^\top$$

where Π is a permutation matrix, and J_n is an $n \times n$ "1" matrix.

Convex sets

Given a partition $\{\Gamma_a\}_{a=1}^k$, the corresponding **X** satisfies

- **X** is positive semidefinite, $X \succeq 0$
- **X** is nonnegative, $X \ge 0$ entrywise
- the constant vector is an eigenvector of $X: X1_N = 1_N$
- the trace of **X** equals k, i.e., $Tr(\mathbf{X}) = k$

RatioCut-SDP - SDP relaxation of RatioCut We relax the ratio cut by

$$\min_{\boldsymbol{Z} \in \mathbb{R}^{N \times N}} \langle \boldsymbol{L}, \boldsymbol{Z} \rangle \text{ s.t. } \boldsymbol{Z} \succeq 0, \ \boldsymbol{Z} \ge 0, \ \mathsf{Tr}(\boldsymbol{Z}) = k, \ \boldsymbol{Z} \boldsymbol{1}_N = \boldsymbol{1}_N.$$

Advantage: It is a semidefinite program² which is solvable in polynomial time.

Question: When does convex relaxation give exact recovery?



²[Peng, Wei, 07], [Mixon, etc, 16], [Xing, Jordan, 03] ...

Intuition



 Intra-cluster: algebraic connectivity

 $\lambda_2(\mathbf{L}_a),$

the second smallest eigenvalue of the Laplacian associated with the *a*-th cluster.

Inter-cluster connectivity: the maximal outer-cluster degree.

$$d_{\text{outer,max}} = \max_{i} \sum_{j: \text{ not share membership} \ \text{with node } i} w_{ij}$$

Note that both $\lambda_2(\mathbf{L}_a)$ and $d_{\text{outer,max}}$ are determined by $\{\Gamma_a\}_{a=1}^k$.

Finding the optimal graph cut via SDP relaxation

Main theorem [Ling, Strohmer, FOCM]

• If a graph cut $\{\Gamma_a\}_{a=1}^k$ satisfies

$$d_{ ext{outer,max}} < rac{1}{4} \min_{1 \leq a \leq k} \lambda_2(\boldsymbol{L}_a)$$

then it is the globally optimal graph cut under Ratiocut.

- The RatioCut convex relaxation recovers \boldsymbol{X} associated to $\{\Gamma_l\}_{l=1}^k$ exactly!
- Even though finding the optimal Ratiocut is NP-hard, there is a regime where a poly-time algorithm works!
- Purely deterministic and depends on the algebraic properties of Laplacian. Near-optimal.
- A similar result holds for normalized cut with an interpretation from a random walk perspective.



A short tour of the proof - Game of Cones

We are dealing with

$$\min_{\boldsymbol{Z} \in \mathbb{R}^{N \times N}} \langle \boldsymbol{L}, \boldsymbol{Z} \rangle \text{ s.t. } \boldsymbol{Z} \succeq 0, \ \boldsymbol{Z} \ge 0, \underbrace{\mathsf{Tr}(\boldsymbol{Z}) = k, \ \boldsymbol{Z} \mathbf{1}_N = \mathbf{1}_N}_{\text{Linear constraints: } \mathcal{A}(\boldsymbol{Z}) = \boldsymbol{b}}.$$

Use Lagrangian duality:

$$\mathcal{L}(oldsymbol{Z},oldsymbol{\lambda}) = \langle oldsymbol{L},oldsymbol{Z}
angle - \langle \mathcal{A}(oldsymbol{Z}) - oldsymbol{b},oldsymbol{\lambda}
angle$$

where \boldsymbol{Z} is in the positive semidefinite and nonnegativity cone, denoted by \mathcal{K} .

Dual program

$$\max \left< oldsymbol{\lambda}, oldsymbol{b} \right>$$
 s.t. $\mathcal{A}^*(oldsymbol{\lambda}) + oldsymbol{L} \in \mathcal{K}^*$

where \mathcal{K}^* is the dual cone³.

 ${}^{3}\mathcal{K}^{*}:=\{\boldsymbol{\boldsymbol{z}}:\langle\boldsymbol{\boldsymbol{z}},\boldsymbol{\boldsymbol{x}}\rangle\geq\boldsymbol{0},\forall\boldsymbol{\boldsymbol{x}}\in\mathcal{K}\}$



Dual program

```
\max \left< oldsymbol{\lambda}, oldsymbol{b} \right> s.t. \mathcal{A}^*(oldsymbol{\lambda}) + oldsymbol{L} \in \mathcal{K}^*
```

where \mathcal{K}^* is the dual cone⁴.

