

SPECTRAL CLUSTERING

Let x_1, \dots, x_n be objects we wish to cluster.

Spectral clustering refers to the following:

- ①
- form a _____
 - construct an $n \times n$ matrix, called the _____, from this graph.
 - infer the clusters from the _____ or _____ decomposition of this matrix.

The ingredients that determine a sim. graph are

-
-

Given these two ingredients, the weighted adjacency matrix is

$$w_{ij} = \begin{cases} s_{ij} & \text{if } x_i, x_j \text{ adjacent} \\ 0 & \text{otherwise} \end{cases}$$

Examples

▷ ϵ -neighborhood graph

- x_i, x_j adjacent $\iff \|x_i - x_j\| \leq \epsilon$

- $s_{ij} = 1$

\implies locality captured entirely by graph

▷ complete (fully connected) graph

- all x_i, x_j adjacent

- $s_{ij} = \exp\left\{-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right\}, \sigma^2 > 0$

⇒ locality captured entirely by similarities

▷ l-nearest neighbor graph

- x_i, x_j adjacent $\iff x_i$ is an l-nearest neighbor of x_j or vice-versa.

- $s_{ij} = \exp\left\{-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right\}$

▷ mutual l -nearest neighbor graph

- x_i, x_j adjacent $\iff x_i$ is an l -nearest neighbor of x_j and vice versa

Unfortunately, these similarity graphs have free parameters (ϵ, σ, l) that can be difficult to tune.

The reason for constructing a similarity graph is that it reduces the problem of clustering to _____ :•

(c)

Find a partition of the graph such that

- edges between clusters have _____ weights
- edges within clusters have _____ weights

Graph Laplacians

Definitions

- The degree of a vertex x_i is

① $d_i :=$

- The degree matrix is the diagonal matrix

$$D :=$$

- The unnormalized graph Laplacian is

$$L :=$$

Note

L is independent of the self-similarity weights w_{ii} , because

$$L_{ii}$$

Proposition | (Properties of L)

1) For every $f \in \mathbb{R}^n$

$$f^T L f = \frac{1}{2} \sum_{i,j=1}^n w_{ij} (f_i - f_j)^2$$

2) L is symmetric and positive semi-definite

3) The smallest eigenvalue of L is 0

with corresponding eigenvector $\underline{1} = [1 \ 1 \ \dots \ 1]^T$

Hence L has n non-negative, real-valued eigenvalues $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$

Proof |

$$1) \quad f^T L f = f^T D f - f^T W f$$

$$= \sum_{i=1}^n d_i f_i^2 - \sum_{i,j=1}^n w_{ij} f_i f_j$$

$$= \frac{1}{2} \left(\sum_i d_i f_i^2 - 2 \sum_{i,j} w_{ij} f_i f_j + \sum_j d_j f_j^2 \right)$$

$$= \frac{1}{2} \sum_i \sum_j w_{ij} (f_i^2 - 2 f_i f_j + f_j^2)$$

$$= \frac{1}{2} \sum_{i,j} w_{ij} (f_i - f_j)^2$$

2) Symmetry follows from symmetry of W

PSD follows from 1) because, for any $f \in \mathbb{R}^n$,

$$\textcircled{E} \quad f^T L f =$$
$$\geq$$

3)

$$L \underline{1} = \begin{bmatrix} d_1 & & & \\ & d_2 & & \\ & & \ddots & \\ & & & d_n \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} - \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1n} \\ w_{21} & w_{22} & \dots & w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n1} & w_{n2} & \dots & w_{nn} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

=

=

The unnormalized graph Laplacian encodes many properties of the graph. The following is one such property that is relevant for clustering.

Notation | Let $A \subseteq \{x_1, \dots, x_n\}$ be a cluster.

Define the indicator vector

$$\underline{1}_A = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_n \end{bmatrix} \in \mathbb{R}^n$$

where $f_i = 1$ if $x_i \in A$ and $f_i = 0$ if $x_i \notin A$.

