

Partial Recovery of Erdős-Rényi Graph Alignment via k -Core Alignment

Daniel Cullina^{*1}, Negar Kiyavash^{†2}, Prateek Mittal^{‡1}, and H. Vincent Poor^{§1}

¹Princeton University, Department of Electrical Engineering

²Georgia Tech, Department of Electrical and Computer Engineering and Department of Industrial & Systems Engineering

Abstract

We determine information theoretic conditions under which it is possible to *partially* recover the alignment used to generate a pair of sparse, correlated Erdős-Rényi graphs. To prove our achievability result, we introduce the k -core alignment estimator. This estimator searches for an alignment in which the intersection of the correlated graphs using this alignment has a minimum degree of k . We prove a matching converse bound. As the number of vertices grows, recovery of the alignment for a fraction of the vertices tending to one is possible when the average degree of the intersection of the graph pair tends to infinity. It was previously known that exact alignment is possible when this average degree grows faster than the logarithm of the number of vertices.

Graph alignment, or graph matching, is the problem of finding a correspondence between the vertex sets of a pair of graphs using structural information from the graphs. It can be thought of as the noisy generalization of the graph isomorphism problem. Graph matching has applications in the privacy of social network data, the analysis of biological protein interaction networks, and in computer vision.

We consider the graph matching problem for random graphs that have been generated in a correlated way, so there is a planted ground-truth alignment of their vertices. In this setting, the combinatorial optimization problem of maximizing edge overlap is also the maximum a posteriori estimator.

0.1 Related work

A number of authors have worked to determine the information theoretic or statistical conditions under which graph alignments can be recovered by any algorithm. Wright determined the conditions under which an Erdős-Rényi graph has a trivial automorphism group, or equivalently under which the isomorphism recovery problem has a unique solution [1]. Pedarsani and Grossglauser obtained achievability conditions for exact recovery in the noisy case [2]. Cullina and Kiyavash obtained matching achievability and converse conditions for exact recovery [3, 4]. Kazemi, Yartseva, and Grossglauser considered alignment of graphs with overlapping but not identical vertex sets [5].

*dcullina@princeton.edu

†negar.kiyavash@ece.gatech.edu

‡pmittal@princeton.edu

§poor@princeton.edu

Shirani, Garg, and Erkip found an achievability condition for partial recovery with a small number of errors was obtained [6]. In all of these cases, the explicit or implicit algorithms require exponential time in the number of vertices. Cullina, Mittal, and Kiyavash obtained analogous limits for the alignment recovery problem for correlated databases [7]. In this case, maximum a posteriori estimation can be done efficiently.

A number of practically motivated and efficient algorithms have been proposed [8, 9, 10, 11, 12, 13, 14, 15, 16]. These have largely been empirically evaluated on a mix of real and synthetic datasets. It is common for these algorithms to return partial matchings of the vertex sets for some graph pairs.

A few efficient algorithms that require some form of initial side information have been rigorously analyzed. Yartseva and Grossglauser used a percolation algorithm to obtain a graph alignment starting with some matched pairs of seed vertices [17]. A number of other works have investigated seeded matching [18, 19, 20]. Feizi et al. used a spectral method to recover an alignment of dense graphs with side information restricting the set of possible alignments [21]. Lyzinski et al. explored the limitations of some convex programming methods, which have presented a barrier to the development of efficient algorithms [22].

Very recently, provably correct quasi-polynomial time algorithms have been obtained. Barak, Chou, Lei, Schramm, and Sheng search for appearance of particular polylogarithmic-sized subgraphs in both graphs [23]. Mossel and Xu use seeds more efficiently than previous algorithms, creating a signature vertex based on the set of seeds in a large neighborhood of the vertex. The number of seed pairs required is small enough that they can be guessed, yielding an algorithm that does not require side information [24].

We intentionally use the terminology “planted alignment” in analogy with “planted clique”, “planted dense subgraph”, “planted coloring”, and “planted partition”. For these settings, there are several basic problems. One is to find the statistical or information theoretic limits of exact recovery, i.e. the conditions under which an algorithm with unlimited resources can with high probability recover the hidden structure with no errors. Another is to find the information theoretic limits of detection, i.e. the conditions under which an object with a planted structure can be distinguished from an object without one. Finally, there are the conditions under which efficient algorithm can succeed at these tasks. There is a large body of work using spectral algorithms, message passing algorithms, and semidefinite optimization to efficiently recover planted structures. See the surveys of Moore [25], Abbe [26], and Wu and Xu [27] for an overview.

In the case of recovering a planted alignment, finding the information-theoretic limits of exact recovery, often the easiest of the standard problems to resolve, is already challenging. In this paper, we investigate the information-theoretic limits of a problem in between exact recovery and detection: recovery of almost all of a planted alignment with one-sided error.

