

On the Complexity of Approximating k -Set Packing

Elad Hazan*

Shmuel Safra[†]

Oded Schwartz[‡]

Abstract

Given a k -uniform hypergraph, the MAXIMUM k -SET PACKING problem is to find the maximum disjoint set of edges. We prove that this problem cannot be efficiently approximated to within a factor of $\Omega(\frac{k}{\ln k})$ unless $P = NP$. This improves the previous hardness of approximation factor of $\frac{k}{2^{\Omega(\sqrt{\ln k})}}$ by Trevisan [Tre01]. This result extends to the problem of k -Dimensional-Matching.

1 Introduction

This paper studies the following basic optimization problem: given a family of sets over a certain domain, find the maximum number of disjoint sets. We consider the case where all sets in the given family are of the same size - k .

For the case where $k = 2$, we can view the sets as edges in a graph whose vertices are the domain, and hence the problem is exactly the famous maximal matching problem which is solvable in polynomial time [Pap94]. For $k \geq 3$, again viewing the sets as hyper-edges in a hypergraph, the problem of finding the maximum matching in k -uniform hyper-graphs is NP-hard. Hence, unless $P = NP$, the best hope is to obtain a polynomial time approximation algorithm with provably good approximation guaranty.

The simple greedy algorithm is the following: iteratively pick an arbitrary set and add it to the collection of sets maintained thus far, while removing all sets intersecting it. Continue as long as there remain edges in the graph. Obviously this algorithm returns a family of pairwise disjoint sets. It is easy to prove that this algorithm provides a k approximation to the optimal solution. A constant improvement in the approximation ratio, to $\frac{k}{2}$ [HS89], can be obtained by a simple local search heuristic, and is the best known approximation to date.

In this work we prove that the latter approximation guaranty is almost tight, proving the following:

Theorem 1 *It is NP-hard to approximate k -SP to within $\Omega(\frac{k}{\ln k})$*

1.1 Previous Results

The general MAXIMUM SET PACKING problem is as follows: given a family $\mathcal{F} = \{S_1, \dots, S_m\}$ of sets over a certain domain $\mathcal{D} = \{x_1, \dots, x_n\}$, the objective is to find a maximum packing, namely a maximum number of pairwise disjoint sets from the given family. This problem is often phrased in graph theory terminology, as a set system is in fact a hyper-graph where the the vertices are

*Computer Science Department, Princeton University. www.cs.princeton.edu/~ehazan.

[†]School of Computer Science, Tel Aviv University, Tel Aviv 69978, Israel. safra@post.tau.ac.il.

[‡]School of Computer Science, Tel Aviv University, Tel Aviv 69978, Israel. odedsc@post.tau.ac.il.

the items in the domain and the edges are the given sets. In graph theory jargon, a disjoint set of edges is called a matching, hence the objective is to find a maximum matching.

As packing problems are of the most fundamental and abundant combinatorial optimization problems, variants of MAXIMUM SET PACKING, including the MAXIMUM INDEPENDENT SET and MAXIMUM CLIQUE problems, have been widely studied [Wig83, BYM84, BH92, Hås99, FGL⁺96, AS98, ALM⁺98]. These general formulations of packing problems are notoriously hard even to approximate: Håstad showed [Hås99] that MAXIMUM CLIQUE (and therefore MAXIMUM INDEPENDENT SET and MAXIMUM SET PACKING as well) cannot be approximated to within $O(N^{1-\varepsilon})$ unless $NP \subseteq ZPP$ (for every $\varepsilon > 0$). The best approximation algorithm for MAXIMUM INDEPENDENT SET achieves an approximation ratio of $O(\frac{N}{\log^2 N})$ [BH92].

In this paper we consider several natural variants of packing problems. The first, and perhaps most natural, is when the size of the hyper-edges is bounded by k . This problem is called MAXIMUM k -SET PACKING (in short k -IS). If in addition we bound the degree of the vertices by two, this becomes the problem of maximum independent-set in graphs of degree at most k .

Another (stronger) natural restriction studied here by is when we impose a bound on the colorability of the input graph. This is the problem of MAXIMUM k -DIMENSIONAL MATCHING (in short k -DM). It is a variant of MAXIMUM k -SET PACKING where the vertices of the input hyper-graph are a union of k disjoint sets, $V = V_1 \cup \dots \cup V_k$, and each hyper-edge contains exactly one vertex from each set, namely, $E \subseteq V_1 \times \dots \times V_k$. In other words, the vertices of the hyper-graph can be colored using k colors, so that no hyper-edge contains the same color twice. A graph having this property is called *k -strongly-colorable*. Thus the color-bounded version of MAXIMUM k -SET PACKING is, given a k -uniform k -strongly colorable hyper-graph, find a matching of maximum size.

These bounded variants of MAXIMUM SET PACKING are known to admit approximation algorithms better than their general versions, the quality of the approximation being a function of the bounds. As mentioned previously, the greedy algorithm guaranties a k -approximation for MAXIMUM k -SET PACKING. A simple local-search heuristic achieves an approximation ratio of $\frac{k}{2}$ [HS89]. This is, to date, the best approximation algorithm for MAXIMUM k -DIMENSIONAL MATCHING as well.

For the special case where $k = 2$, both problems are solvable in polynomial time. MAXIMUM 2-DIMENSIONAL MATCHING is just the problem of finding a maximum matching in a bipartite graph, and can be solved in polynomial time, say by a reduction to network flow problems [Pap94]. MAXIMUM 2-SET PACKING is the problem of finding a maximum matching in a general graph, and for this problem polynomial time algorithms are also known [Edm65] (for recent efficient algorithms see [MM04]).

However, for all $k \geq 3$, MAXIMUM k -DIMENSIONAL MATCHING is NP-hard [Kar72, Pap94]. Furthermore, for $k = 3$, the problem is known to be APX-hard [Kan91]. Alon et al. [AFWD95] proved that for a suitably large k , MAXIMUM k -INDEPENDENT SET (finding an independent set of maximum size in k -regular graphs, in short k -IS) is NP-hard to approximate to within k^c for some $c > 0$. This was later improved to the currently best asymptotical inapproximability factor of $\frac{k}{2^{O(\sqrt{\ln k})}}$ [Tre01]. All hardness factors for MAXIMUM k -INDEPENDENT SET hold in fact for MAXIMUM $k + 1$ -DIMENSIONAL MATCHING as well (by a simple reduction).