Proposition | Suppose the graph has connected

(F) components A_1, \dots, A_k . Then the _____
of L has dimension _____ and is
spanned by

Recall | The nullspace of L is

$$N(L) := \{f : Lf = 0\}$$

It suffices to show:

- $\underline{1}_{A_k} \in N(L)$ for each $k=1, \dots, K$
- If $f \in N(L)$, then

$$f =$$

Ⓒ

for some $\alpha_1, \dots, \alpha_k \in \mathbb{R}$.

Since $\underline{1}_{A_1}, \dots, \underline{1}_{A_K}$ are linearly independent,

it follows that

$$\dim(N(L)) =$$

Proof |

• $K=1$: we have already seen $\underline{1} \in N(L)$.

We need to show that if $Lf = 0$,

then

$$f = \alpha \cdot \underline{1}$$

for some $\alpha \in \mathbb{R}$.

If $Lf = 0$, then

$$0 = f^T L f = \sum_{i,j=1}^n w_{ij} (f_i - f_j)^2$$

If x_i, x_j are adjacent, then

$$w_{ij} > 0 \implies f_i = f_j$$

Since $K=1$, any two points can be connected by a path

$$\implies f_i = \text{constant}$$

Spectral Clustering

How can we use this result to devise a clustering algorithm?

Ideal case

In the ideal case where there are K connected components and K is known, we could proceed as follows:

- compute L
- compute a basis u_1, \dots, u_K for $N(L)$.

$$\hookrightarrow \text{span} \{u_1, \dots, u_K\} = \text{span} \{\mathbf{1}_{A_1}, \dots, \mathbf{1}_{A_K}\}$$

- define

$$y_i = (u_{i1}, \dots, u_{iK})$$

- if x_i, x_j in same component then

In conclusion, we can use the spectral decomposition of the graph Laplacian L to determine the connected components.

However, there are two problems with this approach

- 1) There are simpler ways of finding the connected components of a graph
- 2) It is very difficult to construct a similarity graph such that its connected components are clusters.

Nonideal case

More realistically, the components of our similarity graph will not coincide with clusters.

In practice, we should choose the similarity graph such that

number of connected components \ll desired number of clusters

Then there will be edges connecting points in different clusters, but the weights

⑤ between clusters should be --- --- the weights within clusters.

Then, if points are indexed according to cluster, we have

$$L = \begin{bmatrix} \boxed{\text{large}} & & \\ & \boxed{\text{large}} & \\ \text{small} & & \boxed{\text{large}} \end{bmatrix} = \begin{bmatrix} \boxed{\text{large}} & & 0 \\ & \boxed{\text{large}} & \\ 0 & & \boxed{\text{large}} \end{bmatrix} + \begin{bmatrix} & & \\ & & \text{small} \\ & & \end{bmatrix}$$

