

Clustering

Lecturer: Sushant Sachdeva

Scribe: Liwen Lu

1 Clustering of a Graph and Eigenvalue

Informally, “clustering” is grouping vertices such that there are more connectivity inside each group compare to between groups.

1.1 Conductance

Definition 1.1. Volume: let $S \subseteq V$, $\text{vol}(S) \stackrel{\text{def}}{=} \sum_{v \in S} d_v$ where $d_v = \deg(v) = \mathbf{D}_{v,v}$

Lemma 1.2. Define the indicator of set S as

$$\mathbf{1}_S(v) = \begin{cases} 1, & v \in S \\ 0, & \text{otherwise} \end{cases}$$

Then

$$\mathbf{1}_S^T \mathbf{D} \mathbf{1}_S = \text{vol}(S)$$

Definition 1.3. Define Boundary of S as

$$E(S, \bar{S}) = \{(u, v) | (u, v) \in E, u \in S, v \in V \setminus S\},$$

where \bar{S} is the complement set.

Definition 1.4. Define conductance of S measured in graph G to be

$$\phi_G(S) \stackrel{\text{def}}{=} \frac{|E(S, \bar{S})|}{\min\{\text{vol}(S), \text{vol}(V \setminus S)\}},$$

and define conductance of graph G to be

$$\phi(G) \stackrel{\text{def}}{=} \min_{\substack{S \subseteq V \\ S \neq \emptyset, V}} \phi_G(S).$$

Computing $\phi(G)$ is called “minimum conductance problem”, and it is famously NP hard. However, we can connect conductance and eigenvalues.

1.2 Normalized Laplacian Matrix

Let $\nu_2 = \min_{y^T \mathbf{D} \mathbf{1} = 0} \left(\frac{y^T \mathbf{L} y}{y^T \mathbf{D} y} \right)$.

Define $x = \mathbf{D}^{\frac{1}{2}} y$ (note $\mathbf{D}^{\frac{1}{2}}$ is replacing each of the diagonal entries in diagonal matrix \mathbf{D} with square root), we can write $y^T \mathbf{D} y = x^T x$.

Note we are assuming the graph G is connected, then diagonal of D is strictly non-negative, so the mapping between x and y is a bijection, which gives us $x = \mathbf{D}^{\frac{1}{2}}y \Leftrightarrow y = \mathbf{D}^{-\frac{1}{2}}x \Rightarrow y^\top \mathbf{D}y = x^\top x$. Then we can rewrite ν_2 as

$$\nu_2 = \min_{x^\top (\mathbf{D}^{\frac{1}{2}}\mathbf{1})=0} \frac{x^\top \mathbf{D}^{-\frac{1}{2}}\mathbf{L}\mathbf{D}^{-\frac{1}{2}}x}{x^\top x}.$$

Definition 1.5. Define Normalized Laplacian Matrix as

$$\mathbf{N} \stackrel{\text{def}}{=} \mathbf{D}^{-\frac{1}{2}}\mathbf{L}\mathbf{D}^{-\frac{1}{2}}$$

We also claim that:

$$\begin{aligned} \lambda_1(\mathbf{N}) &= 0 \\ \psi_1 &= \mathbf{D}^{\frac{1}{2}}\mathbf{1} \\ \nu_2 &= \min_{x^\top \psi_1=0} \frac{x^\top \mathbf{N}x}{x^\top x} \\ &= \lambda_2(\mathbf{N}) \end{aligned}$$

1.3 Cheeger's Inequality

Theorem 1.6.

$$\frac{\nu_2}{2} \leq \phi(G) \leq \sqrt{2\nu_2}, \nu_2 \leq 2$$

For example, if $\nu_2 = 0.01 \Rightarrow 0.005 \leq \phi(G) \leq 0.14$
We'll first prove the left hand side, $\frac{\nu_2}{2} \leq \phi(G)$.

Proof. By definition, we know $\exists S$ s.t. $\phi(G) = \phi_G(S) = \frac{|E(S, \bar{S})|}{\text{vol}(S)}$ and $\text{vol}(S) \leq \text{vol}(\bar{S})$.