The best known approximation algorithm for k -IS achieves an approximation ratio of $O(k \log \log k / \log k)$ [Vis96]. For k -IS of low k values, the best approximation algorithm achieves an approximation ratio of $(k+3)/5$ for $k \geq 3$ (see [BF94, BF95]). Berman and Karpinski [BK03] showed an inapproximability factor of $\frac{98}{97}$ for MAXIMUM 3-DIMENSIONAL MATCHING. For more on low-degree inapproximability results see [Haz02].

1.2 Our Contribution.

We improve the inapproximability factor for MAXIMUM k -SET PACKING, and prove that it is NP-hard to approximate k -SP to within $\Omega\left(\frac{k}{\ln k}\right)$. We extend this result to MAXIMUM k -DIMENSIONAL MATCHING. These results also implies the same bound for $(k + 1)$ -claw-free graphs (see [Hal98] for definition of this problem and relation to k -SP). They do not hold, however, for k -IS.

The proof of these lower bound introduces a combinatorial object called Hypergraph-Disperser, and we present a randomized construction of such an object. This object may be of independent interest.

1.3 Outline

Some preliminaries are given in section 2. Section 2.2 presents the notion of hyper-graph-dispersers. Section 3 contains the proof of the asymptotic hardness of approximation for k -SP. Section 4 extends the proof to hold for k -DM. The existence of a good hyper-disperser is proved in section 5. The optimality of its parameters is shown in section 5.1. Section 6 contains a discussion on the implications of our results, the techniques used and some open problems.

2 Preliminaries

In order to prove inapproximability of a maximization problem, one usually defines a corresponding gap problem.

Definition 2 (Gap problems) *Let A be a maximization problem. $\text{gap-}A$ - $[a, b]$ is the following decision problem:*

Given an input instance, decide whether

- *there exists a solution of fractional size at least b , or*
- *every solution of the given instance is of fractional size smaller than a .*

If the size of the solution is located between these values, then the output is unconstrained.

Clearly, for any maximization problem, if $\text{gap-}A$ - $[a, b]$ is NP-hard, then it is NP-hard to approximate A to within any factor smaller than $\frac{b}{a}$.

Our main result in this paper is derived by a reduction from the following problem.

Definition 3 (Linear Equations) *MAX-3-LIN- q is the following optimization problem:*

Input: *A set Φ of linear equations modulo an integer q , each depending on 3 variables.*

Problem: *Find an assignment that satisfies the maximum number of equations.*

The following central theorem stems from a long line of research, using the PCP theorem (see [ALM⁺98, AS98]) and the parallel repetition theorem [Raz98] as a starting point:

Theorem 4 (Håstad [Hås01]) *For every $q \in \mathbb{N}$ and $\varepsilon > 0$ gap-MAX-3-LIN- q - $[\frac{1}{q} + \varepsilon, 1 - \varepsilon]$ is NP-Hard. Furthermore, the result holds for instances of MAX-3-LIN- q in which the number of occurrences of each variable is a constant (depending on ε and on q).*

We denote an instance of MAX-3-LIN- q by $\Phi = \{\varphi_1, \dots, \varphi_m\}$. Φ is over the set of variables $X = \{x_1, \dots, x_n\}$. Let $\Phi(x)$ be the (multi) set of all equations in Φ depending on $x \in X$ (i.e. it can be seen as all the occurrences of x). Denote by $\text{Sat}(\Phi, A)$ the set of all equations in Φ satisfied by the assignment A . For an assignment A , we denote by $A|_x$ the value in $a \in [q]$ that A assigns to x .

2.1 Hyper graphs

A hyper-graph $H = (V, E)$ consists of a set of vertices V and a collection E of subsets of V called hyper-edges (in short, edges).

As usual, the degree of a vertex is the number of edges it appears in. A hyper-graph H is called d -regular if the degree of each of its vertices is exactly d . H is called k -uniform if the size of each of its edges is exactly k .

A matching M is a subset of E , such that all edges of M are pairwise disjoint.

We use the following non-standard definition of an independent set in hyper-graphs:

Definition 5 (Independent-Set) *Let $H = (V, E)$ be a hyper-graph. A subset of vertices $I \subseteq V$ is called an Independent-Set if any edge $e \in E$ contains at most one vertex from I .*

And from it derive the corresponding (but non-standard) definition of colorability:

Definition 6 (strong colorability) *The hyper-graph $H = (V, E)$ is said to be k -strongly-colorable if there is a partition of V into k sets, such that each part is an Independent-Set.*

Hence, a k -uniform k -partite hyper-graph H may be denoted by $H = (V^1, V^2, \dots, V^k, E)$, where we have $E \subseteq V^1 \times \dots \times V^k$. An analogous notion to strong colorability applies to the edges of a hyper-graph:

Definition 7 (strong edge colorability) *A hyper-graph $H = (V, E)$ is said to be d -strongly-edge-colorable if there exists a coloring of the edges $f : E \mapsto [d]$ so that every vertex participates in at most one edge of each color.*

Using these definitions we can formally define the related packing problem studied hereby:

Definition 8 (MAXIMUM k -DIMENSIONAL MATCHING) ***k -Dimensional Matching** is the following optimization problem:*

Input: *A k -uniform k -strongly colorable hyper-graph $H = (V^1, \dots, V^k, E)$.*

Problem: *Find a matching of maximum size in H .*

2.2 Hyper Dispersers

The following definition is a generalization of disperser graphs. For definitions and results regarding dispersers see [RTS00].

Definition 9 ((q, δ) -Hyper-Edge-Disperser) *We call a hyper-graph $H = (V, E)$ a (q, δ) -Hyper-Edge-Disperser if there exists a partition of its edges: $E = E_1 \cup \dots \cup E_q$, $|E_1| = \dots = |E_q|$, such that every large matching M of H is (almost) concentrated in one part of the edges. Formally, there exists i so that*

$$|M \setminus E_i| \leq \delta |E|$$

Lemma 10 *For every $q > 1$ and $t > c(q)$ (where $c(q)$ is a constant depending only on q) there exists a hyper-graph $H = (V, E)$ such that*

- $V = [t] \times [d]$, whereas $d = \Theta(q \ln q)$.
- H is $(q, \frac{1}{q^2})$ -hyper-edge-disperser

- H is d uniform, d -strongly-colorable.
- H is q regular, q -strongly-edge-colorable.