\uparrow \uparrow

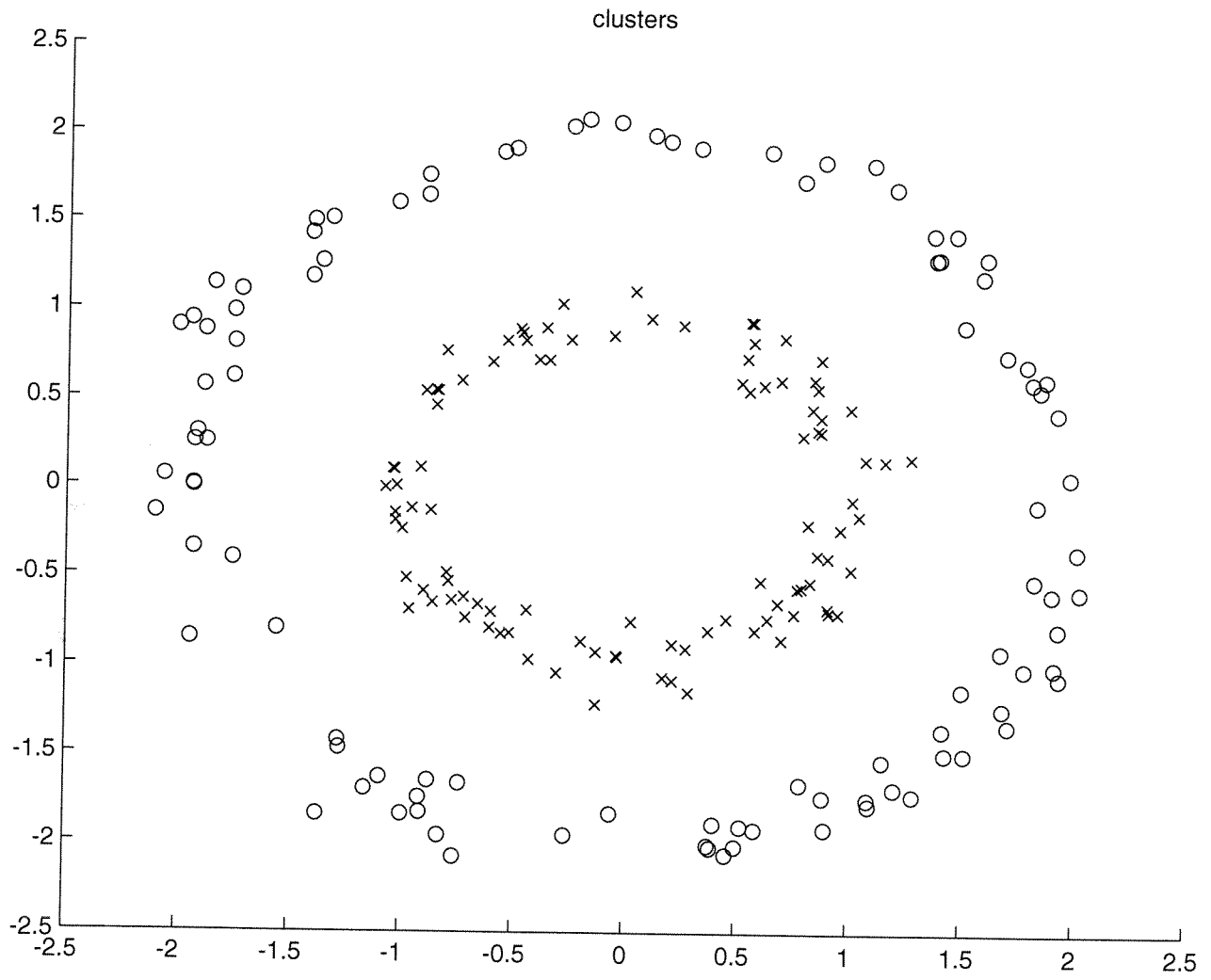
Perturbation Theory

Perturbation theory establishes results that show: if we perturb a matrix by another matrix with small entries, then the eigenvectors and eigenvalues of the matrix are perturbed by a correspondingly small amount.