1 Model

1.1 Notation

A binary relation $\mu \subseteq \mathcal{U}_a \times \mathcal{U}_b$ is a matching if each $i \in \mathcal{U}_a$ and $j \in \mathcal{U}_b$ appears in at most one pair in μ . Define the functions $\alpha : 2^{\mathcal{U}_a \times \mathcal{U}_b} \rightarrow 2^{\mathcal{U}_a}$ and $\beta : 2^{\mathcal{U}_a \times \mathcal{U}_b} \rightarrow 2^{\mathcal{U}_b}$ that find the left and right support of a binary relation:

$$\begin{aligned}\alpha(\mu) &= \{i \in \mathcal{U}_a : \exists j \in \mathcal{U}_b . (i, j) \in \mu\} \\ \beta(\mu) &= \{j \in \mathcal{U}_b : \exists i \in \mathcal{U}_a . (i, j) \in \mu\}.\end{aligned}$$

A matching $\mu \subseteq \mathcal{U}_a \times \mathcal{U}_b$ is a bijection if $\alpha(\mu) = \mathcal{U}_a$ and $\beta(\mu) = \mathcal{U}_b$.

Let \wedge be the minimum or “and” binary operator on $\{0, 1\}$. Let $[n]$ denote the set $\{0, \dots, n-1\}$. For a set \mathcal{U} , let $\binom{\mathcal{U}}{2}$ be the set of unordered pairs of elements of \mathcal{U} . Represent a labeled graph on a vertex set \mathcal{U} by its edge indicator function $G : \binom{\mathcal{U}}{2} \rightarrow [2]$. For a graph G , let $V(G)$ and $E(G)$ be the vertex and edge sets respectively. Throughout, we use boldface letters for random objects and lightface letters for fixed objects.

1.2 Correlated graphs

The correlated Erdős-Rényi graph model has been used in much of the work on alignment recovery for random graphs, beginning with Pedarsani and Grossglauser [2]. The idea is simple: we have two graphs G_a and G_b whose marginal distributions are Erdős-Rényi. Under the true vertex matching, each edge random variable in G_a is aligned with some edge random variable in G_b . These aligned pairs of edge random variables have some joint distribution and this is the only source of correlation between the graphs.

To formalize this, we need the following definition.

Definition 1 (Lifted matching). *A matching $\mu \subseteq \mathcal{U}_a \times \mathcal{U}_b$ gives rise to a lifted matching $\ell(\mu) \subseteq \binom{\mathcal{U}_a}{2} \times \binom{\mathcal{U}_b}{2}$,*

$$\ell(\mu) = \left\{ (\alpha(w), \beta(w)) : w \in \binom{\mathcal{U}}{2} \right\} = \left\{ (\{u_a, v_a\}, \{u_b, v_b\}) : \{(u_a, u_b), (v_a, v_b)\} \in \binom{\mathcal{U}}{2} \right\}.$$

Definition 2. *The distribution of random variables $(\mathbf{X}_a, \mathbf{X}_b) \in \{0, 1\}^2$ is fully specified by a matrix of parameters $p \in \mathbb{R}^{\{0,1\} \times \{0,1\}}$, where $P_{\mathbf{X}_a, \mathbf{X}_b}(i, j) = p_{ij}$. In this case, we say that X_a and X_b have a correlated Bernoulli distribution with parameter p .*

For a matching $\mu \subseteq \mathcal{U}_a \times \mathcal{U}_b$, we define the correlated Erdős-Rényi distribution on pairs of graphs $\mathbf{G}_a : \binom{\mathcal{U}_a}{2} \rightarrow \{0, 1\}$ and $\mathbf{G}_b : \binom{\mathcal{U}_b}{2} \rightarrow \{0, 1\}$, denoted $\text{ER}(\mu, p)$, as follows. For each $(w_a, w_b) \in \ell(\mu)$, $(\mathbf{G}_a(w_a), \mathbf{G}_b(w_b))$ have a correlated Bernoulli distribution with parameter p and these random variables are mutually independent. That is,

$$P_{\mathbf{G}_a, \mathbf{G}_b | \mu}(G_a, G_b | \mu) = \prod_{(w_a, w_b) \in \ell(\mu)} P_{\mathbf{X}_a, \mathbf{X}_b}(G_a(w_a), G_b(w_b)).$$

Because $l(\mu)$ is a matching, each $w_a \in \binom{\mathcal{U}_a}{2}$ appears in exactly one pair $(w_a, w_b) \in l(\mu)$. For a pair $(w_a, w_b) \notin \ell(\mu)$, $G_a(w_a)$ is independent of $G_b(w_b)$.

If $p_{11}p_{00} > p_{10}p_{01}$, then these distributions have *positive correlation*. We will only consider positively correlated graphs in this paper.

