# On a Connection Between Small Set Expansions and Modularity Clustering

Bhaskar DasGupta[1,*]

*Department of Computer Science, University of Illinois at Chicago, Chicago, IL 60607*

Devendra Desai

*Department of Computer Science, Rutgers University, Piscataway, NJ 08854*

## Abstract

In this paper we explore a connection between two seemingly different problems from two different domains: the *small-set expansion* problem studied in unique games conjecture, and a popular community finding approach for social networks known as the *modularity clustering* approach. We show that a sub-exponential time algorithm for the small-set expansion problem leads to a sub-exponential time constant factor approximation for some hard input instances of the modularity clustering problem.

*Keywords:* Small-set expansion, modularity clustering, social network
*2010 MSC:* 68Q25, 68W25

## 1. Introduction and Definitions

All graphs considered in this note are *undirected* and *unweighted*[2]. Let $G = (V, E)$ denote the given input graph with $n = |V|$ nodes and $m = |E|$

---

*Corresponding author.
Email addresses:* `dasgupta@cs.uic.edu` (Bhaskar DasGupta),
`devdesai@cs.rutgers.edu` (Devendra Desai)
*URL:* `www.cs.uic.edu/~dasgupta` (Bhaskar DasGupta),
`paul.rutgers.edu/~devdesai` (Devendra Desai)

[2]Our result can be extended for the more general case of directed weighted graphs by using the correspondence of these versions with unweighted undirected graphs as outlined in [4, Section 5.1].

edges, let $d_v$ denote the degree of a node $v \in V$, and let $A(G) = \left[a_{u,v}(G)\right]$ denote the adjacency matrix of $G$, *i.e.*, $a_{u,v}(G) = \begin{cases} 1, & \text{if } \{u,v\} \in E \\ 0, & \text{otherwise.} \end{cases}$ Since our result spans over two distinct research areas, we summarize the relevant definitions from both research fields [1, 6] below for convenience.

(a) By a "set of $(k)$ *communities*" we mean a partition of the set of nodes $V$ into $(k)$ non-empty parts.

(b) If $G$ is $d$-regular for some given $d$, then its symmetric *stochastic walk* matrix is denoted by $\widehat{A}(G)$, and is defined as the $n \times n$ real symmetric matrix $\widehat{A}(G) = \left[\frac{a_{u,v}(G)}{d}\right]$.

(c) For a real number $\tau \in [\,0,1)$, the $\tau$-*threshold rank* of $G$, denoted by $\operatorname{rank}_\tau(G)$, is the number of eigenvalues $\lambda$ of $\widehat{A}(G)$ satisfying $|\lambda| > \tau$.

(d) For a subset $\emptyset \subset S \subset V$ of nodes, the following quantities are defined:

- The (normalized) *measure* of $S$ is $\mu(S) = \frac{|S|}{n}$.

- The (normalized) *expansion* of $S$ is $\Phi(S) = \dfrac{\left|\,\{\,\{u,v\}\,|\,u \in S,\ v \notin S,\ \{u,v\} \in E\,\}\,\right|}{\displaystyle\sum_{v \in S} d_v}$

- The (normalized) *density* of $S$ is $\mathsf{D}(S) = 1 - \Phi(S)$.

- The *modularity* value of $S$ is $\mathsf{M}(S) = \frac{1}{2m}\left(\displaystyle\sum_{u,v \in S}\left(a_{u,v} - \frac{d_u d_v}{2m}\right)\right)$

(e) The modularity of a set of communities $\mathbf{S}$ is $\mathsf{M}(\mathbf{S}) = \sum_{S \in \mathbf{S}} \mathsf{M}(S)$.

(f) The goal of the *modularity $k$-clustering* problem on an input graph $G$ is to find a set of at most $k$ communities $\mathbf{S}$ that *maximizes* $\mathsf{M}(\mathbf{S})$. Let $\mathsf{OPT}_k(G) = \max\limits_{\mathbf{S}\text{ is a set of at most k communities}} \left\{\,\mathsf{M}(\mathbf{S})\,\right\}$ denote the optimal modularity value for a modularity $k$-clustering; it is easy to verify that $0 \leq \mathsf{OPT}_k(G) < 1$.