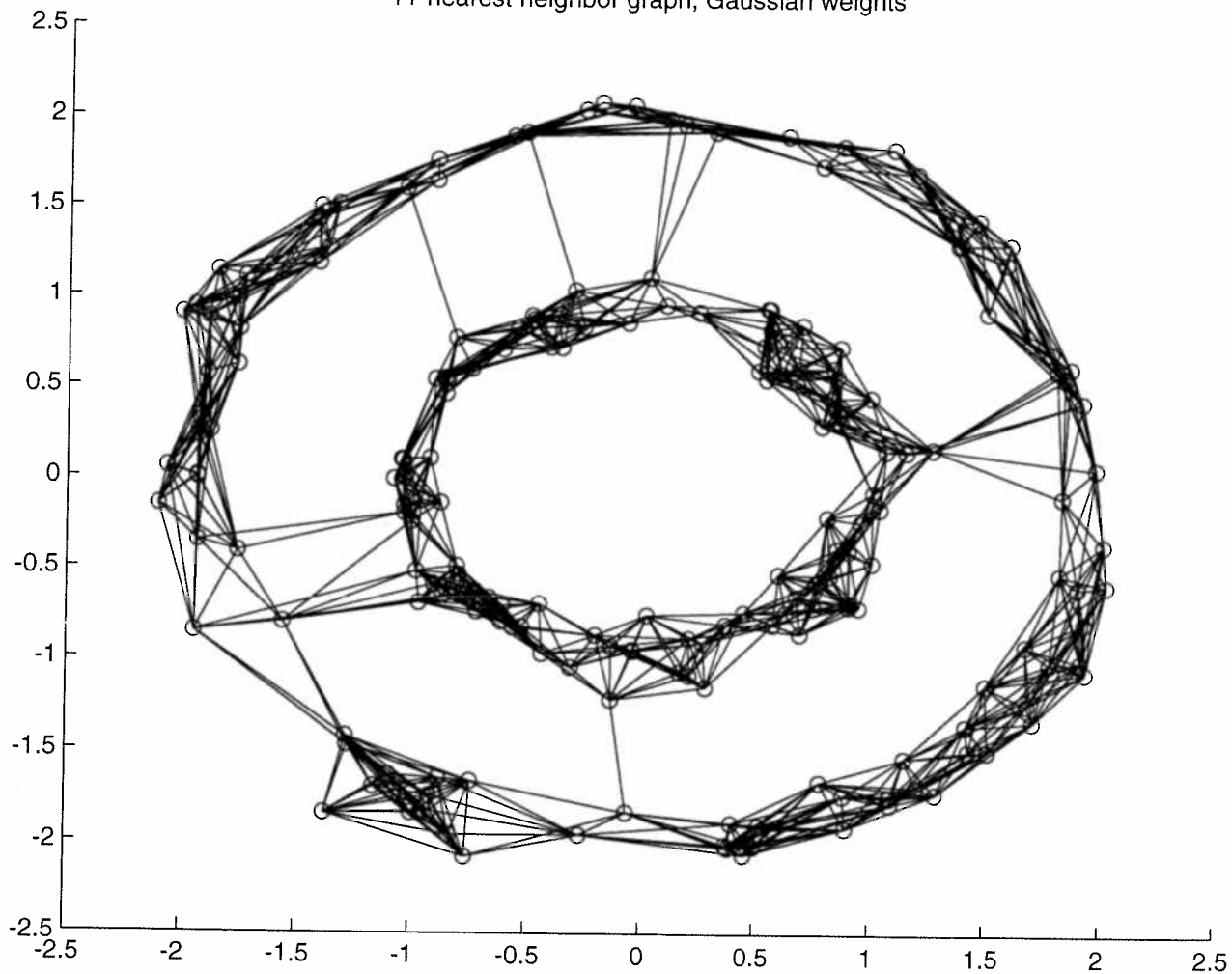
Recall the nullspace of L is also the set of eigenvectors with eigenvalue 0. We can use the eigenvectors of L with the smallest K eigenvalues as an approximation to the nullspace of an idealized L based on the true clusters.

SPECTRAL CLUSTERING

- Construct similarity graph
- Form graph Laplacian $L \in \mathbb{R}^{n \times n}$
- Determine the K smallest eigenvalues of L ,
 $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_K$, and corresponding
eigenvectors $u_1, \dots, u_K \in \mathbb{R}^n$
- Define $y_i = (u_{1i}, u_{2i}, \dots, u_{Ki})$, $i = 1, \dots, n$
- Cluster $\{y_i\}_{i=1}^n$ using K -means clustering
and assign $\{x_i\}_{i=1}^n$ to corresponding clusters