1.3 Estimating a planted alignment

We consider the following estimation problem. Let $|\mathcal{U}_a| = |\mathcal{U}_b| = n$ and let μ be a uniformly random bijection between \mathcal{U}_a and \mathcal{U}_b . Let $(\mathbf{G}_a, \mathbf{G}_b) \sim \text{ER}(\mu, p)$.

The most stringent recovery requirement, $\hat{\mu} = \mu$ or exact recovery, was addressed by Cullina and Kiyavash [3]. Their precise results are discussed in Section 1.5. In that case, there is a clear definition of the optimal estimator: the maximum a posteriori (MAP) estimator: $\hat{\mu}(G_a, G_b) = \arg\max_{\hat{\mu}} \Pr[\mu = \hat{\mu} | \mathbf{G}_a = G_a, \mathbf{G}_b = G_b]$. Because μ has a uniform distribution, by Bayes theorem $\hat{\mu}(G_a, G_b) = \arg\max_{\hat{\mu}} \Pr[\mathbf{G}_a = G_a, \mathbf{G}_b = G_b | \mu = \hat{\mu}]$.

This estimator is closely related to the *aligned intersection* of a pair of graphs. Let G_a and G_b be graphs and let μ be a matching between their vertex sets. Then μ provides an alignment

of the subgraphs $G_a[\alpha(\mu)]$ and $G_b[\beta(\mu)]$. Using this alignment, we can compute the intersection of these two subgraphs. The natural vertex set for this intersection graph is μ . We formalize this construction as follows.

Definition 3. Let G_a and G_b be graphs and let $\mu \subseteq V(G_a) \times V(G_b)$ be a matching. Define $G_a \wedge_\mu G_b$, the aligned intersection of G_a and G_b , to be the graph on the vertex set μ such that

$$(G_a \wedge_\mu G_b) : \binom{\mu}{2} \rightarrow \{0, 1\}$$

$$(G_a \wedge_\mu G_b)(\{(u_a, u_b), (v_a, v_b)\}) = G_a(\{u_a, v_a\}) \wedge G_b(\{u_b, v_b\})$$

or equivalently

$$(G_a \wedge_\mu G_b)(w) = G_a(\alpha(w)) \wedge G_b(\beta(w)).$$

Cullina and Kiyavash [4] observed that for a bijection μ ,

$$\Pr[\mathbf{G}_a = G_a, \mathbf{G}_b = G_b | \boldsymbol{\mu} = \mu] \propto \left(\frac{p_{11}p_{00}}{p_{10}p_{01}} \right)^{|E(G_a \wedge_\mu G_b)|}$$

where the constant of proportionality depends on G_a and G_b but not on μ . Thus, in the case of positive correlation, the MAP estimator is $\hat{\boldsymbol{\mu}}(G_a, G_b) = \operatorname{argmax}_{\hat{\boldsymbol{\mu}}} |E(G_a \wedge_{\hat{\boldsymbol{\mu}}} G_b)|$.

Herein we consider partial recovery of $\boldsymbol{\mu}$ using $(\mathbf{G}_a, \mathbf{G}_b)$, which is interesting when exact recovery is impossible. In particular, we would like to match some of the vertices of \mathbf{G}_a to the corresponding vertices in \mathbf{G}_b without any errors. This means that we want an estimator $\hat{\boldsymbol{\mu}}$ such that $\hat{\boldsymbol{\mu}} \subseteq \boldsymbol{\mu}$ and $|\hat{\boldsymbol{\mu}}|$ is as large as possible. We are interested in estimators that satisfy these conditions with probability $1 - o(1)$.

For a partial matching μ' ,

$$\Pr[\boldsymbol{\mu} \supseteq \mu' | \mathbf{G}_a = G_a, \mathbf{G}_b = G_b] = \sum_{\mu \supseteq \mu', |\mu|=n} \Pr[\boldsymbol{\mu} = \mu | \mathbf{G}_a = G_a, \mathbf{G}_b = G_b].$$

There are two natural generalization of the MAP estimator for partial recovery. The first fixes n' , the size of the estimated matching, and selects $\hat{\boldsymbol{\mu}}$ that maximizes $\Pr[\boldsymbol{\mu} \supseteq \hat{\boldsymbol{\mu}} | \mathbf{G}_a = G_a, \mathbf{G}_b = G_b]$. The second fixes ϵ , an error budget, and selects a $\hat{\boldsymbol{\mu}}$ satisfying $\Pr[\boldsymbol{\mu} \supseteq \hat{\boldsymbol{\mu}} | \mathbf{G}_a = G_a, \mathbf{G}_b = G_b] \geq 1 - \epsilon$ and maximizing $|\hat{\boldsymbol{\mu}}|$. Neither of these are particularly straightforward to analyze, so we introduce the k -core alignment estimator.