(g) The goal of the *modularity clustering* problem on $G$ is to find a set of (unspecified number of) communities $\mathbf{S}$ that *maximizes* $\mathsf{M}(\mathbf{S})$. Let $\mathsf{OPT}(G)$

denote the optimal modularity value for a modularity clustering; obviously, $\mathsf{OPT}(G) = \mathsf{OPT}_n(G)$.

(h) $\exp(\xi)$ denotes $2^{c\xi}$ for some constant $c > 0$ that is independent of $\xi$.

The modularity clustering problems as described above is *extremely popular* in practice in their applications to biological networks [8, 9] as well as to social networks [5–7]. For relevant computational complexity results for modularity maximization, see [2, 4]. The following results from [4] demonstrate the computational hardness of $\mathsf{OPT}_2(G)$ and $\mathsf{OPT}(G)$ even if $G$ is a regular graph.

**Theorem 1.1.** [4]

**(a)** *For every constant $d \geq 9$, there exists a collection of $d$-regular graphs $G$ of $n$ nodes such it is NP-hard to decide if $\mathsf{OPT}_2(G) \geq \frac{1}{2} - \frac{2c}{dn}$ or if $\mathsf{OPT}_2(G) \leq \frac{1}{2} - \frac{2c+2}{dn}$ for some positive $c = O(\sqrt{n})$.*

**(b)** *There exists a collection of $(n-3)$-regular graphs $G$ of $n$ nodes such it is NP-hard to decide if $\mathsf{OPT}(G) > \frac{0.9388}{n-4}$ or if $\mathsf{OPT}(G) < \frac{0.9382}{n-4}$.*

## 2. Our Result

**Theorem 2.1.** *Let $G$ be a $d$-regular graph. Then, for some constant $0 < \varepsilon < 1/2$, there is an algorithm $\mathcal{A}_\varepsilon$ with the following properties:*

- *$\mathcal{A}_\varepsilon$ runs in sub-exponential time, i.e., in time $\exp(\delta n)$ for some constant $0 < \delta = \delta(\varepsilon) < 1$ that depends on $\varepsilon$ only.*

- *$\mathcal{A}_\varepsilon$ correctly distinguishes instances $G$ of modularity clustering with $\mathsf{OPT}(G) \geq 1 - \varepsilon$ from instances $G$ with $\mathsf{OPT}(G) \leq \varepsilon$.*

*(Note that we make no claim if $\varepsilon < \mathsf{OPT}(G) < 1 - \varepsilon$.)*

**Remark 2.2 (usability of the approximation algorithm in Theorem 2.1).** *We prove Theorem 2.1 for $\varepsilon = 10^{-6}$. It is natural to ask if there are in fact infinite families of $d$-regular graphs $G$ that satisfy $\mathsf{OPT}(G) \geq 1 - 10^{-6}$ or $\mathsf{OPT}(G) \leq 10^{-6}$. The answer is affirmative, and we provide below examples of infinite families of such graphs.*

> $\mathsf{OPT}(\mathbf{G}) \geq \mathbf{1} - \mathbf{10^{-6}}$: *Consider, for example, the following explicit bound was demonstrated in [2, Corollary 6.4]:*

> *if $G$ is an union of $k$ disjoint cliques each with $\frac{n}{k} > 3$*
> *nodes then* $\mathsf{OPT}(G) = 1 - \frac{1}{k}$.

*Based on this and other known results on modularity clustering, examples of families of regular graphs $G$ for which $\mathsf{OPT}(G) \geq 1 - 10^{-6}$ include:*

**(1)** *$G$ is an union of $k$ disjoint cliques each with $\frac{n}{k} > 3$ nodes for any $k > 10^6$.*

**(2)** *$G$ is obtained by a local modification from the graph in **(1)** such as:*

- *Start with an union of $k$ disjoint cliques $\mathcal{C}_1, \mathcal{C}_2, \ldots, \mathcal{C}_k$ each with $\frac{n}{k} > 3$ nodes for any $k$ sufficiently large with respect to $10^6$ ($k \geq 10^7$ suffices).*
- *Remove an arbitrary edge $\{u_i, v_i\}$ from each clique $\mathcal{C}_i$. Let $U = \cup_{i=1}^{k}\{u_i\}$ and and $V = \cup_{i=1}^{k}\{v_i\}$.*
- *Add to $G$ the edges corresponding to any perfect matching in the complete bipartite graph with node sets $U$ and $V$.*