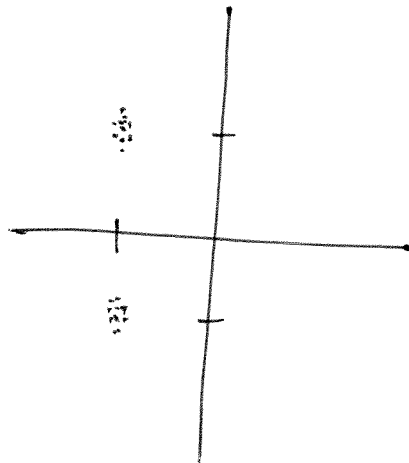


11-nearest neighbor graph, Gaussian weights



-0.070711 -0.07071
-0.070711 -0.07071
-0.070711 -0.070709
-0.070711 -0.07071
-0.070711 -0.070708
-0.070711 -0.070708
-0.070711 -0.070709
-0.070711 0.071153
-0.070711 0.071682
-0.070711 0.070093
-0.070711 0.071059
-0.070711 -0.070708
-0.070711 -0.070709
-0.070711 0.070098
-0.070711 -0.07071
-0.070711 0.071489
-0.070711 -0.070709
-0.070711 -0.07071
-0.070711 0.070098
-0.070711 -0.070708
-0.070711 0.070098
-0.070711 -0.07071
-0.070711 0.071682
-0.070711 -0.070709
-0.070711 -0.07071
-0.070711 0.070098
-0.070711 -0.07071
-0.070711 0.071489
-0.070711 -0.07071
-0.070711 -0.070708
-0.070711 -0.070708
-0.070711 -0.07071
-0.070711 -0.070708
-0.070711 0.071682
-0.070711 0.070098
-0.070711 0.070095
-0.070711 -0.07071
-0.070711 -0.070708
-0.070711 -0.07071
-0.070711 0.070093
-0.070711 0.071682
-0.070711 -0.070709
-0.070711 0.070096
-0.070711 0.071153
-0.070711 0.071153
-0.070711 0.070099
-0.070711 -0.07071
-0.070711 0.070093
-0.070711 0.070093
-0.070711 0.070098

