

Recursive Clustering

Sunday, September 26, 2021 6:11 PM

Julian

We can view the elements of a set as vertices of a graph whose edge weights represent similarities.

$$G = (V, E) \quad (A)_{ij} = a_{ij} = \text{sim}(i, j) \geq 0$$

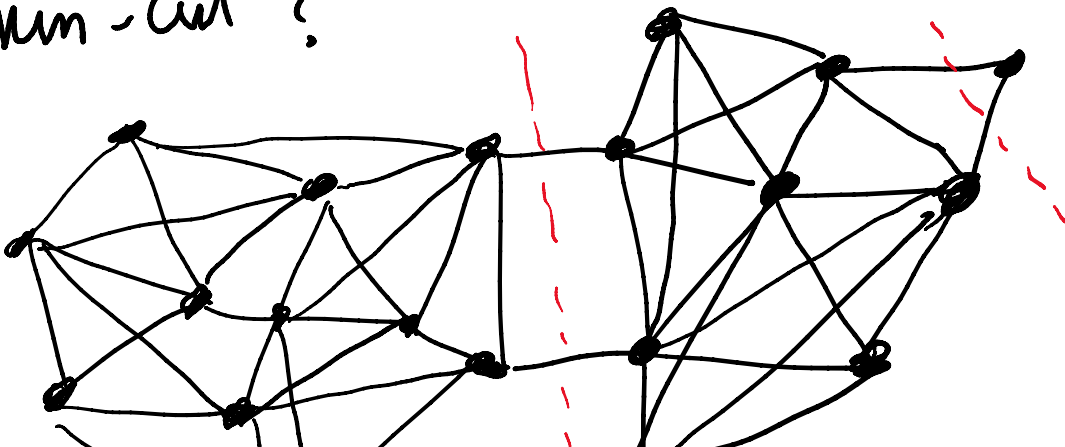
How to split into clusters?

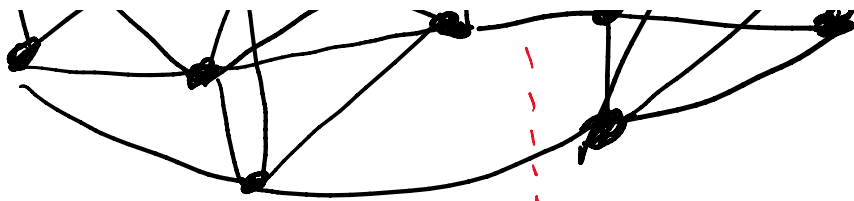
Goal: want to keep similar elements in same cluster; and dissimilar elements in different clusters.

Q.1. How to cut?

Q.2. What is the quality of a cluster?

min-cut?





or is this better \uparrow

$$a(S) = \sum_{i \in S, j \in \bar{S}} a_{ij}$$

Conductance (expansion if unweighted)

$$\phi(S) = \frac{\sum_{i \in S, j \notin S} a_{ij}}{a(S)}$$

Similarity across (S, \bar{S})

$$S \subseteq V$$

$$\min_{S \subseteq V} a(S), a(V \setminus S)$$

Similarity incident to S or \bar{S} .

$$\phi(G) = \min_{S \subseteq V} \phi(S).$$

Since $a_{ij} \geq 0$, we can normalize rows to get B :

$$b_{ij} = \frac{a_{ij}}{\sum_j a_{ij}} = \frac{a_{ij}}{a(i)}$$

$$\text{Then } B \mathbf{1} = \mathbf{1} \quad \text{and} \quad B^T \boldsymbol{\pi} = \boldsymbol{\pi}$$

$$\text{where } \pi(i) \propto a(i)$$

\\ \\ \\

where $\|(\cdot)\| = 1$

$$\left(\begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \end{array} \right) \rightarrow \left(\begin{array}{c} | \\ | \\ | \end{array} \right) \left(\begin{array}{c} \pi \\ \pi \\ \pi \end{array} \right)$$

$$\begin{aligned} (B^T \pi)_i &= \sum_j B_{ji} \pi_j = \sum_j \frac{a_{ij}}{a_{(j)}} \cdot \frac{a_{(j)}}{\|\pi\|} \\ &= \frac{a_{(i)}}{\|\pi\|} = \pi_i \end{aligned}$$

Note: $\pi_i b_{ij} = \pi_j b_{ji}$ (time reversible)

B^T represents the transition matrix of a time-reversible Markov chain.