1.4 k -cores and k -core alignments

Let $\delta(G)$ be the minimum degree of G and for $S \subseteq V(G)$, let $G[S]$ be the subgraph of G induced by S . We adopt the convention that for the null graph, i.e. G such that $V(G) = \emptyset$, $\delta(G) = \infty$. Thus there is always some $S \subseteq V(G)$ such that $\delta(G[S]) \geq k$. If $\delta(G[S]) \geq k$ and $\delta(G[S']) \geq k$, then $\delta(G[S \cup S']) \geq k$. Thus there is a unique maximum among the sets that induce subgraphs with minimum degree at least k . The subgraph induced by this set is the k -core of G [28].¹

We introduce the following related notion.

Definition 4. A k -core alignment of G_a and G_b is a matching $\mu \subseteq V(G_a) \times V(G_b)$ such that $\delta(G_a \wedge_\mu G_b) \geq k$ and for all matchings $\mu' \supset \mu$, $\delta(G_a \wedge_{\mu'} G_b) < k$.

Figure 1 illustrates the concepts of aligned intersection and k -core alignment.

Definition 5. The k -core alignment estimator $\hat{\boldsymbol{\mu}}_k$ selects a k -core alignment of $(\mathbf{G}_a, \mathbf{G}_b)$. If there is more than one k -core alignment, it makes an arbitrary choice.

¹When every nonempty induced subgraph of a graph G has a minimum degree less than k , some authors say that the k -core does not exist. In this case the k -core of G is the null graph under our convention.

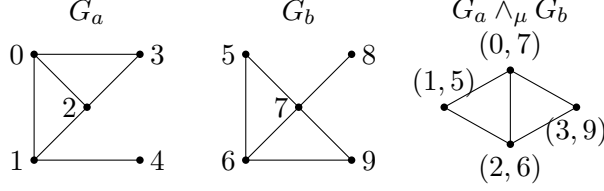


Figure 1: The matching $\mu = \{(0, 7), (1, 5), (2, 6), (3, 9)\}$ is a 2-core alignment of G_a and G_b : $\delta(G_a \wedge_{\mu} G_b) = 2$ and extending μ with $(4, 8)$ leads to a minimum degree of 0.

1.5 Results

We have the following results about the performance of the k -core alignment estimator.

Theorem 1 (Achievability). *Let p satisfy the conditions*

$$p_{11} \geq \omega\left(\frac{1}{n}\right) \quad (1)$$

$$p_{11} \leq \frac{1}{8e^3} \quad (2)$$

$$\frac{p_{01}p_{10}}{p_{00}p_{11}} + p_{01} + p_{10} \leq n^{-\Omega(1)}. \quad (3)$$

Then there is a choice of k such that with probability $1 - o(1)$, the k -core alignment estimator $\hat{\boldsymbol{\mu}}_k$ satisfies $\hat{\boldsymbol{\mu}}_k \subseteq \boldsymbol{\mu}$ and $|\hat{\boldsymbol{\mu}}_k| \geq n(1 - o(1))$. That is, the estimator includes no incorrect pairs and almost all correct pairs.

Section 2 contains the proof. The main condition of Theorem 1, (1), requires $\mathbf{G}_a \wedge_{\boldsymbol{\mu}} \mathbf{G}_b$ to have an average degree that grows with n . Condition (2) is very mild sparsity constraint on $\mathbf{G}_a \wedge_{\boldsymbol{\mu}} \mathbf{G}_b$. Condition (3) requires \mathbf{G}_a and \mathbf{G}_b to have sufficient positive correlation and to be mildly sparse.

Theorem 2 (Converse). *Let*

$$p_{11} \leq \mathcal{O}\left(\frac{1}{n}\right) \quad (4)$$

$$\frac{p_{01}p_{10}}{p_{11}p_{00}} < 1. \quad (5)$$

Then for any estimator $\hat{\boldsymbol{\mu}}$ of $\boldsymbol{\mu}$ given $(\mathbf{G}_a, \mathbf{G}_b)$ and any integer sequence $\epsilon(n) \leq o(n)$, the probability that $\hat{\boldsymbol{\mu}} \subseteq \boldsymbol{\mu}$ and $|\hat{\boldsymbol{\mu}}| \geq n - \epsilon(n)$ is $o(1)$.

Section 3 contains the proof. Observe that the condition (4) matches (1). Condition (5) is weaker than than (3): the converse is valid for any amount of positive correlation.

For the exact recovery problem, Cullina and Kiyavash [3] showed that if p satisfies the conditions

$$p_{11} \geq \frac{\log n + \omega(1)}{n}$$

$$p_{11} + p_{01} + p_{10} \leq \mathcal{O}\left(\frac{1}{\log n}\right)$$

$$\frac{p_{01}p_{10}}{p_{11}p_{00}} \leq \mathcal{O}\left(\frac{1}{(\log n)^3}\right),$$

then the maximum a posteriori estimator for μ given $(\mathbf{G}_a, \mathbf{G}_b)$ is correct with probability $1 - o(1)$. Additionally, if \mathbf{p} satisfies

$$p_{11} \leq \frac{\log n - \omega(1)}{n} \quad \text{and} \quad \frac{p_{01}p_{10}}{p_{11}p_{00}} < 1,$$

then any estimator for μ given $(\mathbf{G}_a, \mathbf{G}_b)$ is correct with probability $o(1)$. In other words, exact recovery of μ requires logarithmic average degree in the intersection graph while recovery of almost all of μ requires only a growing average degree.