$\mathsf{OPT}(\mathbf{G}) \leq \mathbf{10^{-6}}$: *Theorem 1.1 [4] involves infinitely many graphs of $n > 4 + 0.9388 \times 10^6$ nodes satisfying $\mathsf{OPT}(G) < \frac{0.9388}{n-4} < 10^{-6}$ (these graphs are edge complements of appropriate families of 3-regular graphs used in [3]).*

*Proof of Theorem 2.1.*[3] Set $\varepsilon = 10^{-6}$. We assume that $G$ is $d$-regular, and either $\mathsf{OPT}(G) \geq 1 - 10^{-6}$ or $\mathsf{OPT}(G) \leq 10^{-6}$.

*Preliminary Algebraic Simplification*

Let $\mathbf{S} = \big\{S_1, S_2, \ldots, S_k\big\}$ be a set of communities of $G$. The objective function $\mathsf{M}(\mathbf{S})$ can be equivalently expressed as follows via simple algebraic manipulation [2, 5–7]. Let $m_i$ denote the number of edges whose both endpoints are in $S_i$, $m_{ij}$ denote the number of edges one of whose endpoints is in

---

[3]We have made no significant attempts to optimize the constants in Theorem 2.1.

$S_i$ and the other in $S_j$ and $D_i = \sum\limits_{v \in S_i} d_v$ denote the sum of degrees of nodes in $S_i$. Then, $\mathsf{M}(\mathbf{S}) = \sum\limits_{S_i \in \mathbf{S}} \left( \frac{m_i}{m} - \left( \frac{D_i}{2m} \right)^2 \right)$.

We will provide an approximation for $\mathsf{OPT}_2(G)$ and then use the result that $\mathsf{OPT}_2(G) \geq \frac{\mathsf{OPT}(G)}{2}$ proved in [4]. Note that if if $\mathsf{OPT}(G) \leq 10^{-6}$ then obviously $\mathsf{OPT}_2(G) \leq 10^{-6}$, whereas if $\mathsf{OPT}(G) \geq 1 - 10^{-6}$ then $\mathsf{OPT}_2(G) \geq \frac{1}{2} - \frac{10^{-6}}{2}$. Consider a partition $\mathbf{S}$ of $V$ into exactly two sets, say $S$ and $\overline{S} = V \setminus S$ with $0 < \mu(S) \leq \text{\textonehalf}$. By Lemma 2.2 of [4], $\mathsf{M}(S) = \mathsf{M}(\overline{S})$ and thus

$$\mathsf{M}(\mathbf{S}) = 2 \times \left( \frac{m_1}{m} - \left( \frac{|S|}{n} \right)^2 \right) \ = \ 2 \times \left( \frac{\frac{1}{2} \, \mathsf{D}(S) \, d \, |S|}{\frac{1}{2} \, d \, n} - \mu(S)^2 \right)$$

$$= \ 2 \times \left( \mathsf{D}(S) \, \mu(S) - \mu(S)^2 \right)$$

Thus, letting $\mathsf{D} = \mathsf{D}(S), \mu = \mu(S)$ and $\Phi = \Phi(S)$, we have $\Phi = 1 - \mathsf{D}$ as per our notations used in page 2 and the goal of modularity 2-clustering is to maximize the following function $f$ over all possible valid choices of $\mathsf{D}$ and $\mu$:

$$f(\mu, \mathsf{D}) = 2 \times \left( \mu \, \mathsf{D} - \mu^2 \right) = 2 \times \left( \mu(1 - \Phi) - \mu^2 \right)$$

Let $\mathbf{S}^\star = \{ S^\star, \overline{S^\star} \}$ be an optimal solution for modularity 2-clustering of $G$, with $\mathsf{D} = \mathsf{D}^\star, \mu = \mu^\star, \Phi = \Phi^\star$ (and thus $\mathsf{OPT}_2(G) = f(\mu^\star, \mathsf{D}^\star)$). Obviously,

$$\left| \mu^\star - \frac{\mathsf{D}^\star}{2} \right| < \frac{\mathsf{D}^\star}{2}$$

$$f \left( \frac{\mathsf{D}^*}{2} + \delta, \mathsf{D}^* \right) = f \left( \frac{\mathsf{D}^*}{2} - \delta, \mathsf{D}^* \right) \ \text{for any positive } \delta > 0$$