Thm [Cheeger; JS] Let $B \geq 0$ $B \in \mathbb{R}^{N \times N}$

$$B \mathbf{1} = \mathbf{1}, \quad \pi_i b_{ij} = \pi_j b_{ji} \quad \pi_i > 0.$$

Let v be the second (right) eigenvector of B with components $v_{i_1} \geq v_{i_2} \geq \dots \geq v_{i_N}$ and eigenvalue λ_2 . ($\lambda = 1$). Then,

$$1 - \lambda_2 > 1 / \left(\min \phi(\xi_1, \xi_2, \dots, \xi_N) \right)^2.$$

$$2. \min_{S \subseteq V} \phi(S) \geq 1 - \lambda_2 \geq \frac{1}{2} \left(\min_{1 \leq l < N} \phi(\{1, 2, \dots, l\}) \right).$$

Corollary: \exists polytime algorithm to find S

A-t. $\phi(S) \leq \sqrt{2(1-\lambda_2)} \leq 2\sqrt{\text{OPT}}.$

Algorithm:

- 1) compute second eigenvector v
- 2) sort components $v_1 \geq v_2 \geq \dots \geq v_N$
- 3) Output min conductance cut among $\{v_1, \dots, v_i\} / \{v_{i+1}, \dots, v_N\}.$

Pf. (of Thm). B is not symmetric in general.

let D be a diagonal nonnegative matrix

s.t. $D^2 = \Pi.$

Claim: $Q = D B D^{-1}$ is symmetric and has same eigenvalues as $B.$

-1 $\sim -1, T, \dots$

eigenvalues

Note $D^2 B = B^T D^2 \Rightarrow DBD^{-1} = D^{-1} B^T D$
 symmetric.

Next if v $Bv = \lambda v$
 $DBD^{-1}(Dv) = \lambda(Dv)$ same eigenvalues.

We also know $\lambda_1 = 1$ and $\mathbb{1}$ is the eigenvector of B .

$(\pi^T D^{-1}) Q = \pi^T B D^{-1} = \pi^T D^{-1} \leftarrow$ e.v. of Q .

So $\lambda_2 = \max_{x: \pi^T D^{-1} x = 0} \frac{x^T DBD^{-1} x}{x^T x}$

$1 - \lambda_2 = \min \frac{x^T (I - DBD^{-1}) x}{x^T x} = \frac{x^T D (I - B) D^{-1} x}{x^T x}$

let $y = D^{-1} x$

$= \min_{y: \pi^T y = 0} \frac{y^T D^2 (I - B) y}{y^T D^2 y}$

Numerator = $y^T D^2 (I - B) y = -\sum_{i \neq j} \pi_i b_{ij} y_i y_j + \sum \pi_i (1 - b_{ii}) y_i^2$

$$\begin{aligned}
& \text{(F)} + \sum_i \pi_i (1 - b_{ii}) y_i \\
&= - \sum_{i \neq j} \pi_i b_{ij} y_i y_j + \sum_{i \neq j} \pi_i b_{ij} \frac{(y_i^2 + y_j^2)}{2} \\
&= \sum_{i < j} \pi_i b_{ij} (y_i - y_j)^2 = \mathcal{E}(y, y) \quad (\text{say})
\end{aligned}$$

$$1 - \lambda_2 = \frac{\mathcal{E}(y, y)}{\sum_i \pi_i y_i^2}$$

For the first inequality, let (S, \bar{S}) be min conductance cut. Define

$$w_i = \begin{cases} \sqrt{\frac{\pi(\bar{S})}{\pi(S)}} & i \in S \\ -\sqrt{\frac{\pi(S)}{\pi(\bar{S})}} & i \notin S. \end{cases}$$

Then $\pi^T w = 0$

$$\begin{aligned}
\text{and } 1 - \lambda_2 &\leq \frac{\mathcal{E}(y, y)}{\sum_i \pi_i y_i^2} = \frac{\sum_{i \in S, j \notin S} \pi_i b_{ij} \left(\sqrt{\frac{\pi(\bar{S})}{\pi(S)}} + \sqrt{\frac{\pi(S)}{\pi(\bar{S})}} \right)^2}{\sum_{i \in S} \pi_i \frac{\pi(\bar{S})}{\pi(S)} + \sum_{i \notin S} \pi_i \frac{\pi(S)}{\pi(\bar{S})}} \\
&= \leq \pi_i b_{ij} (\pi(\bar{S}) + \pi(S))^2
\end{aligned}$$

$$= \frac{\sum_{i \in S, j \notin S} \pi_i b_{ij} (\pi(S) + \pi(\bar{S}))}{\pi(S) \cdot \pi(\bar{S})}$$

$$\leq \frac{2 \sum_{i \in S, j \notin S} \pi_i b_{ij}}{\min\{\pi(S), \pi(\bar{S})\}} = 2 \phi(S)$$

$$v_1 \geq v_2 \geq \dots \geq v_N$$

Let r be s.t.