Henceforth we denote such a graph by $\mathcal{D}[t, q]$. Because of the regularity, uniformity and colorability conditions on this hyper-graph, the number of edges is exactly $q \cdot t$, and they can be partitioned into q disjoint color sets. We therefore name its edges $e[i, j]$ where $j \in [q]$ is the color of the edge by an arbitrary strong edge coloring (a coloring where no two edges of the same color share a vertex) and $i \in [t]$ is an arbitrary indexing of the t edges of each color. Note that the t edges of any single color, cover all the vertices of $\mathcal{D}[t, q]$.

A proof of the above lemma appears in section 5. Note that $\mathcal{D}[t, 2]$ is the dual graph of a standard disperser.

3 Proof of the Asymptotic Inapproximability Factor for k -SP

This section provides a polynomial time reduction from MAX-3-LIN- q to k -SP. Given an instance Φ of MAX-3-LIN- q , that is a set of equations modulo integer q , we construct an instance of k -SP, namely a k -uniform hyper-graph.

For the hyper-graph we construct, we add hyper-edges corresponding to the equations of Φ and satisfying assignments to them. The main idea of the reduction is to construct the hyper-graph in such a way that a large matching corresponds to a consistent satisfying assignment to Φ . For this purpose, the hyper-graph has common vertices for edges that correspond to assignments that are inconsistent.

In general, the sparsity and uniformity of the constructed graph are strongly related to the quality of the hardness result. In order to obtain a sparse graph with small edge size, while retaining edge-intersection properties, we utilize the hyper-edge-disperser graphs defined in the previous section.

3.1 The construction

Let $\Phi = \{\varphi_1, \dots, \varphi_n\}$ be an instance of MAX-3-LIN- q over the set of variables X , where each variable $x \in X$ occurs a constant number of times $c(x) = O(1)$ (as in theorem 4). Let us now describe how to construct, in polynomial time, an instance of k -SP - the hyper-graph $H_\Phi = (V, E)$.

Let us assume, for every variable $x \in X$, a one-to-one mapping between all indices $i_x \in [c(x)]$ and all occurrences of x in Φ .

Recall lemma 10, that asserts the hyper-edge-disperser $\mathcal{D}[t, q]$ (definition 9) exists. For every vertex $x \in X$ let us consider the graph $\mathcal{D}[c(x), q]$. Each vertex in $V(\mathcal{D}[c(x), q])$ corresponds to an occurrence of x in Φ , and a number in $[d]$, where $d = \Theta(q \log q)$. Each of the edges in $E(\mathcal{D}[c(x), q])$ is in a one-to-one correspondence with an occurrence $i_x \in [c(x)]$ and a value $a \in [q]$ according to the strong-edge-coloring of $\mathcal{D}[c(x), q]$, thus let us denote these edges by $e\langle x, i_x, a \rangle$.

The set of vertices V of H_Φ consists of one copy of the vertices of the disperser graph $\mathcal{D}[c(x), q]$ for every variable $x \in X$. Namely,

$$V \triangleq \{v\langle x, i, j \rangle \mid x \in X, i \in [c(x)], j \in [d]\}$$

Henceforth, for some variable $x \in X$, the copy of $\mathcal{D}[c(x), q]$ over the vertices $V_x \triangleq \{v\langle x, i, j \rangle \mid i \in [c(x)], j \in [d]\}$, shall be denoted \mathcal{D}_x .

Let us now define the set of edges E of H_Φ . The edges of H_Φ will be composed of several edges from the hyper-graphs \mathcal{D}_x . The set E consists of one edge for every equation $\varphi \in \Phi$ over variables x, y, z and assignment A to x, y, z that satisfies φ . Denote by $A|_x, A|_y, A|_z \in [q]$ the values that A assigns to these variables respectively (notice there are q^2 such satisfying assignments). Denote by i_x, i_y, i_z the indices of the occurrences of x, y, z respectively in φ . The edge corresponding to φ and A is a union of three edges from the copies of hyper-graphs $\mathcal{D}_x, \mathcal{D}_y, \mathcal{D}_z$ of the variables in the equation φ :

$$e\langle\varphi, A\rangle = e\langle x, i_x, A|_x\rangle \cup e\langle y, i_y, A|_y\rangle \cup e\langle z, i_z, A|_z\rangle$$

Clearly, the cardinality of each edge $e\langle\varphi, A\rangle$ is $3d$, as it is the disjoint union of three edges of cardinality d . Note that each of the three edges $e\langle x, i_x, A|_x\rangle, e\langle y, i_y, A|_y\rangle, e\langle z, i_z, A|_z\rangle$, which compose $e\langle\varphi, A\rangle$, participates in q edges of the hypergraph H_Φ .

Altogether, the edges of H_Φ are

$$E = \{e\langle\varphi, A\rangle \mid \varphi \in \Phi, A \text{ is a satisfying assignment to } \varphi\}$$

This concludes the construction of the k -SP instance H_Φ .

Notice that for every constant q , the construction can be carried out in deterministic polynomial time. To this end, each disperser \mathcal{D}_x should be constructed in deterministic polynomial time. As each variable x occurs $c(x) = O(1)$ times (by Theorem 4), the size of \mathcal{D}_x is constant as well. According to lemma 10, we know that \mathcal{D}_x exists. Therefore, we can enumerate all possible hyper-graphs of the required size and verify whether they are indeed hyper-edge-dispersers with the required parameters.

3.2 Proof of correctness

We next show that the size of a maximum matching in H_Φ is proportional to the maximum number of equations of Φ that can be simultaneously satisfied. That is, if there exists an assignment that satisfies almost all equations of Φ then there exists a matching that covers almost all vertices of H_Φ . On the other hand, if every assignment satisfies at most a small fraction of the equations of Φ , then every matching of H_Φ is small.