a few randomly selected y_i 's



Normalized Spectral Clustering

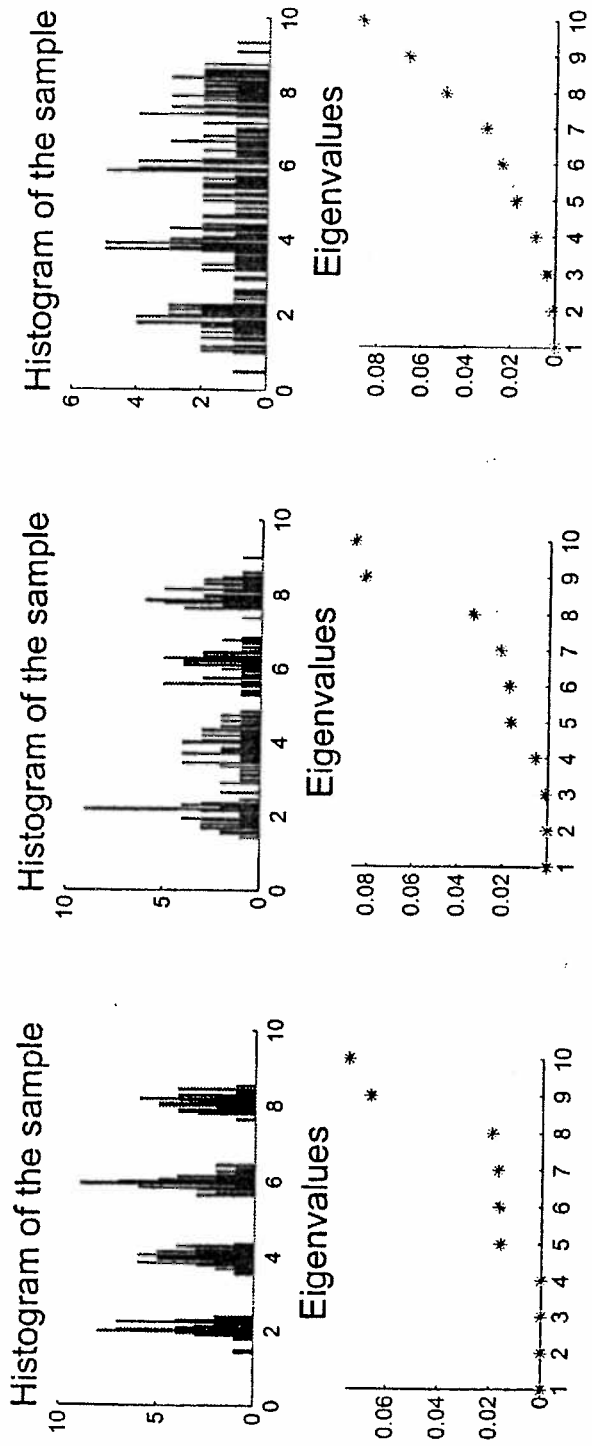
(K) Define the _____

$$\tilde{L} :=$$

Proposition 1

- 1) \tilde{L} is positive semidefinite with real, nonnegative eigenvalues $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$
- 2) 0 is an eigenvalue of \tilde{L} with corresponding eigenvector $\underline{1} = [1 \ 1 \ \dots \ 1]^T$.
- 3) Suppose the graph has connected components A_1, \dots, A_K . Then the nullspace of \tilde{L} has dimension K and is spanned by $\underline{1}_{A_1}, \dots, \underline{1}_{A_K}$

\implies we can substitute \tilde{L} for L and get another spectral clustering algorithm.



The first large gap in the spectrum can be used to infer the number of clusters automatically. Results shown are for normalized spectral clustering.

Graph Cuts

Spectral clustering attempts to find a partition A_1, \dots, A_K of the similarity graph such that

- (L)
- w_{ij} large if x_i, x_j in _____ cluster
 - w_{ij} small if x_i, x_j in _____ clusters

using properties of graph Laplacians.

It is also possible to approach this goal directly.

The _____ problem is to minimize

$$\text{cut}(A_1, \dots, A_K) :=$$

with respect to A_1, \dots, A_K , where

$$W(A, B) :=$$

The min cut problem is efficiently solvable when $K=2$.
Unfortunately, it tends to lead to small
(and often singleton) clusters.

Therefore researchers have examined modified
criteria that favor larger clusters:

- (M)
- Ratio Cut (Hagen and Kahng, 1992)

$$\text{Ratio Cut}(A_1, \dots, A_K) :=$$

where $|A| :=$

- Normalized Cut (Shi and Malik, 2000)

$$\text{Ncut}(A_1, \dots, A_K) :=$$

=

where $\text{vol}(A_i) :=$

Unfortunately, minimizing these criteria is NP hard.

(N) Remarkably, however, spectral clustering may be used to solve _____ of these problems

In particular,

• unnormalized spectral clustering \implies

• normalized spectral clustering \implies

Approximating Ratio Cut for $K=2$

We wish to solve

$$\min_A \text{Ratio Cut}(A, \bar{A}) = \min_A \left[\frac{\text{cut}(A, \bar{A})}{|A|} + \frac{\text{cut}(A, \bar{A})}{|\bar{A}|} \right]$$

Given a subset $A \subseteq \{1, 2, \dots, n\}$, define $f_A := (f_{A_1}, \dots, f_{A_n})^T$ by

$$f_{A_i} := \begin{cases} \sqrt{|\bar{A}|/|A|} & \text{if } i \in A \\ -\sqrt{|A|/|\bar{A}|} & \text{if } i \notin A. \end{cases}$$