- Construct λ to certify the optimality of a graph cut, with help of spectral graph theory.
- Construction of λ is rather technical and involves methods derived by X. Li, Y. Li, S. Ling, T.S. and K. Wei in [When Do Birds of a Feather Flock Together? K-Means, Proximity, and Conic Programming, Mathematical Programming, 2018.]

$${}^{4}\mathcal{K}^{*}:=\{\boldsymbol{\textit{z}}:\langle\boldsymbol{\textit{z}},\boldsymbol{\textit{x}}\rangle\geq0,\forall\boldsymbol{\textit{x}}\in\mathcal{K}\}$$

$$\int_{T}^{S}$$

Some observations

Known theoretical bounds for SDP relaxation of k-means clustering have the undesirable property that they do depend on the number of clusters.

This is not the case for our SDP relaxation of spectral clustering. The bounds are independent of the number of clusters

Our theorem also yields a certificate of optimality: Assume someone gives us a graph partition $\{\Gamma_a\}_{a=1}^k$ and claims it is the optimal RatioCut. This is in general very difficult to verify. Easy test: If the partition $\{\Gamma_a\}_{a=1}^k$ satisfies

$$d_{\text{outer,max}} < \frac{1}{4} \min_{1 \le a \le k} \lambda_2(\boldsymbol{L}_a),$$

then it is indeed optimal.



Spectral clustering for two concentric circles

We consider

$$\boldsymbol{x}_{1,i} = \begin{bmatrix} \cos(\frac{2\pi i}{n}) \\ \sin(\frac{2\pi i}{n}) \end{bmatrix}, \ 1 \le i \le n; \qquad \boldsymbol{x}_{2,j} = \frac{m}{n} \begin{bmatrix} \cos(\frac{2\pi j}{m}) \\ \sin(\frac{2\pi j}{m}) \end{bmatrix}, \ 1 \le j \le m$$

where m > n.





Spectral clustering for two concentric circles

Corollary (Ling, Strohmer, 2018) The RatioCut-SDP recovers the underlying two clusters exactly if

$$\Delta = \Omega(n^{-1}).$$

minimal separation

where n^{-1} is the distance of two adjacent points on one circle.

It is near-optimal.





Community detection under stochastic block model

Each community has n/2 nodes

- if member i and j are in the same community, pair them with probability p;
- if member *i* and *j* are in different communities, pair them with probability *q*

Assume $\mathbf{p} > \mathbf{q}$. Given the adjacency matrix, how to find out the underlying communities?



р



Corollary (Ling, Strohmer, 2018) Let $p = \frac{\alpha \log n}{n}$ and $q = \frac{\beta \log n}{n}$. The RatioCut-SDP recovers the underlying communities exactly if

$$lpha > \mathbf{13} \Big(\sqrt{eta} + \sqrt{\mathbf{2}}\Big)^{\mathbf{2}}$$

with high probability.

The information theoretic lower bound for exact recovery⁵ is

$$\alpha > (\sqrt{\beta} + \sqrt{2})^2.$$

If $\alpha < (\sqrt{\beta} + \sqrt{2})^2$, exact recovery is impossible.





Open problem

Suppose there are *n* data points drawn from a probability density function p(x) supported on a manifold \mathcal{M} . How can we estimate the second smallest eigenvalue of the graph Laplacian (either normalized or unnormalized) given the kernel function Φ and σ ?



The solution will require tools from empirical process, differential geometry, spectral graph theory, etc⁶.



⁶[Singer, 06], [Belkin, etc, 08], [Trillos, etc, 16], ...

Assume we know for a small number of points to which cluster they belong. How can we incorporate this knowledge?

The SDP approach provides a natural framework for clustering when a few labels are known.