Lemma 11 *[Completeness] If there is an assignment to Φ which satisfies $1 - \varepsilon$ of its equations, then there is a matching in H_Φ of size $\left(\frac{1-\varepsilon}{q^2}\right) |E|$.*

PROOF: Let $A: X \mapsto [q]$ be an assignment that satisfies $1 - \varepsilon$ of the equations. Consider the matching $M \subseteq E$ consisting of all edges corresponding to A , namely

$$M = \{e\langle\varphi, A\rangle \mid \varphi \in \text{Sat}(\Phi, A)\}$$

Since M contains one edge corresponding to each satisfied equation, and for each equation there are q^2 satisfying assignments, we have $|M| = \left(\frac{1-\varepsilon}{q^2}\right) |E|$. To see that these edges are indeed a matching, consider any two edges of M . If they do not relate to the same variables then they do not contain vertices from a joint hyper-edge-disperser. On the other hand, if they do relate to a joint variable $x \in X$, then they relate to different occurrences $i_{x,1}, i_{x,2} \in [c(x)]$, but the same assignment $a \in [q]$ to it. Hence they contain vertices of the same hyper-edge-disperser \mathcal{D}_x , but from two distinct edges of the same color, therefore they do not share a vertex. \square

Lemma 12 [Soundness] *If every assignment to Φ satisfies at most $\frac{1}{q} + \varepsilon$ fraction of its equations, then every matching in H_Φ is of size $O\left(\frac{1}{q^3}|E|\right)$.*

The proof idea is as follows: given a matching M , each edge in it corresponds to an assignment to three variables. Given a matching M , we use it to define a global-assignment A_{maj} to the variables of Φ : every variable is assigned the value which agrees with the maximal number of the hyper-edges of M . We then partition the edges of M into two sets: those that agree with the global assignment (named M_{maj}) and the complement set (named M_{min}). The size of M_{maj} is bounded, as it corresponds to the set equations satisfied by A_{maj} (which is small). We then proceed to bound the size of M_{min} using the expansion properties of the hyper-edge-dispersers.

PROOF: Denote by E_x the edges of H_Φ corresponding to equations that contain the variable x , namely,

$$E_x = \{e\langle\varphi, A\rangle \mid \varphi \in \Phi(x), e\langle\varphi, A\rangle \in E\}$$

Denote by $E_{x=a}$ the subset of E_x corresponding to an assignment of a to x , that is,

$$E_{x=a} = \{e\langle\varphi, A\rangle \mid e\langle\varphi, A\rangle \in E_x, A|_x = a\}$$

Let M be a matching of maximum size in H_Φ . According to the matching M we define the majority assignment A_{maj} as follows: for every $x \in X$, the assignment $A_{maj}(x)$ is the value $a \in [q]$ such that $|E_{x=a} \cap M|$ is maximized. Let M_{maj} be the set of edges in M that agree with A_{maj} , and M_{min} be all the other edges in M , namely

$$M_{maj} = M \cap \{e\langle\varphi, A_{maj}\rangle \mid \varphi \in \Phi\}$$

$$M_{min} = M \setminus M_{maj}$$

As the number of equations satisfied by the assignment A_{maj} satisfies $|Sat(\Phi, A_{maj})| \leq \frac{1}{q} + \varepsilon$, and for each equation there are q^2 edges corresponding to all satisfying assignments for this equation, we have:

$$|M_{maj}| < \left(\frac{1}{q} + \varepsilon\right) \frac{|E|}{q^2} \quad (1)$$

We next bound the size of M_{min} . The idea is as follows: we decompose each of edges in M_{min} into the three constructing edges. At least one of those three edges corresponds to an assignment to the variable which is not the most popular one. We then exploit the disperser property to bound their number.

Consider a certain variable $x \in X$. Then \mathcal{D}_x is a $(q, \frac{1}{q^2})$ -hyper-edge-disperser (recall definition 9). That is, in any subset of edges of \mathcal{D}_x which is a matching, all but at most a fraction $\frac{1}{q^2}$ of the edges are of one color (which corresponds to a single assignment to the variable x). Clearly, if two edges of \mathcal{D}_x intersect, then so do any pair of edges of H_Φ containing these two edges. Therefore,

$$\sum_{a \neq A_{maj}(x)} |M_{min} \cap E_{x=a}| \leq \frac{1}{q^2} |E(\mathcal{D}_x)| \quad (2)$$

where $|E(\mathcal{D}_x)|$ is the number of edges of \mathcal{D}_x .

Consider an edge $e\langle x, i_x, a \rangle$ of \mathcal{D}_x , where i_x is the index of x when it appears in the equation $\varphi \in \Phi$. This edge was used in the construction for q edges of H_Φ , namely those edges that correspond to the satisfying assignments of φ that assign to x the value a .

Hence, every edge of \mathcal{D}_x is a subset of q hyper edges in E_x . However, no more than one of these q edges may participate in M (as M is a matching). Plugging this observation to (2) we obtain:

$$\sum_{a \neq A_{maj}(x)} |M_{min} \cap E_{x=a}| \leq \frac{1}{q^3} |E_x| \quad (3)$$

Summing up over the variables of Φ we obtain:

$$|M_{min}| \leq \sum_{x \in X, a \neq A_{maj}(x)} |M_{min} \cap E_{x=a}| \leq \frac{1}{q^3} \sum_{x \in X} |E_x| = \frac{3}{q^3} |E| \quad (4)$$

where the last equality follows from the fact that each equation contains three variables. Thus, from (1) and (4):

$$|M| = |M_{min}| + |M_{maj}| \leq \left(\frac{4}{q^3} + \varepsilon\right) |E|$$

□

By lemmas 11 and 12 we showed that $\text{Gap-}k\text{-SP-}\left[\frac{4}{q^3} + \varepsilon, \frac{1}{q^2} - \varepsilon\right]$ is NP-hard. Since each edge is of size $k = 3d = \Theta(q \log q)$ it is NP-hard to approximate $k\text{-SP}$ to within $\Omega(\frac{k}{\ln k})$.

4 Extending the Proof for $k\text{-DM}$

The proof for $k\text{-DM}$ is similar to the $k\text{-SP}$ proof, yet we take additional care to ensure that the graph H_Φ we construct has the required structure (namely, that H_Φ is not only k -uniform, but is also k -strongly-colorable).

The construction for $k\text{-DM}$ takes into account the *location* of the variables in the equations they appear in. As there are three variables per equation, there are three possible locations. We use the following notation: $\Phi(x, l)$ to be the subset of $\Phi(x)$ where x is the l 'th variable in the equation ($l \in [3]$). W.l.o.g we may assume that for every variable, it appears the same number of times in every location ($\Phi(x, 1) = \Phi(x, 2) = \Phi(x, 3)$), as we can take three copies of each equation, and shift the location of the variables.

