

Lecture 22: Cheeger's inequality

We will prove Cheeger's inequality which relates the second eigenvalue of the Laplacian matrix to the graph expansion, and also discuss some recent generalizations.

Graph expansion

Cheeger's inequality will show that λ_2 is "small" if and only if the graph G is "close" to disconnected.

First, let us make precise what it means for a graph to be close to disconnected.

There are different definitions to measure how well a graph is connected.

The expansion of a graph is defined as $\Phi(G) := \min_{S \subseteq V, |S| \leq |V|/2} \Phi(S)$, where $\Phi(S) := \frac{|E(S, V \setminus S)|}{|S|}$,

the ratio of the number of edges cut to the number of vertices in the set.

The conductance of a graph is defined as $\phi(G) = \min_{S \subseteq V, \text{vol}(S) \leq |E|} \phi(S)$, where $\phi(S) := \frac{|E(S, V \setminus S)|}{\text{vol}(S)}$ and

$\text{vol}(S) := \sum_{v \in S} \deg(v)$, the ratio of the number of edges cut to the total degree in the set.

There is a related problem known as the uniform sparsest cut problem, whose objective is to find $S \subseteq V$

that minimizes $\phi'(S) = \frac{|E(S, V \setminus S)|}{|S||V \setminus S|}$, the ratio of the number of edges cut to the number of pairs cut.

These definitions are more or less equivalent when the graph is d -regular (i.e. $\Phi(G) = d\phi(G)$ and $\frac{n}{2}\phi(S) \leq \Phi(S) \leq n\phi(S)$).

In general graphs, we will relate the graph conductance to the second eigenvalue.

We say a graph is an expander graph if $\phi(G)$ is large, and we say S is a sparse cut if $\phi(S)$ is small. Notice that $0 \leq \phi(S) \leq 1$ for every $S \subseteq V$.

Both concepts are very useful. As we have seen, sparse expander graphs are "magical" and have algorithmic applications, and we will also see that they can be used in derandomization (next time).

Finding a sparse cut is useful in designing divide-and-conquer algorithms, and have applications in image segmentation, data clustering, community detection in social networks, VLSI design, etc.

The Spectral Partitioning Algorithm

A popular heuristic in finding a sparse cut in practice is the following spectral partitioning algorithm.

- ① Compute the second eigenvector x of \mathcal{L} (the eigenvector corresponding to the second largest eigenvalue)
- ② Sort the vertices so that $x_1 \geq x_2 \geq \dots \geq x_n$ (where $n = |V|$ is the number of vertices)
- ③ Let $S_i = \begin{cases} \{1, \dots, i\} & \text{if } i \leq n/2 \\ \{i+1, \dots, n\} & \text{if } i > n/2 \end{cases}$.

Return $\min_i \{\phi(S_i)\}$.

That's the algorithm.

First, there is an almost linear time algorithm (in terms of number of edges) to compute the ^{approximately} second eigenvector of the adjacency matrix. It is known as the "power method", which we won't discuss today. So, the whole algorithm can be implemented in near linear time, quite easily especially if you use some mathematical software (e.g. MATLAB). This is one reason that this heuristic is popular.

Another reason is that it performs very well in various applications, especially in image segmentation and clustering, and it was considered a breakthrough in image segmentation about 15 years ago.

The proof of Cheeger's inequality will provide some performance guarantee of this algorithm.

Normalized Matrices

To state Cheeger's inequality nicely, we will use the "normalized" Laplacian matrix, which allows us to remove the dependency on the maximum degree of the graph.

Given an adjacency matrix A , let $\mathcal{A} = D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$ be the normalized adjacency matrix, and let $\mathcal{L} = I - \mathcal{A}$ be the normalized Laplacian matrix, where D is the diagonal matrix whose i -th entry is the degree of vertex i . Note that $\mathcal{L} = I - \mathcal{A} = D^{-\frac{1}{2}} (D - A) D^{-\frac{1}{2}} = D^{-\frac{1}{2}} L D^{-\frac{1}{2}}$.

