

Graph Sparsification

- Sparse graphs:

Graphs with less # of edges

⇒ less space to store the graphs

⇒ less processing time

- IDEA: Given an undirected weighted graph $G=(V, E, w)$, where $|V|=n$, $|E|=n$, the goal is to output the graph H , a subgraph of G with fewer edges, where H may be reweighted, while preserving "interesting quantities"

Interesting quantities:

e.g. extremal (min, max) cuts,

eigenvalues, random walk properties.

(typically captured by graph Laplacian)

Graph Sparsification w.r.t. cuts

- For simplicity, let's assume we are given an undirected unweighted graph G . The goal is to approximate G by a sparse graph while preserving the cut size for all possible cuts with small errors.
- More precisely, a cut (S, \bar{S}) is a partition of V into two subsets $S, \bar{S} = V \setminus S$. Let $E_G(S, \bar{S})$ denote the set of edges crossing the cut (S, \bar{S}) in G , i.e., $E_G(S, \bar{S}) = \{u, v \in E(G) \mid u \in S \text{ and } v \notin S\}$. The capacity (or size) of the cut (S, \bar{S}) is denoted by $|E_G(S, \bar{S})|$. If G is weighted by a weight function $w: E \rightarrow \mathbb{R}$, then $|E_G(S, \bar{S})| = \sum_{e \in E_G(S, \bar{S})} w(e)$.

Goal: Construct a graph $H = (V, E')$, where $E' \subset E$ and $|E'| \ll |E|$, and H is potentially reweighted by a function $w': E' \rightarrow \mathbb{R}$ s.t. $\forall U \subseteq V$

$$|E_H(U, \bar{U})| = (1 \pm \epsilon) |E_G(U, \bar{U})|$$

for a small fixed $\epsilon > 0$. Such a graph H is called a "cut sparsifier" of G .

Cut Sparsifier of Complete Graphs

- Given a complete graph $G = K_n$, how can we construct a cut sparsifier H of K_n ? (not that $|E(K_n)| = O(n^2)$.)

• Sub sampling: Consider the following process:

- Sample (keep) every edge independently with some probability p . Then,

$$\mathbb{E}[|E(H)|] = p|E(K_n)|, \text{ and}$$

$$\forall u \in V(K_n), \mathbb{E}[|E_H(u, \bar{u})|] = p|E_{K_n}(u, \bar{u})|$$

- Let's assign weight $\frac{1}{p}$ to each edge of H so that

$$\forall u \in V(K_n), \mathbb{E}[|E_H(u, \bar{u})|] = p \cdot \frac{1}{p} |E_{K_n}(u, \bar{u})|.$$

- So in expectation, cuts are preserved !!

Let's analyze how likely the cut capacity is close to its expectation,

Theorem [Chernoff-Hoeffding Concentration Bound]

Let $X = \sum_{i \in [n]} X_i$, where each $X_i \in [0, 1]$ is

an indicator random variable, and $(X_i:$

$i \in [n])$ are independently distributed. Then,

$$\bullet \forall t > 0, \Pr\{|X - \mathbb{E}[X]| > t\} \leq e^{-\frac{2t^2}{n}}$$

$$\bullet \forall 0 < \epsilon < 1, \Pr\{X < (1 - \epsilon)\mathbb{E}[X]\} \leq e^{-\frac{\epsilon^2 \mathbb{E}[X]}{2}}$$

$$\bullet \forall 0 < \epsilon < 1, \Pr\{X > (1 + \epsilon)\mathbb{E}[X]\} \leq e^{-\frac{\epsilon^2 \mathbb{E}[X]}{3}}$$

$$\bullet \forall t > 2\epsilon \mathbb{E}[X], \Pr\{X > t\} \leq 2^{-t}.$$

Analysis of Subsampling

- To simplify the analysis, we consider H in the unweighted version.

For a subset $U \subseteq V$, let $q = |U|$, $C_H = |E_H(U, \bar{U})|$, and $C_G = |E_G(U, \bar{U})| \geq \frac{q}{2}n$, thus

$E[C_H] = pC_G \geq \frac{pq}{2}n$. Using the concentration bound from above,

$$\Pr\{C_H > (1+\varepsilon)E[C_H]\} \leq e^{-\frac{\varepsilon^2 E[C_H]}{3}} = e^{-\frac{\varepsilon^2 pC_G}{3}} \leq e^{-\frac{\varepsilon^2 pqn}{6}}$$

- Suppose we want the RHS to be at most $1/n^{d/2}$ for some fixed $d > 1$ (say $d=5$)

Then, we should set $p \geq \frac{6d \log n}{\varepsilon^2 n}$

$$\Rightarrow e^{-\frac{\varepsilon^2 pqn}{6}} \leq e^{-\frac{\varepsilon^2 (\frac{6d \log n}{\varepsilon^2 n}) qn}{6}} = e^{-dq \log n} = \frac{1}{n^{dq}}$$

- Note that a similar bound applies to deviation in the other direction, we get

$$\Pr\{C_H \notin (1 \pm \varepsilon)E[C_H]\} \leq \frac{2}{n^{d/2}}$$

- Also, note that the failure probability above is for a single cut. The probability for any cut to fail is obtained by the following analysis.

$$\Pr\{\text{any cut fails}\} \leq \Pr\{\text{any cut fails independently}\}$$

$$= \sum_{1 \leq q \leq N} (\# \text{ of cut of size } q) \cdot \Pr\{\text{cut of size } q \text{ fails}\}$$

$$\leq \sum_{1 \leq q \leq N} \binom{n}{q} \cdot \frac{2}{n^{dq}} \quad \theta(n^q)$$

$$\leq \sum_{1 \leq q \leq N} n^q \cdot \frac{2}{n^d \cdot n^{dq}} \leq \frac{2n}{n^d} = \frac{2}{n^{d-1}}$$