1.6 Product graphs

The aligned intersection of G_a and G_b has another interpretation. Let $G_a \times G_b$ be the tensor product of G_a and G_b . This is the graph with $V(G_a \times G_b) = V(G_a) \times V(G_b)$ and

$$(G_a \times G_b)(\{(u_a, u_b), (v_a, v_b)\}) = G_a(\{u_a, v_a\}) \wedge G_b(\{u_b, v_b\}).$$

In other words, the adjacency matrix of $G_a \times G_b$ is the tensor product of the adjacency matrices of G_a and G_b . Then $G_a \wedge_{\mu} G_b = (G_a \times G_b)[\mu]$.

From this point of view, exact recovery of μ corresponds to finding a dense n -vertex subgraph inside the n^2 -vertex graph $G_a \times G_b$. This looks superficially like recovering a single dense community is a stochastic block model, a problem which has been extensively studied. There are two important differences. First, we only need to search over subgraphs induced by matchings, not all n -vertex subgraphs. This does not significantly reduce the total number of candidate subgraphs, but it has a bigger effect on the number of subgraphs that are nearly equal to the true matching. Second, the edge random variables in $G_a \times G_b$ are not jointly independent. Because of this, bounding the probability of each error event requires some care.

The fact that we are searching for a subgraph of size $\sqrt{|V(G_a \times G_b)|}$ has potential implications for the computational tractability of this problem due to the planted clique hypothesis [27] and associated conditional hardness results for planted dense subgraph problems. However, in this paper, we focus only on information-theoretic thresholds.

Note that the k -core alignment of G_a and G_b is not the k -core of $G_a \times G_b$, which in general will be much larger and not induced by a matching.

1.7 Proof outline

Our achievability proof has the following structure. First, we need to show that the true alignment of \mathbf{G}_a and \mathbf{G}_b yields a large k -core alignment. This follows easily from known results about the k -core of an Erdős-Rényi graph. The majority of our work is devoted to the second task: showing that this is the only k -core alignment. For each matching $\mu \not\subseteq \mu^*$, we need to bound the probability that it is a k -core alignment of G_a and G_b . A large number of matchings can be immediately ruled by the maximality part of the definition of k -core alignment. We define an error event for each of the remaining matching and use a union bound over them. There are exponentially many of these error events, so it is crucial that our bound on error probability depends on the distance between the imposter matching and the true matching. This part of the argument is in Section 2.1 and is summarized in Lemma 1.

In order for μ to be a k -core alignment, each vertex in $\mathbf{G}_a \wedge_{\mu} \mathbf{G}_b$ must have degree at least k . This is much more difficult for vertices in μ that do not appear in the true matching μ^* . To bound the probability that μ is a k -core alignment, we consider the sum of the degrees of the vertices in $\mu \setminus \mu^*$.

The main technical task is to obtain large deviations upper bounds for these sum-of-degrees random variables. Each of these random variables is the sum of many correlated indicator random variables. To obtain tight bounds, we take advantage of the structure of the correlation by analyzing the cycle-path decomposition of the imposter matching relative to the true matching. This decomposition is explained in Section 2.2. In Lemma 2, whose proof is in Appendix A, we bound the generating function for the sum-of-degrees random variable using combinatorial arguments. This allows us to find conditions under which the tails of their distributions behave like the tails of Poisson random variables.

Our converse proof is based on the concept of list estimation or list decoding. It has two main components. First, we find a relationship between the number of automorphisms of $\mathbf{G}_a \wedge_{\mu} \mathbf{G}_b$ and the list length for any list estimator that succeeds with high probability. Second, we use the fact that sufficiently sparse Erdős-Rényi graphs have many isolated vertices to obtain a lower bound on the number of automorphisms of $\mathbf{G}_a \wedge_{\mu} \mathbf{G}_b$.

2 Achievability

2.1 Weak k -core alignments

The property discussed at the start of Section 1.4 that leads to the uniqueness of the k -core has the following analogue for k -core alignments. If $\delta(G_a \wedge_{\mu} G_b) \geq k$, $\delta(G_a \wedge_{\mu'} G_b) \geq k$, and $\mu \cup \mu'$ is a matching, then $\delta(G_a \wedge_{\mu \cup \mu'} G_b) \geq k$. The graphs G_{μ} and $G_{\mu'}$ are induced subgraphs of $G_{(\mu \cup \mu')}$. Thus each vertex in $G_{(\mu \cup \mu')}$ has a degree that is at least as large as its degree in the subgraph.