Let $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_n$ be the eigenvalues of \mathcal{A} , and let $\beta_1 \leq \beta_2 \leq \dots \leq \beta_n$ be the eigenvalues of \mathcal{L} .

Claim $1 = \alpha_1 \geq \alpha_n \geq -1$ and $0 = \beta_1 \leq \beta_n \leq 2$.

Proof We prove the result for normalized adjacency, and the result for normalized Laplacian follows easily.

Note that 0 is an eigenvalue for \mathcal{L} , as $\mathcal{L} (D^{\frac{1}{2}} \vec{1}) = (D^{-\frac{1}{2}} L D^{-\frac{1}{2}}) (D^{\frac{1}{2}} \vec{1}) = D^{-\frac{1}{2}} L \vec{1} = 0$

To prove $\beta_1 = 0$, we will show that \mathcal{L} is a positive semidefinite matrix.

To see it, observe that $x^T \mathcal{L} x = x^T D^{-\frac{1}{2}} L D^{-\frac{1}{2}} x = \sum_{i,j \in E} x_i^T D^{-\frac{1}{2}} L_e D^{-\frac{1}{2}} x = \sum_{i,j \in E} \left(\frac{x_i}{\sqrt{d_i}} - \frac{x_j}{\sqrt{d_j}} \right)^2 \geq 0$,

where $L_e = b e b^T$ that we defined last time.

This implies that $I - A \succeq 0$, and thus $\alpha_1 \leq 1$.

Also, we can write $x^T(I+A)x = x^T L x + 2x^T A x = \sum_{e=ij \in E} \left(\left(\frac{x_i}{\sqrt{d_i}} - \frac{x_j}{\sqrt{d_j}} \right)^2 + \frac{2x_i x_j}{\sqrt{d_i d_j}} \right) = \sum_{e=ij \in E} \left(\frac{x_i}{\sqrt{d_i}} + \frac{x_j}{\sqrt{d_j}} \right)^2 \geq 0$.

and this implies that $I+A \succeq 0$, and thus $\alpha_n \geq -1$, and thus $\beta_n = 1 - \alpha_n \leq 2$. \square

Cheeger's inequality: $\frac{1}{2} \lambda_2 \leq \phi(G) \leq \sqrt{2 \lambda_2}$, where λ_2 is the second eigenvalue of L .

For simplicity, we assume the graph is d -regular. The proof for the general case is similar;

see e.g. my previous course notes for the proof.

The first inequality is called the "easy" direction, and the second inequality is called the "hard" direction.

So, naturally we prove the easy direction first.

(assuming d -regular)
↓

One nice thing about the Laplacian matrix is that we know that the first eigenvector is the all-one vector,

(unlike for adjacency matrix), and by the characterization of λ_2 using Rayleigh quotient, we have

$$\lambda_2 = \min_{x \perp \vec{1}} \frac{x^T L x}{x^T x} = \min_{x \perp \vec{1}} \frac{x^T L x}{d x^T x} = \min_{x \perp \vec{1}} \frac{\sum_{i,j \in E} (x_i - x_j)^2}{d \sum_{i \in V} x_i^2}.$$

To upper bound λ_2 , we just need to find a vector $x \perp \vec{1}$ and compute its Rayleigh quotient.

To get some intuition, let say $\phi(G) = \phi(S)$ and $|S| = n/2$.

We consider the "binary" solution: $x_i = \begin{cases} +1 & \text{if } i \in S \\ -1 & \text{if } i \notin S \end{cases}$.