Note that we need to show that, if $\mathsf{OPT}_2(G) = f(\mu^\star, \mathsf{D}^\star) > \frac{1}{2} - \frac{10^{-6}}{2}$, then there is an algorithm $\mathcal{A}_\varepsilon$ as described in Theorem 2.1 that outputs a valid choice of $\mu$ and $\mathsf{D}$, say $\mu'$ and $\mathsf{D}'$, such that $f(\mu', \mathsf{D}') > 10^{-6}$.

*Guessing* $\mathsf{D}^\star$

Note that there are at most $\mathrm{O}\left( d \, n^2 \right)$ choices for $\mathsf{D}^\star$ since $\mathsf{D}^\star$ is of the form $i/(j\,d)$ for $j \in \{1, 2, \ldots, n/2\}$ and $i \in \{1, 2, \ldots, j\,d\}$. In the sequel, we will run our algorithm for each choice of $\mathsf{D}^\star$ and take the best of these solutions. Thus, it will suffice to prove our approximation bound assuming we have guessed $\mathsf{D}^\star$ exactly.

In the remainder of the proof, we will make use of results for small-set expansion from [1]. The description is self-contained, and the reader *will not need any prior knowledge of expansion properties of graphs.* Remember that we assume that $f(\mu^\star, \mathsf{D}^\star) > \frac{1}{2} - \frac{10^{-6}}{2}$ and thus $\mu^\star > \frac{1}{2} - \frac{10^{-3}}{2}$ since otherwise

$$
\begin{aligned}
& \mu^\star \leq \tfrac{1}{2} - \tfrac{10^{-3}}{2} \\
\Rightarrow \quad & \mu^\star = \tfrac{1}{2} - \tfrac{10^{-3}}{2} - \xi \qquad [\text{ for some } \xi \geq 0 \,] \\
\Rightarrow \quad & f(\mu^\star, \mathsf{D}^\star) = 2 \times \left( \mu^\star \mathsf{D}^\star - (\mu^\star)^2 \right) \\
& \qquad\qquad \leq 2 \times \left( \mu^\star - (\mu^\star)^2 \right) \qquad [\text{ since } 0 \leq \mathsf{D}^\star \leq 1 \,] \\
& \qquad\qquad = 2 \times \mu^\star \times (1 - \mu^\star) \\
& \qquad\qquad = 2 \times \left( \tfrac{1}{2} - \tfrac{10^{-3}}{2} - \xi \right) \times \left( \tfrac{1}{2} + \tfrac{10^{-3}}{2} + \xi \right) \\
& \qquad\qquad = 2 \times \left( \tfrac{1}{4} - \left( \tfrac{10^{-3}}{2} + \xi \right)^2 \right) \\
& \qquad\qquad < \tfrac{1}{2} - \tfrac{10^{-6}}{2}
\end{aligned}
$$

which contradicts $f(\mu^\star, \mathsf{D}^\star) > \frac{1}{2} - \frac{10^{-6}}{2}$. Similarly, we also get:

$$
\mathsf{D}^\star = \frac{f(\mu^\star, \mathsf{D}^\star)}{2\,\mu^\star} + \mu^\star > \left( \frac{1 - 10^{-6}}{4} \right) \frac{1}{\mu^\star} + \mu^\star
$$

Consider the function $g(\mu) = \frac{a}{\mu} + \mu$ where $a = \frac{1 - 10^{-6}}{4}$. Since $\mu > 0$, $\frac{\mathrm{d}^2 g(\mu)}{\mathrm{d}\,\mu^2} = \frac{2a}{\mu^3} > 0$ and thus the minimum of $g(\mu)$ is attained at $\mu = b$ that satisfies $\frac{\mathrm{d}\,g(b)}{\mathrm{d}\,b} = -\frac{a}{b^2} + 1 = 0$, giving $b = \sqrt{a}$. Thus, we have

$$
\mathsf{D}^\star > \left( \frac{1 - 10^{-6}}{4} \right) \left( \frac{1}{\sqrt{\frac{1 - 10^{-6}}{4}}} \right) + \sqrt{\frac{1 - 10^{-6}}{4}} = \sqrt{1 - 10^{-6}} > 1 - 10^{-6}
$$

which implies $\Phi^\star = 1 - \mathsf{D}^\star < 10^{-6}$.