Recall $\nu_2 = \min_{y^\top \mathbf{D}\mathbf{1}=0} \left(\frac{y^\top \mathbf{L}y}{y^\top \mathbf{D}y} \right)$ and $\mathbf{1}_S^\top \mathbf{D}\mathbf{1}_S = \text{vol}(S)$, we have

$$\mathbf{1}_S^\top \mathbf{L}\mathbf{1}_S = \sum_{(u,v) \in E} (\mathbf{1}_S(u) - \mathbf{1}_S(v))^2$$

note

$$(\mathbf{1}_S(u) - \mathbf{1}_S(v))^2 = \begin{cases} 1, & \mathbf{1}_S(u) \neq \mathbf{1}_S(v) \Leftrightarrow (u,v) \in E(S, \bar{S}) \\ 0, & \text{otherwise} \end{cases}$$

therefore

$$\mathbf{1}_S^\top \mathbf{L}\mathbf{1}_S = \sum_{(u,v) \in E} \mathbf{1}[(u,v) \in E(S, \bar{S})] = |E(S, \bar{S})|$$

If we set $y = \mathbf{1}_S$, it would minimize $\frac{y^\top \mathbf{L}y}{y^\top \mathbf{D}y}$, however it may not satisfy $y^\top \mathbf{D}\mathbf{1} = 0$. In order to make y satisfy $y^\top \mathbf{D}\mathbf{1} = 0$, let $y = \mathbf{1}_S + c\mathbf{1}$, note we still have $y^\top \mathbf{L}y = |E(S, \bar{S})|$.

$$y^\top \mathbf{D}\mathbf{1} = 0 \Leftrightarrow c = \frac{-\mathbf{1}_S^\top \mathbf{D}\mathbf{1}}{\mathbf{1}^\top \mathbf{D}\mathbf{1}} = -\frac{\text{vol}(S)}{\text{vol}(V)} \quad (\text{by previous lemma})$$

$$\begin{aligned}
y^\top \mathbf{D}y &= \mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S + 2c \mathbf{1}_S^\top \mathbf{D} \mathbf{1} + c^2 \mathbf{1}^\top \mathbf{D} \mathbf{1} \\
&= \mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S - \frac{(\mathbf{1}_S^\top \mathbf{D} \mathbf{1})^2}{\mathbf{1}^\top \mathbf{D} \mathbf{1}} \\
&= \text{vol}(S) - \frac{\text{vol}(S)^2}{\text{vol}(V)} \\
&= \text{vol}(S) \left(1 - \frac{\text{vol}(S)}{\text{vol}(V)} \right) \\
&\geq \frac{\text{vol}(S)}{2} \tag{since \text{vol}(S) \leq \text{vol}(\bar{S})}
\end{aligned}$$

Therefore

$$\begin{aligned}
\frac{y^\top \mathbf{L}y}{y^\top \mathbf{D}y} &= \frac{|E(S, \bar{S})|}{y^\top \mathbf{D}y} \\
&\leq \frac{2|E(S, \bar{S})|}{\text{vol}(S)} \\
&= 2\phi_G(S) \\
&= 2\phi(G) \tag{by assumption} \\
\Rightarrow \nu_2 &= \min_{y^\top \mathbf{D} \mathbf{1} = 0} \frac{y^\top \mathbf{L}y}{y^\top \mathbf{D}y} \leq 2\phi(G) \\
\Rightarrow \frac{\nu_2}{2} &\leq \phi(G)
\end{aligned}$$

□

Lemma 1.7. *Given y s.t. $y^\top \mathbf{D} \mathbf{1} = 0$, we can find a distribution on t with $S_t \subseteq V$ s.t. $\text{vol}(S_t) \leq \frac{\text{vol}(V)}{2}$, and*

$$\frac{\mathbb{E}_t |E(S_t, \bar{S}_t)|}{\mathbb{E}_t \text{vol}(S_t)} \leq \sqrt{\frac{2y^\top \mathbf{L}y}{y^\top \mathbf{D}y}}$$

Let t be independent choices, P_t be distribution of t , X_t, Y_t are variables depend on t with $Y_t \geq 0$.

We can prove that $\exists t$ s.t.

$$\frac{X_t}{Y_t} \leq \frac{\mathbb{E}_t X_t}{\mathbb{E}_t Y_t} = \frac{\sum_t P_t X_t}{\sum_t P_t Y_t}$$

The proof is left as exercise. This implies that $\exists t$, s.t.

$$\frac{|E(S_t, \bar{S}_t)|}{\text{vol}(S_t)} \leq \sqrt{\frac{2y^\top \mathbf{L}y}{y^\top \mathbf{D}y}}$$

Furthermore, if y was the minimizing vector for ν_2 , then

$$\phi_S(S_t) = \frac{|E(S_t, \bar{S}_t)|}{\text{vol}(S_t)} \leq \sqrt{2\nu_2}$$