Since $|S| = n/2$, $\sum_{i \in V} x_i = 0$, and thus $x \perp \vec{1}$.

$$\text{Then } \lambda_2 \leq \frac{\sum_{i,j \in E} (x_i - x_j)^2}{d \sum_{i \in V} x_i^2} = \frac{4|E(S)|}{d|V|} = \frac{2|E(S)|}{d|S|} = 2\phi(S).$$

For general S , we consider the binary solution: $x_i = \begin{cases} +\frac{1}{|S|} & \text{if } i \in S \\ -\frac{1}{|V-S|} & \text{if } i \notin S \end{cases}$.

$$\text{Then } x \perp \vec{1}, \text{ and } \lambda_2 \leq \frac{\sum_{i,j \in E} (x_i - x_j)^2}{d \sum_{i \in V} x_i^2} = \frac{|E(S)| \cdot \left(\frac{1}{|S|} + \frac{1}{|V-S|} \right)^2}{d \left(|S| \cdot \frac{1}{|S|^2} + |V-S| \cdot \frac{1}{|V-S|^2} \right)} = \frac{|E(S)| \cdot |V|}{d \cdot |S| \cdot |V-S|} \leq 2\phi(S)$$

This proves the easy direction.

To summarize, if there is a sparse cut, then λ_2 is small.

A consequence is that if λ_2 is large, then we know that G has no sparse cut.

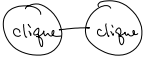
This direction is useful in deterministic construction of expander graphs.

The Hard Direction: Intuition

In the minimization problem $\min_{x \perp \vec{1}} \frac{\sum_{i,j \in E} (x_i - x_j)^2}{d \sum_{i \in V} x_i^2}$, if we can only search for "binary" solutions,

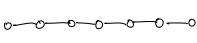
then we are essentially optimizing over the conductances.

Unfortunately, we are optimizing over a much larger domain (otherwise the problem is not efficiently solvable), and there could be some very non-binary solutions (very "smooth" vector), for which it is not clear how to find a sparse cut from it.

To get some feeling, suppose we are given a graph like .

Observe that the optimizer tries to minimize the average $(x_i - x_j)^2 / (x_i^2 + x_j^2)$.

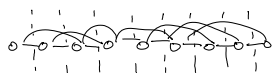
In this case, it is not good to "split" the vertices in a clique, because there are so many edges within it. So, we would expect that the values in each clique are very similar, while the two cliques would have different values so that $x \perp \vec{1}$. Hence, we expect that the minimizer would look very similar to a binary vector, and we can easily find a good cut with $\phi(S) \approx \lambda_2$.

Now, suppose we are given a graph like , then the minimizer can do much better by making each edge very short, while the values decrease smoothly from +1 to -1, in which case $\lambda_2 \ll \phi(G)$.

The key of Cheeger's inequality is to show that λ_2 cannot be much smaller than $\phi(G)$.

In other words, if λ_2 is small, then we can extract a somewhat sparse cut from the eigenvector.

We can think of the optimizer "embedding" the graph into a line, while most edges are short.



Then it should be the case that some threshold gives a sparse cut (e.g. row and column argument).

The Hard Direction : Proof

The first step is to preprocess the second eigenvector so that at most half the entries are nonzero.

This would guarantee that the output set S satisfies $|S| \leq |V|/2$.

This step is simple. Without loss of generality we assume there are fewer positive entries in x than negative entries.

∩ . . .

than negative entries.

Consider the following vector y : $y_i = \begin{cases} x_i & \text{if } x_i \geq 0 \\ 0 & \text{if } x_i < 0 \end{cases}$.

Define $R(x)$ as $x^T L x / x^T x$.

Claim $R(y) \leq R(x)$.

proof $(Ly)_i = y_i - \sum_{j \in N(i)} \frac{y_j}{d} = x_i - \sum_{j \in N(i)} \frac{x_j}{d} = (Lx)_i = \lambda_2 \cdot x_i \quad \forall i \text{ with } y(i) > 0$.

Therefore, $y^T L y = \sum_{i \in V} y_i \cdot (Ly)_i \leq \sum_{i: y(i) > 0} \lambda_2 x_i^2 = \sum_i \lambda_2 y_i^2$, proving the claim. \square

There is a very elegant argument to make the above intuition precise: just pick a random threshold.