## Case I: $G$ has a small threshold rank, *i.e.*, $\mathrm{rank}_{1-10^{-6}}(G) < n^{10^{-1}}$

The following result, restated below under the assumption of this case in our terminologies after instantiation of parameters with specific values and trivial algebraic simplification, was proved by Arora, Barak and Steurer in [1] in the bigger context of obtaining sub-exponential algorithms for unique games in PCP theory.

**Theorem 2.3.** [1][4] *There exists a* $\left(\exp\left(n^{10^{-1}}\right)\operatorname{poly}(n)\right)$-*time algorithm that outputs a subset* $\emptyset \subset S \subset V$ *such that* $0.92\,|S^\star| \le |S| \le 1.08\,|S^\star|$, *and* $\Phi(S) \le \Phi(S^\star) + 0.08$.

We run the algorithm in Theorem 2.3, and return $\left\{\, S, \overline{S} \,\right\}$ as our solution. Note that:

$$\Phi(S) \le \Phi^\star + 0.08 < 0.080001 \implies \mathsf{D}(S) > 1 - 0.080001 = 0.919999$$
$$0.92\mu^\star \le \mu(S) \le 1.08\mu^\star \implies 0.4599 \le \mu(S) \le 0.54$$

and thus

$$
\begin{aligned}
f(\mu(S), \mathsf{D}(S)) \;&=\; 2 \times \mu(S) \times \big(\mathsf{D}(S) - \mu(S)\big) \\
&>\; 2 \times 0.4599 \times \big(0.919999 - 0.54\big) > 10^{-6}
\end{aligned}
$$

### Case II: Remaining Case, *i.e.*, $\operatorname{rank}_{1-10^{-6}}(G) \ge n^{10^{-1}}$

The following result, restated below in our terminologies after instantiation of parameters with specific values, was again proved in [1].

**Theorem 2.4.** [1][5] *Let* $H$ *be a regular graph of* $r$ *nodes with* $\operatorname{rank}_{1-10^{-5}}(H) \ge r^{10^{-1}}$. *Then, there is an algorithm that*

- *runs in* $\operatorname{poly}(r)$ *time, and*

- *finds a subset* $S$ *of nodes of* $H$ *with* $|S| \le r^{1-10^{-3}}$ *and* $\Phi(S) \le 10^{-2}$.

Our strategy is to use the algorithm in Theorem 2.4 *repeatedly*[6] to extract "high-rank parts" from $G$. Namely, we compute in polynomial time an ordered partition of nodes $\left(T_1, T_2, \ldots, T_k, V \setminus \cup_{i=1}^k T_i\right)$ such that each $T_i$ is obtained by using the algorithm in Theorem 2.4 on graph $G_i$ induced by the set of nodes $V \setminus \cup_{j=1}^{i-1} T_i$, and the last (possibly empty) graph $G''$ induced by

---

[4]Instantiate Theorem 2.2 in [1] with $\eta = 10^{-4}$ and $\varepsilon = 10^{-6}$.
[5]Instantiate Theorem 2.3 in [1] with $\eta = 10^{-4}$ and $\gamma = 10^{-1}$.
[6][1] points out how to "re-regularize" the remaining graph each time a set of nodes have been extracted by adding appropriate number of self-loops of weight $1/2$.

the set of nodes $V'' = V \setminus \cup_{i=1}^{k} T_i$ satisfy $\text{rank}_{1-10^{-6}}(G'') < |V''|^{10^{-1}}$. Let $G'$ be the graph induced by the set of nodes $V' = \cup_{i=1}^{k} T_i$.