Claim $f_A^T L f_A = n \text{Ratio Cut}(A, \bar{A})$

Proof

$$f_A^T L f_A = \frac{1}{2} \sum_{ij} w_{ij} (f_{A_i} - f_{A_j})^2$$

$$= \frac{1}{2} \sum_{i \in A, j \in \bar{A}} w_{ij} \left(\sqrt{\frac{|\bar{A}|}{|A|}} + \sqrt{\frac{|A|}{|\bar{A}|}} \right)^2 + \frac{1}{2} \sum_{i \in \bar{A}, j \in A} w_{ij} \left(-\sqrt{\frac{|A|}{|\bar{A}|}} - \sqrt{\frac{|\bar{A}|}{|A|}} \right)^2$$

$$= \frac{1}{2} \text{cut}(A, \bar{A}) \left(\frac{|\bar{A}|}{|A|} + \frac{|A|}{|\bar{A}|} + 2 \right) + \frac{1}{2} \text{cut}(A, \bar{A}) \left(\frac{|A|}{|\bar{A}|} + \frac{|\bar{A}|}{|A|} + 2 \right)$$

$$= \text{cut}(A, \bar{A}) \left(\frac{|\bar{A}|}{|A|} + \frac{|A|}{|\bar{A}|} + 2 \right)$$

$$= \text{cut}(A, \bar{A}) \left(\frac{|A| + |\bar{A}|}{|A|} + \frac{|\bar{A}| + |A|}{|\bar{A}|} \right)$$

$$= \underbrace{(|A| + |\bar{A}|)}_n \cdot \underbrace{\left[\frac{\text{cut}(A, \bar{A})}{|A|} + \frac{\text{cut}(\bar{A}, A)}{|\bar{A}|} \right]}_{\text{Ratio cut}(A, \bar{A})}$$

Ratio cut (A, \bar{A})

Furthermore, f_A satisfies

$$\textcircled{0} \quad \bullet \quad \underline{1}^T f_A = \sum_{i=1}^n f_{A_i}$$

=

=

$$\bullet \quad \|f_A\|^2 = \sum_{i=1}^n f_{A_i}^2$$

=

=

Therefore, the RatioCut problem can be written as the following optimization problem:

$$\begin{aligned} \min_A \quad & f_A^T L f_A \\ \text{s.t.} \quad & \underline{1}^T f_A = \\ & \|f_A\| = \end{aligned}$$

If we allow $f \in \mathbb{R}^n$ we have the following relaxation :

$$\begin{aligned} \min_{f \in \mathbb{R}^n} \quad & f^T L f \\ \text{s.t.} \quad & \underline{1}^T f = 0 \\ & \|f\| = \sqrt{n} \end{aligned} \quad \leftarrow \begin{cases} \text{No longer} \\ \text{require } f = f_A \\ \text{for some } A \end{cases}$$

The solution is

① $f =$

To recover a solution to the original discrete problem, we can use K-means, $K=2$, to cluster the vectors

$$y_i := (1 \quad f_i)$$

Therefore, the approximate solution is given by unnormalized spectral clustering.

A similar analysis applies to $K > 2$ and to Ncut.

Final comments

- The k -nn graph with Gaussian kernel similarity is most common, although the choice of similarity graph is largely an art.
- Another normalized graph Laplacian is

$$\tilde{L} := D^{-\frac{1}{2}} L D^{-\frac{1}{2}}$$

(Ng, Jordan, + Weiss, 2002). It can also be used for spectral clustering but the algorithm requires some modification.

- Which method is preferred? \tilde{L} is recommended by

U. von Luxburg, "A Tutorial on Spectral Clustering," 2007.

- b -matching is an interesting alternative to nearest neighbor graphs: it ensures that each node has the same number of incident edges (unweighted degree). See Tebara et al, ICML 2009

graph Laplacians are also used in
② _____ learning.