Lemma Given any y , there exists a subset $S \subseteq \text{supp}(y)$ such that $\phi(S) \leq \sqrt{2R(y)}$, where $\text{supp}(y) = \{i \mid y_i \neq 0\}$.

Proof We can assume that $0 \leq y_i \leq 1$ for all i , by scaling y if necessary.

Let $t \in (0, 1]$ be chosen uniformly at random.

Let $S_t = \{i \mid y_i^2 \geq t\}$. Then $S_t \subseteq \text{supp}(y)$ by construction.

We analyze the expected value of $|\delta(S_t)|$ and the expected value of $|S_t|$.

$$\begin{aligned} E(|\delta(S_t)|) &= \sum_{ij \in E} [\Pr(\text{the edge } ij \text{ is cut})] \quad \text{by linearity of expectation} \\ &= \sum_{ij \in E} [\Pr(y_i^2 < t \leq y_j^2)] \\ &= \sum_{ij \in E} |y_j^2 - y_i^2| \\ &= \sum_{ij \in E} |y_i - y_j| |y_i + y_j| \\ &\leq \sqrt{\sum_{ij \in E} (y_i - y_j)^2} \sqrt{\sum_{ij \in E} (y_i + y_j)^2} \quad \text{by Cauchy-Schwarz } \langle a, b \rangle \leq \|a\| \cdot \|b\| \\ &\leq \sqrt{\sum_{ij \in E} (y_i - y_j)^2} \sqrt{2 \sum_{ij \in E} (y_i^2 + y_j^2)} \\ &= \sqrt{\sum_{ij \in E} (y_i - y_j)^2} \sqrt{2d \sum_{i \in V} y_i^2} \\ &= \sqrt{2R(y)} \left(d \sum_{i \in V} y_i^2 \right) \end{aligned}$$

$$E[|S_t|] = \sum_{i \in V} \Pr[y_i^2 \geq t] = \sum_{i \in V} y_i^2$$

$$\text{Therefore, } \frac{E[|\delta(S_t)|]}{E[d|S_t|]} \leq \sqrt{2R(y)}.$$

$$\text{This means that } E[|\delta(S_t)| - \sqrt{2R(y)} \cdot d \cdot |S_t|] \leq 0.$$

$$\text{Hence, there exists } t \text{ such that } \frac{|\delta(S_t)|}{d \cdot |S_t|} \leq \sqrt{2R(y)}. \quad \square$$

Combining the claim and the lemma proves Cheeger's inequality.

And the proof shows that the spectral partitioning algorithm achieves the performance guarantee, because the output set S_t is a "threshold" set that the algorithm searches.

Discussions

① The proof can be generalized to weighted non-regular graphs, with minor modifications.

② Both sides of Cheeger's inequality are tight - even the constants are tight.

To see an example where the hard direction is (almost) tight, consider a cycle of length n .

One can compute the spectrum of the cycle exactly, but we won't do it here.

Recall that $\lambda_2 = \min_{x \perp \mathbf{1}} \frac{x^T L x}{x^T x}$, so to give an upper bound on λ_2 , we just need to demonstrate one vector.

Consider $x = (1, 1 - \frac{1}{n}, 1 - \frac{2}{n}, \dots, \frac{1}{n}, 0, -\frac{1}{n}, \dots, -1 + \frac{1}{n}, -1, -1 + \frac{1}{n}, \dots, 0, \frac{1}{n}, \dots, 1)$.

Then $\lambda_2 \leq \frac{\sum_{i,j} (x_i - x_j)^2}{2 \sum_i x_i^2} = O\left(\frac{n(\frac{1}{n})^2}{n}\right) = O(\frac{1}{n^2})$.

On the other hand, it is easy to verify that the expansion of a cycle is $\Omega(\frac{1}{n})$.

Therefore, in this example, $\phi(G) = \Omega(\sqrt{\lambda_2})$.

You may think that it is an artificial example in which the second eigenvalue clearly underestimates the expansion, but let's consider the following related example.