SDP with additional linear constraints.

$$\begin{array}{l} \min \left< \boldsymbol{L}, \boldsymbol{Z} \right> \text{ subject to } \boldsymbol{Z} \succeq \boldsymbol{0}, \\ \boldsymbol{Z} \geq \boldsymbol{0}, \\ \mathrm{Tr}(\boldsymbol{Z}) = \boldsymbol{k}, \\ \boldsymbol{Z} \boldsymbol{1}_N = \boldsymbol{1}_N \\ \boldsymbol{Z}_{j,k} = \boldsymbol{Z}_{m,n} \quad j,k,m,n \in \boldsymbol{\mathcal{I}} \end{array}$$



Semisupervised clustering



No good theory yet



(Borrowing from slides by Stefan Steinerberger)

t-SNE:

Main problem: given a set of high-dimensional points $\{x_1, \ldots, x_n\} \in \mathbb{R}^d$ we would like to get an 'equivalent' representation $\{y_1, \ldots, y_n\} \in \mathbb{R}^2$ so we can have a look. Somewhat ill-posed but 'clusters should remain clusters.'

Answer: **t-distributed stochastic neighborhood embedding** (van der Maaten & Hinton, 2008)

Main idea: Turn configurations of points into probability distributions and force these distributions to be similar.



MNIST



Laplacian eigenmaps









- Create a probability distribution in high dimensions.
- Create another on a set of points in two dimensions.
- Measure the KL-divergence between them and move the points in low dimensions around so that KL becomes small.

KL Divergence (Kullback & Leibler, 1951)

$$D_{\mathit{KL}}(\mathit{P}||\mathit{Q}) = \sum_{i} \mathit{P}(i) \log rac{\mathit{P}(i)}{\mathit{Q}(i)}$$



t-SNE algorithm

Given a set of $\{x_1, \ldots, x_n\} \in \mathbb{R}^d$, t-SNE searches for $\{y_1, \ldots, y_n\} \in \mathbb{R}^2$ that minimizes the loss

$$L(y_1,...,y_n) = -\sum_{i=1}^n \sum_{j=1,j\neq i}^n p(x_i,x_j) \log \frac{q(y_i,y_j)}{p(x_i,x_j)},$$

where the functions p and q are given by

$$p(x_i, x_j) = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{\ell \neq i} \exp(-\|x_i - x_\ell\|^2 / 2\sigma_i^2)}$$

and

$$q(y_i, y_j) = \frac{\left(1 + (\|y_i - y_j\|^2)^{-1}\right)}{\sum_{\ell \neq i} \left(1 + (\|y_i - y_\ell\|^2)^{-1}\right)}$$



Connections between t-SNE and Laplacians

t-SNE uses a gradient descent algorithm to compute a solution to this non-convex problem.

- Can show: [Steinerberger, Clustering with t-SNE, provably]
 For large parameters in the algorithm that computes the t-SNE solution, one obtains Laplacian Eigenmaps
- Observation: Laplacian Eigenmaps give good initialization of t-SNE algorithm

Questions:

- Can we modify Laplacian Eigenmaps by somehow adding "repulsion" to obtain better embedding in ℝ² without taking the t-SNE detour?
- Discard Kullback-Leibler interpretation of t-SNE, instead think of it as particle systems. Adopting this viewpoint, can we design good visualization methods based on particle systems & dynamical systems?



Conclusion and Outlook

- First meaningful theory for performance of spectral clustering
- Provides deterministic bounds for optimal graph cuts

Conclusion and Outlook

- First meaningful theory for performance of spectral clustering
- Provides deterministic bounds for optimal graph cuts
- Solved: Strong consistency of Laplacians for community detection [Shaofeng Deng, Oct.2019]
- Open problem: optimal estimates for random data models
- Open problem: theory for semisupervised clustering
- Open problem: how to modify Laplacian eigenmap to obtain better visualization tool

Conclusion and Outlook

- First meaningful theory for performance of spectral clustering
- Provides deterministic bounds for optimal graph cuts
- Solved: Strong consistency of Laplacians for community detection [Shaofeng Deng, Oct.2019]
- Open problem: optimal estimates for random data models
- Open problem: theory for semisupervised clustering
- Open problem: how to modify Laplacian eigenmap to obtain better visualization tool

S.Ling and T.Strohmer. Certifying Global Optimality of Graph Cuts via Semidefinite Relaxation: A Performance Guarantee for Spectral Clustering. *Foundations of Comp. Math.*, to appear. [and on the arxiv]

Interested in data science? Check out **cedar.ucdavis.edu** Center for Data Science and Artificial Intelligence Research