Similar to the $k\text{-SP}$ construction, we associate a vertex with each appearance of a variable. For every variable $x \in X$, we now have three copies of a hyper-edge disperser (instead of just one we had for $k\text{-SP}$) - a different disperser for each location in the equations. For every location $l \in [3]$, we have a hyper disperser $\mathcal{D}[\frac{c(x)}{3}, q]$ which is denoted by $\mathcal{D}_{x,l}$. The vertices of H_Φ are the union of the vertices of all these hyper-dispersers corresponding to all variables in the equation set and all locations.

Since $c(x)$ is exactly the number of appearances of the variable x in the equation set Φ , we can enumerate the vertices of H_Φ according to the variable $x \in X$ and equation $\varphi \in \Phi$ they correspond to (and these two parameters determine the location of the variable in the equation as well):

$$V = \{v\langle x, \varphi, j \rangle \mid x \in X, \varphi \in \Phi(x), j \in [d]\}$$

The construction of the edges of H_Φ is almost identical to that of the k -SP instance, the differences being the distinction between the three dispersers for each variable. Notice there is a bijection between an occurrence of a variable in a certain equation, and the corresponding vertex in one of the three hyper-graphs corresponding to this variable. Therefore, there is no ambiguity in the edge construction process, which is otherwise identical to the one for k -SP.

The notation we use for the edges is identical to the k -SP construction as well: the edges correspond to the satisfying assignments to the equations, and are composed of three disperser edges each $e\langle\varphi, A\rangle = e\langle x, i_x, A|_x\rangle \cup e\langle y, i_y, A|_y\rangle \cup e\langle z, i_z, A|_z\rangle$ (where these three edges are taken from $\mathcal{D}_{x,1}$, $\mathcal{D}_{y,2}$ and $\mathcal{D}_{z,3}$, respectively). The set of all edges is denoted:

$$E = \{e\langle\varphi, A\rangle \mid \varphi \in \Phi, A \text{ is a satisfying assignment to } \varphi\}$$

This concludes the construction for k -DM. We first show that the graph constructed is indeed a k -DM instance:

Proposition 13 *H_Φ is $3d$ -strongly-colorable.*

PROOF: We show how to partition V into $3d$ independent sets of equal size. Let the sets be $P_{l,i}$ whereas $i \in [d]$ and $l \in [3]$:

$$P_{l,i} = \{v\langle x, \varphi, i\rangle \mid x \in X, \varphi \in \Phi(x, l)\}$$

$P_{l,i}$ is clearly a partition of the vertices, as each vertex belongs to at least one part.

We now explain why each part is an independent set. Let $P_{l,i}$ be an arbitrary part, and let $e\langle\varphi, A\rangle \in E$ be an arbitrary edge, where $\varphi \equiv x + y + z = a \pmod q$. By construction, this edge is a disjoint union of three hyper-disperser edges corresponding to the three variables x, y, z , namely,

$$e\langle\varphi, A\rangle = e\langle x, i_x, A|_x\rangle \cup e\langle y, i_y, A|_y\rangle \cup e\langle z, i_z, A|_z\rangle$$

$P_{l,i} \cap e\langle\varphi, A\rangle$ may contain vertices corresponding only to one of the variables x, y, z , since it contains variables corresponding to a single location (first, second or third). Let that variable be, w.l.o.g, x . Since the hyper-graph $D_{x,1}$ is d -partite and d -strongly-edge-colorable, the edge $e\langle x, i_x, A|_x\rangle$ (and hence $e\langle\varphi, A\rangle$) contains exactly one vertex from each of the d parts. Therefore, the set $P_{l,i} \cap e\langle\varphi, A\rangle$ contains exactly one vertex. Since $|P_{l,i} \cap e\langle\varphi, A\rangle| = 1$ for every edge and every set $P_{l,i}$, the graph H_Φ is $3d$ -strongly-colorable. \square

We proceed to prove the gap in maximum matching size between the cases in which the equation set Φ is almost satisfiable and very unsatisfiable. For the case in which Φ is almost satisfiable, the completeness lemma (lemma 11) is valid for the current construction as well. We prove the appropriate soundness lemma, which is very similar to lemma 12.

Lemma 14 *[Soundness] If every assignment to Φ satisfies at most $\frac{1}{q} + \varepsilon$ fraction of its equations, then every matching in H_Φ is of size $O\left(\frac{1}{q^3}|E|\right)$.*

PROOF: Denote by $E_{x,l}$ the edges of H_Φ corresponding to equations φ containing the variable x in location l , namely,

$$E_{x,l} = \{e\langle\varphi, A\rangle \mid \varphi \in \Phi(x, l), A \in [q^2]\}$$

Denote by $E_{x=a,l}$ the subset of $E_{x,l}$ corresponding to an assignment of $a \in [q]$ to x , that is,

$$E_{x=a,l} = \{e\langle\varphi, A\rangle \mid \varphi \in \Phi(x, l), A|_x = a\}$$

Let M be a matching of maximum size in H_Φ . According to the matching M we define the majority assignment A_{maj} , taking into account the locations of the variables, as follows: for every $x \in X$, let $\hat{l}(x)$ be the location for which $|E_{x,\hat{l}(x)} \cap M|$ is maximized. The assignment $A_{maj}(x)$ is the value $a \in [q]$ such that $|E_{x=a,\hat{l}(x)} \cap M|$ is maximized.