**Case II(a)** $\left| S^{\star} \cap V'' \right| \geq |S^*|/2$.

Let $S_1^*$ be the set containing an arbitrary $|S^*|/2$ elements from the set $S^{\star} \cap V''$. Note that $\mu\left(S_1^*\right) = \mu^*/2$ and $\Phi\left(S_1^*\right) \leq 2\,\Phi^*$. We now use Theorem 2.3 on the graph $G''$ with $|S^*|$ replaced by $|S^*|/2$ to output a set $S \subseteq V''$ of nodes such that

$$\Phi(S) \leq 2\,\Phi^{\star} + 0.08 < 0.080002 \implies \mathsf{D}(S) > 1 - 0.080002 = 0.919998$$
$$0.46\mu^{\star} \leq \mu(S) \leq 0.54\mu^{\star} \implies 0.229 < \mu(S) \leq 0.27$$

and thus

$$
\begin{aligned}
f(\mu(S), \mathsf{D}(S)) \;&=\; 2 \times \mu(S) \times \big(\mathsf{D}(S) - \mu(S)\big) \\
&>\; 2 \times 0.229 \times (0.919998 - 0.27) > 10^{-6}
\end{aligned}
$$

**Case II(b)** $\left| S^{\star} \cap V'' \right| < |S^*|/2$.

Since $|S^*| \geq \left(\frac{1}{2} - \frac{10^{-3}}{2}\right) n$ and $|T_j| \leq \left| V \setminus \cup_{\ell=1}^{j-1} T_\ell \right|^{1-10^{-3}} < n^{1-10^{-3}}$ for any $j$, there exists an index $i$ such that $\frac{|S^*|}{2} - n^{1-10^{-3}} < \left| \cup_{j=1}^{i} T_j \right| < \frac{|S^*|}{2} + n^{1-10^{-3}}$. Notice that the graph induced by the set of nodes $S = \cup_{j=1}^{i} T_j$ satisfy $\Phi(S) \leq 10^{-2}$ and, since $\left(\frac{1}{2} - \frac{10^{-3}}{2}\right) n \leq |S^*| \leq n$, we have

$$\frac{|S^*|}{2} - n^{1-10^{-3}} < |S| = \left| \cup_{j=1}^{i} T_j \right| < \frac{|S^*|}{2} + n^{1-10^{-3}} \implies 0.24 < \mu(S) < 0.51$$

and thus, 
$$
\begin{aligned}
f(\mu(S), \mathsf{D}(S)) \;&=\; 2 \times \mu(S) \times \big(\mathsf{D}(S) - \mu(S)\big) \\
&>\; 2 \times 0.24 \times (0.99 - 0.51) > 10^{-6}
\end{aligned}
$$

$\square$

### Further Research

An interesting open question is whether it is possible to prove the converse of Theorem 2.1, *i.e.*, can we use a sub-exponential approximation algorithm for modularity maximization to design a sub-exponential algorithm for small-set expansion problems ? If possible, this may lead to an alternate interpretation of unique games via communities in social networks.

# References

[1] S. Arora, B. Barak and D. Steurer. *Subexponential Algorithms for Unique Games and Related Problems*, 51[st] annual IEEE Symposium on Foundations of Computer Science, 563-572, 2010.

[2] U. Brandes, D. Delling, M. Gaertler, R. Görke, M. Hoefer, Z. Nikoloski and D. Wagner, *On Modularity Clustering*, IEEE Transaction on Knowledge and Data Engineering, 20 (2), 172-188, 2007.

[3] M. Chlebík and J. Chlebíková. *Complexity of approximating bounded variants of optimization problems*, Theoretical Computer Science, 354 (3), 320-338, 2006.

[4] B. DasGupta and D. Desai. *Complexity of Newman's Community Finding Approach for Social Networks*, Journal of Computer & System Sciences, 79, 50-67, 2013.

[5] E. A. Leicht and M. E. J. Newman, *Community Structure in Directed Networks*, Physical Review Letters, 100, 118703, 2008.

[6] M. E. J. Newman, *Modularity and community structure in networks*, PNAS, 103 (23), 8577-8582, 2006.

[7] M. E. J. Newman and M. Girvan, *Finding and evaluating community structure in networks*, Physical Review E, 69, 026113, 2004.

[8] R. Guimer'a, M. Sales-Pardo and L. A. N. Amaral. *Classes of complex networks defined by role-to-role connectivity profiles*, Nature Physics, 3, 63-69, 2007.

[9] E. Ravasz, A. L. Somera, D. A. Mongru, Z. N. Oltvai and A.-L. Barabási. *Hierarchical Organization of Modularity in Metabolic Networks*, Science, 297 (5586), 1551-1555, 2002.