- The subsampling will fail with probability at most $\frac{2}{n^{d-1}}$, if $p \geq \frac{6d \log n}{\varepsilon^2 n}$
- If we set $d=5$ then, the subsampling will success with high probability (i.e., the success probability is at least $1 - \frac{2}{n^4}$), and the graph sparsifier H will have $p(E(K_n)) = \frac{(6 \times 5) \log n}{\varepsilon^2 n} \times \Theta(n^2) = \Theta\left(\frac{n \log n}{\varepsilon^2}\right)$
 $\Rightarrow E[|E(H)|] = \tilde{\Theta}\left(\frac{n}{\varepsilon^2}\right)$
- Clearly, the subsampling of H takes time $O(m)$.

Theorem 1: The subsampling can construct a cut sparsifier H of a complete graph K_n in time $O(m)$, where $|E(H)| = \tilde{\Theta}\left(\frac{n}{\varepsilon^2}\right)$, with high probability of success.

Cut Sparsifier of Arbitrary Graphs

- Let's check if the same subsampling also works for arbitrary graph $G = (V, E)$

- Same as before,

$$\mathbb{E}[|E(H)|] = p |E(G)|,$$

$$\forall U \subseteq V, \mathbb{E}[|E_H(U, \bar{U})|] = \frac{1}{p} |E_G(U, \bar{U})|$$

\Rightarrow All cuts are preserved up to scaling by a factor $1/p$ in expectation!!

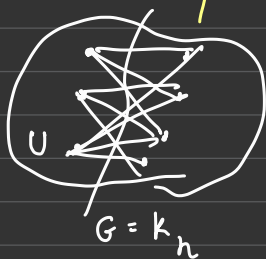
- Next, we use the concentration bound to analyze the followings:

$$\Rightarrow \Pr \{ C_H > (1+\epsilon) \mathbb{E}[C_H] \} \leq e^{-\frac{\epsilon^2 \mathbb{E}[C_H]}{2}} = e^{-\frac{\epsilon^2 p C_G}{2}}$$

- Note that this failure probability is exponentially decreasing in C_G .

- In the case of G being a complete graph, we can simply obtain the bound

$$C_G \geq \frac{nq}{2}$$



$$C_G = |U| \cdot |\bar{U}| \text{ and}$$

we know that either $|U| = q \geq \frac{n}{2}$

or $|\bar{U}| = n - q \geq \frac{n}{2}$.

However, for the case of arbitrary graphs, we cannot obtain this kind of bound because no assumption can be made about the structure of G .

- Instead, we will use the fact about the number of cuts of small size.

Lemma 2 [Karger]: Let G be a graph with n vertices.

Let C_G^* denote the size of minimum cuts of G . Then, for every $\alpha \geq 1$, the number of cuts of size at most αC_G^* is at most $n^{2\alpha}$.

Observation: A cut of size αC_G^* fails with probability at most $e^{-\frac{\epsilon^2 p \alpha C_G^*}{3}}$

- Suppose among all cuts of size αC_G^* we want each cut to deviate from its mean more than a factor of $1 \pm \epsilon$, independently, with probability at most $\frac{1}{n^{d+2}}$. Then,

$$n^{2\alpha} \cdot e^{-\frac{\epsilon^2 p \alpha C_G^*}{3}} \leq \frac{1}{n^{d+2}}$$

$$\Leftrightarrow n^{2\alpha + d+2} \leq e^{\frac{\epsilon^2 p \alpha C_G^*}{3}}$$

$$\Leftrightarrow \alpha(d+2) \ln n \leq \frac{\epsilon^2 p \alpha C_G^*}{3}$$

$$\Leftrightarrow p \geq \frac{(d+2) \ln n}{\epsilon^2 C_G^*}$$

Note that for each cut to deviate more than

a factor of $(1 \pm \epsilon)$, independently, we have $\frac{2}{n^{d+2}}$ for the bounded probability.

- Now, with $p \geq \frac{(d+2) \ln n}{\varepsilon^2 C_G^*}$

$$\Pr\{\text{any cut fails}\} \leq \Pr\{\text{all cut fails independently}\}$$

$$\leq \int_1^\infty 2 \left(n^{2\alpha} \cdot e^{-\frac{\varepsilon^2 p \cdot \alpha C_G^*}{2}} \right) 2\alpha \leq \int_1^\infty \frac{2}{n^{d\alpha}} 2\alpha = O\left(\frac{1}{n^d}\right)$$

Therefore, the subsampling succeeds w.h.p.

- In addition, given the value of C_G^* beforehand, the subsampling works in $O(m)$ time, and the resulting sparsifier H has

$$\mathbb{E}[|E(H)|] \geq \frac{(d+2) \log n \cdot m}{\varepsilon^2 C_G^*}$$

$$= \tilde{\Theta}\left(\frac{m}{\varepsilon^2 C_G^*}\right)$$

So, the number of edges decreases roughly by

a factor of $\tilde{\Theta}\left(\frac{\log n}{C_G^*}\right)$ ← We shall discuss the limitation due to this result later.

Theorem 3:

The subsampling can construct a sparsifier H of an arbitrary graph G , where $|E(H)| = \tilde{\Theta}\left(\frac{m}{\varepsilon^2 C_G^*}\right)$, with high probability of success.