As before, let M_{maj} be the set of edges in M that agree with A_{maj} , and M_{min} be all the other edges in M , namely

$$M_{maj} = M \cap \{e \langle \varphi, A_{maj} \rangle \mid \varphi \in \Phi\}$$

$$M_{min} = M \setminus M_{maj}$$

For the exact same reasons as in the previous soundness proof (lemma 12), we have the analogous equations to 1 and 3:

$$|M_{maj}| < \left(\frac{1}{q} + \varepsilon\right) \frac{|E|}{q^2} \quad (5)$$

$$\forall x \in X \quad \sum_{a \neq A_{maj}(x)} |M_{min} \cap E_{x=a,\hat{l}(x)}| \leq \frac{1}{q^3} |E_{x,\hat{l}(x)}| \quad (6)$$

And proceeding according to the lines of lemma 12,

$$\begin{aligned} |M| &\leq \sum_{x,l} |M \cap E_{x,l}| \\ &\leq \sum_{x,l} |M_{maj} \cap E_{x,l}| + \sum_{x,l,a \neq A_{maj}(x)} |M_{min} \cap E_{x=a,l}| \end{aligned}$$

since $|E_{x,\hat{l}(x)} \cap M|$ is maximized by $\hat{l}(x)$, we have

$$\begin{aligned} &\leq 3 \cdot \sum_x |M_{maj} \cap E_{x,\hat{l}(x)}| + 3 \cdot \sum_{x,a \neq A_{maj}(x)} |M_{min} \cap E_{x=a,\hat{l}(x)}| \\ &\leq 3 \cdot |M_{maj}| + 3 \cdot \sum_{x,a \neq A_{maj}(x)} |M_{min} \cap E_{x=a,\hat{l}(x)}| \end{aligned}$$

plugging in (5) and (6)

$$\begin{aligned} &< 3\left(\frac{1}{q} + \varepsilon\right) \frac{|E|}{q^2} + \frac{3}{q^3} \sum_x |E_{x,\hat{l}(x)}| \\ &\leq \left(\frac{4}{q^3} + 3\varepsilon\right) |E| \end{aligned}$$

□

Finally, by the completeness lemma from the previous section and lemma 14 we conclude that $\text{Gap-}k\text{-DM-}\left[\frac{4}{q^3} + 3\varepsilon, \frac{1}{q^2} - \varepsilon\right]$ is NP-hard, thus it is NP-hard to approximate $k\text{-DM}$ to within $\Omega(\frac{k}{\ln k})$.

5 Hyper-Dispersers

In this section, we prove lemma 10. As stated before, these are generalizations of disperser graphs. In section 5.1, we prove that these are the best (up to a constant) parameters for a hyper-disperser one can hope to achieve.

Lemma 10 *For every $q > 1$ and $t > c(q)$ (where $c(q)$ is a constant depending only on q) there exists a hyper-graph $H = (V, E)$ such that*

- $V = [t] \times [d]$, whereas $d = \Theta(q \ln q)$.
- H is $(q, \frac{1}{q^2})$ -hyper-edge-disperser
- H is d uniform, d -strongly-colorable.
- H is q regular, q -strongly-edge-colorable.

We denote this graph by $\mathcal{D}[t, q]$.

PROOF: We follow the probabilistic method to prove that the probability that a randomly generated graph is not a $\mathcal{D}[t, q]$ graph, is strictly smaller than 1, from which follows the existence of such graphs. Let

$$V = [t] \times [d]$$

and denote $V_i = [t] \times \{i\}$.

We next randomly construct the edges of the hyper-graph, so that it is d -uniform and q -regular. Let S_t be the set of all permutation over t elements. For every $(i_1, i_2) \in [q] \times [d]$ choose a permutation from S_t uniformly at random:

$$\Pi_{i_1, i_2} \in_R S_t$$

Define

$$e[i, j] = \{ (\Pi_{j,1}(i), 1), (\Pi_{j,2}(i), 2), \dots, (\Pi_{j,d}(i), d) \} \quad (7)$$

and let

$$E = \{e[i, j] \mid (i, j) \in [t] \times [q]\}$$

Hence $|E| = tq$. Define a partition of the edges as follows: $E_i = \{e[j, i] \mid j \in [t]\}$. Thus $|E_1| = \dots = |E_q| = t$ and each set of edges E_j covers every vertex exactly once. Therefore, H is q strongly-edge-colorable. On the other hand, every edge contains exactly one vertex from each set of vertices V_i . Thus H is d -strongly-colorable.

We next show that with high probability H has the disperser property, namely, every matching M of H is concentrated on a single part of the edges, except for maybe $\frac{1}{q^2}|E| = \frac{t}{q}$ edges of M . Denote by P the probability that H does *not* have the disperser property.

Definition 15 *Define \mathcal{M}_k as the following family:*

$$\mathcal{M}_k = \{M \mid M \subseteq E, |M \cap E_k| \leq \frac{t}{q}, |M \setminus E_k| = \frac{t}{q}, \forall i |M \cap E_k| \geq |M \cap E_i|\}$$

Proposition 16 *If H is not a $(q, \frac{1}{q^2})$ -hyper-edge-disperser, then there exists a $k \in [q]$ and a set $M \in \mathcal{M}_k$ that is a matching.*

PROOF: Suppose that H is not a $(q, \frac{1}{q^2})$ -hyper-edge-disperser. Namely, there exists a matching $M' \subseteq E$ such that it is not concentrated on one color of edges: $\forall i, |M' \setminus E_i| > \frac{1}{q^2}|E| = \frac{t}{q}$. Let $k \in [q]$ be such that $|M' \cap E_k|$ is maximal. As any subset of a matching is a matching, we can remove edges from $M' \setminus E_k$ until we are left with exactly $\frac{t}{q}$ edges. Likewise, we can remove edges of $M' \cap E_k$ until this set contains at most $\frac{t}{q}$ edges. Notice that this new set is a matching in \mathcal{M}_k . \square

Following the above proposition, we proceed with the proof considering only sets in \mathcal{M}_1 . Denote by $\Pr[M]$ the probability (over the random choice of H) that M is a matching. By union bound, symmetry with respect to k , and the above proposition,

$$\begin{aligned} P &\leq \Pr_H[\exists k, M \in \mathcal{M}_k, M \text{ is a matching}] \leq \\ &\leq q \sum_{M \in \mathcal{M}_1} \Pr[M] \leq q|\mathcal{M}_1| \Pr[\hat{M}] \end{aligned} \quad (8)$$

where $\hat{M} \in \mathcal{M}_1$ is the set which maximizes $\Pr[M]$. The size of \mathcal{M}_1 is bounded from above by the number of possibilities to choose $\frac{t}{q}$ edges from E_1 and another $\frac{t}{q}$ edges from the rest of the edge color sets. Therefore, using the known inequality $\binom{n}{k} \leq (\frac{en}{k})^k$ and assuming $t \gg q \gg 2$ we obtain:

$$|\mathcal{M}_1| \leq \binom{(q-1)t}{\frac{t}{q}} \binom{t}{\frac{t}{q}} \leq (eq^2)^{\frac{t}{q}} (eq)^{\frac{t}{q}} \leq (eq)^{\frac{3t}{q}} \quad (9)$$

We next bound $\Pr[\hat{M}]$. Denote by \hat{M}_i the event that \hat{M} restricted to the vertices of V_i is a matching (that is, the edges of \hat{M} do not share a vertex in V_i). According to the independent choice of permutations in the construction of H (recall (7)), the events \hat{M}_i are independent and identically distributed. Hence,

$$\Pr[\hat{M}] = \prod_{i=1}^d \Pr[\hat{M}_i] \quad (10)$$

and we proceed to bound $\Pr[\hat{M}_1]$. Henceforth we shall only consider vertices of V_1 .