Because $\mu \cup \mu'$ is not guaranteed to be a matching, there may be more than one k -core alignment of G_a and G_b . For example, if both G_a and G_b are complete graphs with n vertices, every bijection between $V(G_a)$ and $V(G_b)$ is an n -core alignment.

The k -core alignment estimator can make an error when some matching other than μ is a k -core alignment. To analyze this event, we introduce weak k -core alignments.

Definition 6. Let $\deg_G : V(G) \rightarrow \mathbb{N}$ be the degree function for the graph G . Let μ and μ^* be matchings of $V(G_a)$ and $V(G_b)$ and let

$$M(\mu, \mu^*; G_a, G_b) = \sum_{v \in \mu \setminus \mu^*} \deg_{G_a \wedge_{\mu} G_b}(v).$$

A matching μ is a weak k -core alignment of G_a and G_b relative to μ^* if $M(\mu, \mu^*; G_a, G_b) \geq k|\mu \setminus \mu^*|$.

Let $\mathbf{M}_{\mu, \mu^*} = M(\mu, \mu^*; \mathbf{G}_a, \mathbf{G}_b)$.

Observe that if μ is a k -core alignment of G_a and G_b , then it is also a weak k -core alignment of G_a and G_b relative to any matching μ^* . We have relaxed the property in two ways: first by only checking the vertices that are matched differently in μ than in μ^* and second by checking the average degree of these vertices rather than the minimum.

All $\mu \subseteq \mu^*$ are trivially weak k -core alignments relative to μ^* (the sum is empty). We will show that under certain conditions, every weak k -core alignment of \mathbf{G}_a and \mathbf{G}_b relative to μ is a subset of μ .

Definition 7. A matching μ is μ^* -maximal if no pairs from μ^* can be added to μ without destroying the matching property. More precisely, for all $(i, j) \in \mu^*$, either $i \in \alpha(\mu)$ or $j \in \beta(\mu)$. Let

$$\mathcal{M}(\mu^*, d) = \{\mu : \mu \text{ is } \mu^* \text{ maximal and } |\mu \setminus \mu^*| = d\}.$$

This property allows us to show that a large number of matchings cannot be k -core alignments of $(\mathbf{G}_a, \mathbf{G}_b)$ because they are not μ -maximal. If it is possible to add any pairs from the true matching μ to μ , then $\hat{\mu}_k \neq \mu$.

The following lemma allows us to bound the probability of error of the k -core alignment estimator.

Lemma 1. *Let μ^* be a bijection. Then*

$$\Pr \left[\bigvee_{\mu} (\delta(\mathbf{G}_a \wedge_{\mu} \mathbf{G}_b) \geq k) \wedge (\mu \not\subseteq \mu^*) \right] \leq \exp(n^2 \xi) - 1$$

where the disjunction is over all $\mu \subseteq \mathcal{U}_a \times \mathcal{U}_b$ and

$$\log \xi = \max_{d \geq 1} \max_{\mu \in \mathcal{M}(\mu, d)} \frac{1}{d} \log \Pr[\mathbf{M}_{\mu, \mu^*} \geq kd] \quad (6)$$

Proof. If $|\mu \setminus \mu^*| = d$ and $\delta(\mathbf{G}_a \wedge_{\mu} \mathbf{G}_b) \geq k$, directly from the definition of \mathbf{M}_{μ, μ^*} we have $\mathbf{M}_{\mu, \mu^*} \geq kd$. Any matching μ has a unique extension to a μ^* -maximal matching that is produced by adding as many elements of μ^* as possible: $\mu' = \mu \cup (\mu^* \setminus (\alpha(\mu) \times \beta(\mu)))$. Then $\mu \setminus \mu^* = \mu' \setminus \mu^*$, $\mu' \in \mathcal{M}(\mu^*, d)$, and $\mathbf{M}_{\mu', \mu^*} \geq kd$. Thus the event in the statement of the lemma is equivalent to

$$\begin{aligned} \Pr \left[\bigvee_{d \geq 1} \bigvee_{\mu' \in \mathcal{M}(\mu^*, d)} \mathbf{M}_{\mu', \mu^*} \geq kd \right] &\leq \sum_{d \geq 1} \sum_{\mu' \in \mathcal{M}(\mu^*, d)} \Pr[\mathbf{M}_{\mu', \mu^*} \geq kd] \\ &\stackrel{(a)}{\leq} \sum_{d \geq 1} \frac{n^{2d}}{d!} \xi^d \\ &\leq \exp(n^2 \xi) - 1 \end{aligned}$$

where (a) uses the bound $|\mathcal{M}(\mu^*, d)| \leq \frac{n^{2d}}{d!}$ and the definition of ξ in (6). Because μ^* is a bijection, each $\mu \in \mathcal{M}(\mu^*, d)$ is fully specified by $\mu \setminus \mu^*$. There are $\binom{n}{d}$ choices of $\alpha(\mu \setminus \mu^*)$, $\binom{n}{d}$ choices of $\beta(\mu \setminus \mu^*)$, and $d!$ bijections between these sets, and $\binom{n}{d} \leq \frac{n^d}{d!}$. \square