$$\pi_1 + \pi_2 + \dots + \pi_{r-1} \leq \frac{1}{2} < \pi_1 + \dots + \pi_r$$

and $z_i = v_i - v_r$ so that

$$z_1 \geq z_2 \geq \dots \geq z_r = 0 \geq z_{r+1} \geq \dots \geq z_N$$

$$\frac{E(v, v)}{\sum_i \pi_i v_i^2} = \frac{E(z, z)}{-v_r^2 + \sum_i \pi_i z_i^2} \geq \frac{E(z, z)}{\sum_i \pi_i z_i^2}$$

$$\begin{aligned} \sum_i \pi_i v_i^2 &= \sum_i \pi_i (z_i + v_r)^2 = \sum_i \pi_i z_i^2 + (\sum_i \pi_i) \cdot v_r^2 \\ &\quad + \sum_i 2\pi_i z_i v_r \\ &= \sum_i \pi_i z_i^2 + v_r^2 - 2v_r^2 \end{aligned}$$

$$= \sum \pi_i z_i^2 - V_r^2 - 2V_r$$

$$\pi^T \bar{z} = \pi^T V - \pi_r^2 = -V_r^2 \quad = \sum \pi_i z_i^2 - V_r^2$$

$$\frac{\sum_{i < j} \pi_i b_{ij} (z_i - z_j)^2 \left(\sum_{i < j} \pi_i b_{ij} (|z_i| + |z_j|)^2 \right)}{\sum \pi_i z_i^2 \left(\sum_{i < j} \pi_i b_{ij} (|z_i| + |z_j|)^2 \right)^2}$$

$$\text{Num} \geq \left(\sum \pi_i b_{ij} |z_i - z_j| (|z_i| + |z_j|) \right)^2$$

$$\text{Claim} \quad |z_i - z_j| (|z_i| + |z_j|) \geq \sum_{k=i}^{j-1} |z_{k+1}^2 - z_k^2|$$

z_i, z_j have same sign

$$\text{LHS} = |z_i^2 - z_j^2| = \text{RHS.}$$

$$\text{else} \quad \text{LHS} = (|z_i| + |z_j|)^2 > z_i^2 + z_j^2 = \text{RHS.}$$

$$i \leq r \leq j$$

$$z_r = 0.$$

$$\text{So Numerator} \geq \left(\sum_{i < j} \pi_i b_{ij} \sum_{k=i}^{j-1} |z_{k+1}^2 - z_k^2| \right)^2$$

$$\text{Denominator} \leq \sum_i \pi_i z_i^2 \cdot 2 \sum_{i < j} \pi_i b_{ij} (z_i^2 + z_j^2)$$

$$\leq 2 \left(\sum \pi_i z_i^2 \right)^2$$

$$\leq 2 \left(\sum_i \pi_i z_i^2 \right)^2$$

$$1 - \lambda_2 \geq \frac{1}{2} \cdot \left(\frac{\sum_{i < j} \pi_i b_{ij} \sum_{k=i}^{j-1} |z_{k+1}^2 - z_k^2|}{\sum_i \pi_i z_i^2} \right)^2$$

Let $\hat{\alpha} = \min_{1 \leq k < N} \frac{\sum_{i \leq k < j} \pi_i b_{ij}}{\min \{ \pi(\{1, \dots, k\}), \pi(\{k+1, \dots, N\}) \}}$

Then

$$\sum_{i < j} \pi_i b_{ij} \sum_{k=i}^{j-1} |z_{k+1}^2 - z_k^2| = \sum_{k=1}^{N-1} |z_{k+1}^2 - z_k^2| \sum_{i \leq k < j} \pi_i b_{ij}$$

$$\geq \sum_{k=1}^{N-1} |z_{k+1}^2 - z_k^2| \cdot \hat{\alpha} \cdot \min \{ \pi(\{1, \dots, k\}), \pi(\{k+1, \dots, N\}) \}$$

$$= \hat{\alpha} \left(\sum_{k=1}^{N-1} (z_k^2 - z_{k+1}^2) \pi(S_k) + \sum_{k=1}^{N-1} (z_{k+1}^2 - z_k^2) (1 - \pi(S_k)) \right)$$

$$= \hat{\alpha} \left(\sum_{k=1}^{N-1} (z_k^2 - z_{k+1}^2) \pi(S_k) + z_N^2 - z_1^2 \right)$$

$$= \hat{\alpha} \left(\sum_{k=1}^{N-1} z_k^2 (\pi(S_k) - \pi(S_{k+1})) + z_N^2 \right)$$

$$= \hat{\alpha} \left(\sum_{k=1}^N \pi_k z_k^2 \right)$$

$$= \alpha \left(\sum_{k=1}^K \tau_k \right)$$

$$\therefore -\lambda_2 \geq \frac{1}{2} \hat{\alpha}^2$$

Multinway Cheeger

Recursive Partitioning.

— cut into two parts

— Recurse while cluster is not high enough quality.

α = conductance of cluster

ϵ = fraction of similarity between clusters.

Thm. $\exists (\alpha, \epsilon)$ - bicriteria clustering,

Recursive Spectral Partitioning finds a

$\left(c \frac{\alpha^2}{\log^2\left(\frac{n}{\epsilon}\right)}, c \sqrt{\epsilon \log \frac{n}{\epsilon}} \right)$ - clustering.

^ can be used with any

Can be used with any
approximate conductance cut algorithm.

Thm. \exists algorithm that finds S s.t.
 $\phi(S) \leq C\sqrt{\log n} \cdot \text{OPT}.$