For example, suppose we are in a regression setting with data

$$(x_i, y_i)_{i=1}^m$$

$$(x_i)_{i=m+1}^n$$

To make sure the estimated function doesn't "wobble" too much, we could minimize

$$\frac{1}{m} \sum_{i=1}^m (y_i - f(x_i))^2 + \underbrace{\frac{1}{2} \sum_{i,j=1}^n w_{ij} (f(x_i) - f(x_j))^2}_{f^T L f}$$

Key

A. similarity graph, graph Laplacian, eigenvalue/spectral

B. a graph, a similarity matrix $S = [s_{ij}]_{i,j=1}^n$

C. graph partitioning, low, high

D. $d_i = \sum_{j=1}^n w_{ij}$, $D = \begin{bmatrix} d_1 & & 0 \\ & d_2 & \\ 0 & & \dots & \\ & & & d_n \end{bmatrix}$, $L = D - W$

$$L_{ii} = d_i - w_{ii} = \sum_{j \neq i} w_{ij}$$

$$E. f^T L f = \sum w_{ij} (f_i - f_j)^2 \geq 0$$

$$L \underline{1} = \begin{bmatrix} d_1 \\ \vdots \\ d_n \end{bmatrix} - \begin{bmatrix} \sum w_{ij} \\ \vdots \\ \sum w_{nj} \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} = 0 \cdot \underline{1}$$

L is PSD

F. nullspace (0-eigenspace); K ; $\underline{1}_{A_1}, \dots, \underline{1}_{A_K}$

G. $f = \sum_{k=1}^K \alpha_k \underline{1}_{A_k}$, $\dim(N(L)) = K =$ multiplicity of 0 as an eigenvalue

$$H. y_i = y_j$$
$$I. \underline{1}_{A_1} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}, \underline{1}_{A_2} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \\ y_6 \\ y_7 \\ y_8 \end{bmatrix} = \begin{bmatrix} \alpha_1 & \alpha_2 \\ \alpha_1 & \alpha_2 \\ \beta_1 & \beta_2 \\ \beta_1 & \beta_2 \\ \alpha_1 & \alpha_2 \\ \beta_1 & \beta_2 \\ \alpha_1 & \alpha_2 \\ \alpha_1 & \alpha_2 \end{bmatrix}$$

J. less than, ideal case, noise/perturbation

K. normalized graph Laplacian $\tilde{L} = D^{-1} L = I - D^{-1} W$

L. same, different, mincut, $\text{cut}(A_1, \dots, A_k) = \frac{1}{2} \sum_{k=1}^K W(A_k, \bar{A}_k)$, $W(A, B) = \sum_{i \in A} \sum_{j \in B} w_{ij}$

M. Ratio Cut $(A_1, \dots, A_k) = \frac{1}{2} \sum_{k=1}^K \frac{W(A_k, \bar{A}_k)}{|A_k|} = \sum_{k=1}^K \frac{\text{cut}(A_k, \bar{A}_k)}{|A_k|}$

$|A_k| = \#$ of nodes in A

$N\text{cut}(A_1, \dots, A_k) = \frac{1}{2} \sum_{k=1}^K \frac{W(A_k, \bar{A}_k)}{\text{vol}(A_k)} = \sum_{k=1}^K \frac{\text{cut}(A_k, \bar{A}_k)}{\text{vol}(A_k)}$

$\text{vol}(A_k) = \sum_{i \in A} \sum_{j \in A} w_{ij}$

N. relaxations, Ratio Cut, Ncut

O. $\mathbb{1}^T f_A = \sum_{i \in A} \sqrt{\frac{|\bar{A}|}{|A|}} - \sum_{i \in \bar{A}} \sqrt{\frac{|A|}{|\bar{A}|}} = \sqrt{|A| \cdot |\bar{A}|} - \sqrt{|A| \cdot |\bar{A}|} = 0$

$\|f_A\|^2 = \sum_{i \in A} \frac{|\bar{A}|}{|A|} + \sum_{i \in \bar{A}} \frac{|A|}{|\bar{A}|} = |\bar{A}| + |A| = n$

discrete, 0, \sqrt{n}

P. eigenvector of L corresponding to second smallest eigenvalue

Q. semi-supervised