Let M_i be the set of edges in $\hat{M} \cap E_i$ restricted to the vertices of V_1 . Let A_i be the event that the sets of edges $\{M_j | j \leq i\}$ are all disjoint. Then:

$$\Pr[\hat{M}_1] = \Pr[\cap_{i=2}^q A_i] = \prod_{i=2}^q \Pr[A_i | A_{i-1}] \quad (11)$$

The probability of the event $A_i | A_{i-1}$ is the probability of picking at random $|M_i|$ different vertices from a set of t vertices (the set V_1), and avoiding all vertices from the set $\cup_{l=1}^{i-1} M_l$. Naturally, this probability is smaller than the probability of picking $|M_i|$ vertices from a set of t vertices with repetition (one is allowed to choose the same vertex more than once). The assumption A_{i-1} implies that the sets M_l for all $l < i$ are disjoint, and hence $|\cup_{l=1}^{i-1} M_l| = \sum_{l=1}^{i-1} |M_l|$. Therefore,

$$\Pr[A_i | A_{i-1}] \leq (1 - \frac{\sum_{l < i} |M_l|}{t})^{|M_i|} \leq e^{-\frac{1}{t} |M_i| \sum_{l < i} |M_l|}$$

Where for the last inequality we used $1 - x \leq e^{-x}$. Thus by equations (10) and (11) we have:

$$\Pr[\hat{M}] \leq e^{-\frac{d}{t} \sum_{i=2}^q (|M_i| \sum_{j=1}^{i-1} |M_j|)} = e^{-\frac{d}{t} \sum_{i < j} |M_i| |M_j|} \quad (12)$$

For ease of notation, denote $x_i = |M_i|$. We would like an upper bound on the previous probability, and that is obtained when the term $\sum_{i < j} |M_i| |M_j|$ is minimized. In our case, the constraint that $\hat{M} \in \mathcal{M}_1$ implies that $|M_1| \geq \max_{i=2}^q |M_i|$; $\sum_{i=2}^q |M_i| = \frac{t}{q}$. Lemma 20 in the appendix shows that the minimum of this expression under these constraints is at least $\frac{t^2}{4q^2}$. Therefore, from equation 12 we obtain the following bound on the probability:

$$\Pr[\hat{M}] \leq e^{-\frac{d}{t} \cdot \frac{t^2}{4q^2}} = e^{-\frac{dt}{4q^2}} \quad (13)$$

Therefore by equations (8),(9),(13),

$$P \leq q(eq)^{\frac{3t}{q}} e^{-\frac{dt}{4q^2}}$$

Any d which guarantees that $q(eq)^{\frac{3t}{q}} e^{-\frac{dt}{4q^2}} \ll 1$ (for example $d \geq 20q \ln q$) suffices to conclude that $P < 1$, and therefore that there exists H with the disperser properties.

□

5.1 Optimality of Hyper-Disperser Construction

We now turn to see why the hyper-disperser from lemma 10 has optimal parameters. We base our observation on the following lemma from [RTS00]:

Definition 17 *A bipartite graph $G = (V_1, V_2, E)$ is called a δ -disperser if for every $U_1 \subseteq V_1, U_2 \subseteq V_2, |U_1|, |U_2| \geq \delta|V_1| = \delta|V_2|$, the subset $U_1 \cup U_2$ is not an independent set.*

Lemma 18 (RTS) *Every bipartite d -regular $\frac{1}{k}$ -disperser must satisfy $d = \Omega(k \ln k)$.*

Using this lemma we prove:

Lemma 19 *Every d -uniform, d -strongly colorable, q -regular, q -strongly-edge-colorable, $(q, \frac{1}{q^2})$ -hyper-edge-disperser, must satisfy $d = \Omega(q \ln q)$.*

PROOF: We prove that if there exists such a hyper-graph which satisfies $d = o(q \ln q)$, then there exists a bipartite $o(q \ln q)$ -regular $\frac{1}{q}$ -disperser, in contradiction to lemma 18. We transform a d -partite d -uniform q -regular q -strongly-edge colorable $(q, \frac{1}{q^2})$ -hyper-disperser $H = (V_H, E_1, E_2, \dots, E_q)$ into a bipartite d -regular $\frac{1}{q}$ -disperser $G = (V_1, V_2, E_G)$ in the following way. Let

$$\begin{aligned} V_1 &= E_1 \\ V_2 &= E_2 \\ E_G &= \{(e_1, e_2) \mid e_1 \cap e_2 \neq \emptyset\} \end{aligned}$$

Obviously G is a bipartite d -regular graph (we allow multi-edges). In addition, suppose that two sets of fractional sizes:

$$S_1 = \frac{1}{q} V_1, S_2 = \frac{1}{q} V_2$$

are an independent set in G . The corresponding sets of edges in H are disjoint and are of fractional size $\frac{2}{q^2}$, thus contradicting the fact that H is a $(q, \frac{1}{q^2})$ -hyper-disperser. □

6 Discussion

An interesting property of our construction is the *almost perfect completeness*. This property refers to the fact that the matching proved to exist in the completeness lemma 11 is an almost perfect matching. That is, it covers $1 - \varepsilon$ of the vertices. Knowing the location of a gap is interesting by itself and may prove useful (in particular if it is extreme on either the completeness or the soundness parameters, see for example [Pet94]). In fact, applying our reduction on other PCP variants instead of Max-3-Lin-q (e.g. parallel repetition of 3-SAT) yields perfect completeness for k -SP and for k -DM (but with weaker hardness factors).