Two cycles of length n , and there is a perfect matching between the two cycles, where each edge has weight $100/n^2$.



Clearly, the optimal sparse cut is the perfect matching, with $\phi(G) = O(\frac{1}{n^2})$.

On the other hand, one can show that the second eigenvector would still be the same as the cycle example (with two nodes in the perfect matching identified as one node).

Therefore, λ_2 is still $O(\frac{1}{n^2})$ and the value is correct, but the optimal cut is lost and every threshold cut is bad.

These examples show how Cheeger's inequality got cheated, both in terms of the value and in terms of the actual cut returned.

- ③ Cheeger's inequality gives an $O(\frac{1}{\sqrt{\lambda_2}})$ -approximation algorithm for computing $\phi(G)$. When λ_2 is large, then it is a pretty good approximation. But λ_2 could be as small as $O(\frac{1}{n^2})$, and so it could be an $\Omega(n)$ -approximation. It doesn't quite explain the good performance in practice. We will come back to this question later.
- ④ The second eigenvalue is closely related to the mixing time of random walks, and so Cheeger's inequality provides a combinatorial approach to bound the mixing time, and we will see it next lecture.
- ⑤ The proof can be modified to show that any vector with Rayleigh quotient $\approx \lambda_2$ could be used to produce a sparse cut of conductance $O(\sqrt{\lambda_2})$, i.e. it doesn't have to be an eigenvector. We will use this fact later.
- Note that our rounding step works fine, but our truncation step uses the fact that it is an eigenvector. This can be handled by a more complicated argument; see e.g. my previous notes on Cheeger's inequality.
-

Last Eigenvalue

Now we survey the recent developments relating other eigenvalues to graph partitioning problems.


In the following we will just state the results and discuss some main ideas.

Recall that Cheeger's inequality is a robust generalization of the basic spectral characterization of connectedness. Can we generalize the spectral characterization of bipartiteness? It turns out that the proof in Cheeger's inequality can be adapted in this setting.

$$\text{We define } \beta(G) = \min_{y \in \{-1, 0, +1\}^V} \frac{\sum_{ij \in E} |y_i + y_j|}{d \sum_{i \in V} |y_i|} = \min_{\substack{S \subseteq V \\ (L, R) \text{ partition of } S}} \frac{2|\# \text{ edges within } L| + 2|\# \text{ edges within } R| + |\delta(S)|}{d|S|}.$$

This is called the bipartiteness ratio of G . Note that $\beta(G)$ is small if and only if G

contains a subset $S \subseteq V$ which is close to a bipartite component, with most edges in S

crossing L and R .  (edges within are counted twice and edges going out S are counted once)

Note that we still assume the graph is d -regular for simplicity.

Consider the matrix $I + A = I + \frac{1}{d}A$. Let $2 = \alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_n \geq 0$ be its eigenvalues.

Recall that $\alpha_n = 0$ if and only if G has a bipartite component.

Trevisan proved the following generalization.

Theorem $\frac{1}{2} \alpha_n \leq \beta(G) \leq \sqrt{2\alpha_n}$

The proof is very similar to the proof of Cheeger's inequality ("randomized rounding"), and it is left in the problem set.

Applying this theorem recursively, Trevisan obtained an approximation algorithm for the maximum cut problem with worst case approximation ratio 0.531.

Note that getting a 0.5-approximation is trivial, and there is a 0.878-approximation algorithm by semidefinite programming. The spectral algorithm is the only known alternative method to do better than 0.5.

The k-th Eigenvalue

Let $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n \leq 2$ be the eigenvalues of the normalized Laplacian matrix.

Recall the spectral characterization that $\lambda_k = 0$ if and only if G has k connected components.

It turns out that there are two meaningful ways to generalize this basic fact.

① Small sparse cut: If λ_k is small, then there is a sparse cut S with $|S| \approx |V|/k$.

② Many sparse cuts: If λ_k is small, then there are k disjoint sparse cuts.

It may appear that ② is more general than ①, but the results obtained are incomparable as we will explain soon.