2.2 Lifted matchings

Given two matchings $\mu, \mu' \subseteq \mathcal{U}_a \times \mathcal{U}_b$, let $(\mu + \mu') : \mathcal{U}_a \times \mathcal{U}_b \rightarrow \mathbb{N}$ be the multisubset of $\mathcal{U}_a \times \mathcal{U}_b$ in which the multiplicity of each element is the sum of its multiplicity in μ and its multiplicity in μ' . In the bipartite multigraph $(\mathcal{U}_a, \mathcal{U}_b, \mu + \mu')$, all vertices have degree 0, 1, or 2, so the multigraph is the union of paths and even-length cycles (including cycles of length two, which are pairs of parallel edges). More precisely, the degree of $u_a \in \mathcal{U}_a$ is $\mathbb{1}(u_a \in \alpha(\mu)) + \mathbb{1}(u_a \in \alpha(\mu'))$ and the degree of $u_b \in \mathcal{U}_b$ is $\mathbb{1}(u_b \in \beta(\mu)) + \mathbb{1}(u_b \in \beta(\mu'))$.

A vertex pair $w \in \binom{\mu}{2}$ contains 0, 1 or 2 elements from $\mu \setminus \mu^*$. The others are from $\mu^* \cap \mu$.

This matters because

$$\begin{aligned} \mathbf{M}_{\mu, \mu^*} &= \sum_{v \in \mu \setminus \mu^*} \sum_{w \in \binom{\mu}{2}} \mathbb{1}(v \in w) \cdot (\mathbf{G}_a \wedge_{\mu} \mathbf{G}_b)(w) \\ &= \sum_{w \in \binom{\mu}{2}} |w \cap (\mu \setminus \mu^*)| (\mathbf{G}_a \wedge_{\mu} \mathbf{G}_b)(w) \end{aligned} \quad (7)$$

$$\mu^* = \{(0, 5), (1, 6), (2, 7), (3, 8), (4, 9)\}$$

$$\mu = \{(0, 6), (2, 8), (3, 7), (4, 9)\}$$

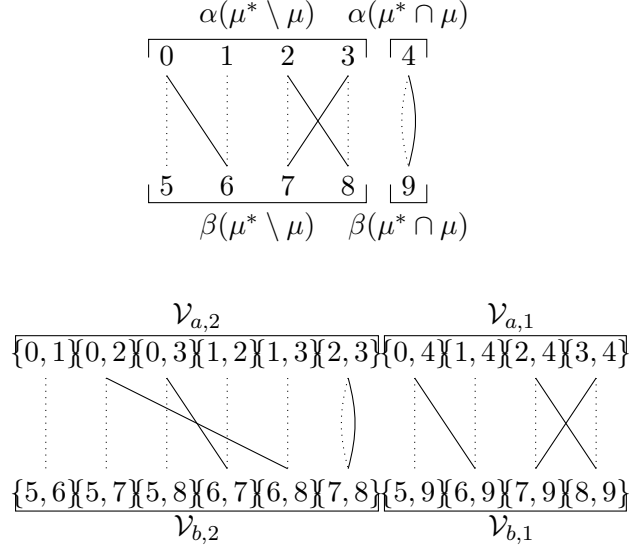


Figure 2: Illustration of the decomposition of $\mu + \mu^*$ and $\ell(\mu) + \ell(\mu^*)$ into cycles and paths. The matchings μ and $\ell(\mu)$ are drawn with solid lines and the bijections μ^* and $\ell(\mu^*)$ are drawn with dotted lines. The sets $\mathcal{V}_{a,0}$ and $\mathcal{V}_{b,0}$ are empty.

Because μ^* is a bijection and $\alpha(\mu^*) = \mathcal{U}_a$, $\alpha(\mu \cap \mu^*)$ and $\alpha(\mu^* \setminus \mu)$ partition the vertex set \mathcal{U}_a . Similarly $\beta(\mu \cap \mu^*)$ and $\beta(\mu^* \setminus \mu)$ partition the vertex set \mathcal{U}_b . At the level of vertex pairs, we have partitions of $\binom{\mathcal{U}_a}{2}$ and $\binom{\mathcal{U}_b}{2}$ into three regions each:

$$\mathcal{V}_{a,i} = \left\{ w_a \in \binom{\mathcal{U}_a}{2} : |w_a \cap \alpha(\mu \cap \mu^*)| = 2 - i \right\}$$

$$\mathcal{V}_{b,i} = \left\{ w_b \in \binom{\mathcal{U}_b}{2} : |w_b \cap \beta(\mu \cap \mu^*)| = 2 - i \right\}.$$

for $i \in \{0, 1, 2\}$. The important fact for us is that both $\ell(\mu)$ and $\ell(\mu^*)$ match elements of $\mathcal{V}_{a,i}$ with elements of $\mathcal{V}_{b,i}$.