The ratio between the asymptotic inapproximability factor presented herein for k -SP and k -DM, and the tightest approximation algorithm known, was reduced to $O(\ln k)$. The open question of where in the range, from $\frac{k}{2}$ to $O(\frac{k}{\ln k})$ is the approximability threshold, is interesting by itself, as well as its implications to the difference between k -DM and k -IS. The current asymptotic inapproximability factor of $\Omega(\frac{k}{\ln k})$ for k -DM approaches the tightest approximation ratio known for k -IS, namely $O(k \log \log k / \log k)$ [Vis96]. Thus, a small improvement in either the approximation ratio or the inapproximability factor will show these problems to be of inherently different complexity.

7 Acknowledgements

We would like to thank Adi Akavia and Dana Moshkovitz for their helpful comments.

References

- [AFWD95] N. Alon, U. Fiege, A. Wigderson, and D. Zuckerman. Derandomized graph products. *Computational Complexity*, 5:60–75, 1995.
- [ALM⁺98] S. Arora, C. Lund, R. Motwani, M. Sudan, and M. Szegedy. Proof verification and intractability of approximation problems. *Journal of the ACM*, 45:501–555, 1998.
- [AS98] S. Arora and S. Safra. Probabilistic checking of proofs: A new characterization of NP. *Journal of the ACM*, 45:70–122, 1998.
- [BF94] P. Berman and M. Furer. Approximating maximum independent set in bounded degree graphs. *SODA*, pages 365–371, 1994.
- [BF95] P. Berman and T. Fujito. On the approximation properties of independent set problem in degree 3 graphs. *WADS*, pages 449–460, 1995.
- [BH92] R. Boppana and M. M. Halldorsson. Approximating maximum independent sets by excluding subgraphs. *Bit* 32, pages 180–196, 1992.
- [BK03] P. Berman and M. Karpinski. Improved approximation lower bounds on small occurrence optimization. *ECCC TR03-008*, 2003.
- [BYM84] R. Bar-Yehuda and S. Moran. On approximation problems related to the independent set and vertex cover problems. *Discrete Applied Mathematics*, 9:1–10, 1984.
- [Edm65] J. Edmonds. Paths, trees and flowers. *Canadian Journal of Mathematics*, 17:449–467, 1965.
- [FGL⁺96] U. Feige, S. Goldwasser, L. Lovász, S. Safra, and M. Szegedy. Interactive proofs and the hardness of approximating cliques. *J. ACM*, 43(2):268–292, 1996.
- [Hal98] M. M. Halldorsson. Approximations of independent sets in graphs. *APPROX*, 1998.
- [Hås99] J. Håstad. Clique is hard to approximate within $n^{1-\epsilon}$. *Acta Math.*, 182(1):105–142, 1999.
- [Hås01] J. Håstad. Some optimal inapproximability results. *Journal of the ACM*, 48(4):798–859, July 2001.
- [Haz02] E. Hazan. On the hardness of approximating k-dimensional matching. Master’s thesis, Tel-Aviv University, 2002.
- [HS89] C. A. J. Hurkens and A. Schrijver. On the size of systems of sets every t of which have an sdr, with an application to the worst-case ratio of heuristics for packing problems. *SIAM Journal Discrete Math*, 2:68–72, 1989.
- [Kan91] V. Kann. Maximum bounded 3-dimensional matching is MAXSNP-complete. *Information Processing Letters*, 37:27–35, 1991.
- [Kar72] R. M. Karp. Reducibility among combinatorial problems. *Complexity of Computer Computations*, pages 83–103, 1972.

- [MM04] P. Sankowski M. Mucha. Maximum matchings via gaussian elimination. In *Proc. 45nd IEEE Symp. on Foundations of Computer Science*, pages 248–255, 2004.
- [Pap94] C. Papadimitriou. *Computational Complexity*. Addison Wesley, 1994.
- [Pet94] E. Petrank. The hardness of approximation: Gap location. *Computational Complexity*, 4(2):133–157, 1994.
- [Raz98] R. Raz. A parallel repetition theorem. *SIAM Journal on Computing*, 27(3):763–803, June 1998.
- [RTS00] J. Radhakrishnan and A. Ta-Shma. Bounds for dispersers, extractors, and depth-two superconcentrators. *SIAM Journal on Discrete Mathematics*, 13:2–24, 2000.
- [Tre01] L. Trevisan. Non-approximability results for optimization problems on bounded degree instances. *Proc. 33rd ACM Symp. on Theory of Computing*, 2001.
- [Vis96] S. Vishwanathan. Personal communication to M. Halldorsson cited in [Hal98]. 1996.
- [Wig83] A. Wigderson. Improving the performance guarantee for approximate graph coloring. *Journal of the Association for Computing Machinery*, 30(4):729–735, 1983.

A Additional Lemmas

Lemma 20 *Under the constraints $\forall_{i \in [m]} x_i \geq 0$ $x_1 \geq \max_{i=2}^q x_i$ $\sum_{i=2}^q x_i = T$, it holds that:*

$$\sum_{1 \leq i < j \leq q} x_i x_j \geq \frac{1}{4} T^2$$

PROOF: If $x_1 \geq \frac{T}{2}$, then we directly obtain:

$$\sum_{1 \leq i < j \leq q} x_i x_j \geq x_1 \cdot \sum_{i=2}^q x_i \geq \frac{T}{2} \cdot T \geq \frac{1}{4} T^2$$

Otherwise, we know that:

$$\sum_{1 \leq i < j \leq q} x_i x_j \geq \sum_{2 \leq i < j \leq q} x_i x_j = \frac{1}{2} \left(\sum_{i=2}^q x_i \right)^2 - \frac{1}{2} \sum_{i=2}^q x_i^2 \geq \frac{1}{2} T^2 - \sum_{i=2}^q x_i^2$$

The function $\sum_{i=2}^q x_i^2$ is convex, and hence under the constraints $\sum_{i=2}^q x_i = T$; $\max_{i=2}^q x_i \leq \frac{T}{2}$, it is maximized where $x_2 = x_3 = \frac{T}{2}$ and the rest of the variables are zero. We obtain that $\sum_{i=2}^q x_i^2 \leq \frac{1}{4} T^2$, and finally $\sum_{1 \leq i < j \leq q} x_i x_j \geq \frac{1}{2} T^2 - \sum_{i=2}^q x_i^2 \geq \frac{1}{4} T^2$. \square