The informal intuition is that each eigenvector defines a sparse cut, and since the eigenvectors are orthogonal, the sparse cuts should look quite different, and thus cut the graph into pieces.

Proving this intuition formally is another story.

Small Sparse Cut

Arora, Barak, and Steurer proved:

Theorem For $k = \Omega(n^\varepsilon)$ for $\varepsilon \in (0, 1)$, there is a set S with $\phi(S) = O(\sqrt{\lambda_k})$ and $|S| \approx n/k$.

The proof has some nice new ideas. Consider the matrix $W := \frac{1}{2}I + \frac{1}{2}A = \frac{1}{2}I + \frac{1}{2d}A$ (d -regular).

Let $1 \geq \alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_n \geq 0$ be the eigenvalues of this "lazy random walk matrix".

So the condition that $\lambda_k \approx 0$ is translated to $\alpha_k \approx 1$.

Consider $\text{Tr}(W^k)$.

Recall that $\text{Tr}(W^k) = \text{sum of eigenvalues of } W^k = \sum_i \lambda_i^k = \text{"large" by our assumption.}$

On the other hand, $\text{Tr}(W^k) = \text{sum of diagonal entries of } W^k = \text{sum of "returning probabilities" after } k \text{ steps of random walk.}$

If every small set has large conductance, then we expect that $\text{Tr}(W^k)$ is small, contradicting to our assumption. So there must exist a small sparse cut.

We will discuss this result in details next week.

Many Sparse Cuts

Two research groups proved essentially the same results (with Trevisan in one group, and Raghavendra in another)

Theorem Let $\Phi_k(G) = \min_{S_1, \dots, S_k} \max_{i \text{ disjoint}} \phi(S_i)$. Then $\frac{1}{2} \lambda_k \leq \Phi_k(G) \leq O(k^2) \sqrt{\lambda_k}$.

The first inequality is the easy direction, and may leave to you in the problem set.

Both use the spectral embedding: Let x_1, \dots, x_k be the first k eigenvectors. Map each vertex i to a k -dimensional point $(x_1(i), x_2(i), \dots, x_k(i))$.

Since x_1, \dots, x_k are orthonormal, it can be proved that the points are "well-spread out".

The algorithm of one group is very simple: just pick k random directions, each point is put to the cluster defined by its closest direction.

Both proofs are quite non-trivial. We won't discuss them in this course.

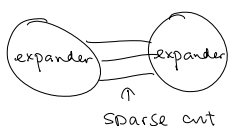
Improved Cheeger's Inequality

Finally, I would like to mention a recent result on analysis of spectral partitioning through higher eigenvalues.

Let $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n \leq 2$ be the eigenvalues. Suppose λ_2 is small.

If λ_2 is small, then we know that there will be three disjoint sparse cuts.

What if λ_2 is large? Then we know that there is a good way to cut the graph into two pieces, but there is no good way to cut the graph into three pieces.

Then we expect the graph to look like , and the second eigenvector should look like a "binary" vector,

sparse cut

and thus $\lambda_2 \approx \phi(G)$, and thus λ_2 is a better approximation to $\phi(G)$ when λ_k is large.

In general, for any $k \geq 2$,

Theorem $\phi(G) \leq O(k) \frac{\lambda_2}{\sqrt{\lambda_k}}.$

The proof is to show that when λ_k is large, then the second eigenvector looks like a k -step function,

and thus λ_2 cannot be much smaller than $\phi(G)$, i.e. Cheeger's inequality can't be tight.

It shows that spectral partitioning is a constant factor approximation when λ_k is large for a small k .

In image segmentation and clustering, there are usually only a few outstanding sparse cuts, thus

λ_k is large. Thus, the result shows that spectral partitioning actually gives a good approximation

in those instances. It gives some theoretical justification of the empirical performance of spectral

partitioning. The result can be generalized to the multiway partitioning problem.

References

You are referred to the course notes and publications of Luca Trevisan for more information.