The matchings $\ell(\mu)$ and $\ell(\mu^*)$ also decompose into a union of paths and cycles. Because μ^* is a bijection, the length of each path is odd and the edges from μ are never the initial or final edges in a path. Each of these paths and cycles stays within one of the three regions.

Definition 8. For a matching μ and a bijection μ^* , define the following. For $i \in \{0, 1, 2\}$, let $\nu_i = \ell(\mu) \cap (\mathcal{V}_{a,i} \times \mathcal{V}_{b,i})$ and $\nu_i^* = \ell(\mu^*) \cap (\mathcal{V}_{a,i} \times \mathcal{V}_{b,i})$. For $\ell \geq 1$, let $t_{i,\ell}^\circ$ be the number of cycles of length 2ℓ and let $t_{i,\ell}$ be the number of paths of length $2\ell + 1$ in the decomposition of $\nu_i + \nu_i^*$.

We have $\ell(\mu) = \nu_0 \cup \nu_1 \cup \nu_2$ and $\ell(\mu^*) = \nu_0^* \cup \nu_1^* \cup \nu_2^*$, so t and t° capture the whole structure of $\ell(\mu) + \ell(\mu^*)$.

An example of this decomposition is illustrated in Figure 2. The matchings $\ell(\mu)$ and $\ell(\mu^*)$ always have the same structure between $\mathcal{V}_{a,0}$ and $\mathcal{V}_{b,0}$. Thus $t_{0,\ell}^\circ = 0$ for $\ell \geq 2$ and $t_{0,\ell} = 0$ for all ℓ . The structure between $\mathcal{V}_{a,1}$ and $\mathcal{V}_{b,1}$ is $|\mu^* \cap \mu|$ copies of the structure between $\alpha(\mu^* \setminus \mu)$ and $\beta(\mu^* \setminus \mu)$ and consequently contains no cycles of length two, i.e. $t_{1,1}^\circ = 0$. Observe that between

$\mathcal{V}_{a,2}$ and $\mathcal{V}_{b,2}$, a cycle of length two can only be produced by a cycle of length four in the region between $\alpha(\mu^* \setminus \mu)$ and $\beta(\mu^* \setminus \mu)$. Thus $t_{2,1}^{\circ} \leq |\mu \setminus \mu^*|/2$.

2.3 Generating functions

Definition 9. Let $A_{\mu,\mu^*}(z)$ be the generating function for the random variable \mathbf{M}_{μ,μ^*} :

$$A_{\mu,\mu^*}(z) = \mathbb{E}[z^{\mathbf{M}_{\mu,\mu^*}}].$$

Let $p_{1*} = p_{10} + p_{11}$ and $p_{*1} = p_{01} + p_{11}$.

Lemma 2. For a matching μ and a bijection μ^* , if $\frac{p_{11}p_{00}}{p_{10}p_{01}} \geq 1$ then

$$\log A_{\mu,\mu^*}(z) \leq t_{2,1}^{\circ} p_{11}(z^2 - 1) + \frac{\tilde{t}}{4}(2p_{1*}p_{*1}(z^2 - 1) + p_{11}^2(z^2 - 1)^2)$$

where $\tilde{t} = d(n-1) - 2t_{2,1}^{\circ}$ and $d = |\mu \setminus \mu^*|$.

The proof is given in Appendix A.

Lemma 3. Let $q_1 \geq 0$ and $q_2 \geq 0$. Then

$$\begin{aligned} \operatorname{argmin}_{z \geq 0} \exp(q_2(z^2 - 1) + q_1(z - 1))z^{-\tau} &\leq \zeta^{\tau} \\ \zeta &= \max\left(\sqrt{2}e\frac{q_1}{\tau}, 4e\left(\frac{q_2}{\tau}\right)^{1/2}\right). \end{aligned} \tag{8}$$

The proof is given in Appendix B.

Lemma 4. If p satisfies conditions (1), (2), and (3) and $k \geq \Omega(np_{11})$, then

$$\max_{d \geq 1} \max_{\mu \in \mathcal{M}(\mu^*, d)} \frac{1}{d} \log \Pr[\mathbf{M}_{\mu,\mu^*} \geq kd] \leq -\omega(\log n) \tag{9}$$

The proof is given in Appendix C.

Theorem 3 ([29] Theorem 2). Let $c = c(n) = (n-1)p(n)$. For every $\epsilon > 0$ there is a constant d , such that for all $c(n) > d$ and $k(n) \leq c - c^{\frac{1}{2} + \epsilon}$, has the size of the k -core of a graph $G \sim G(n, p)$ is at least $n - n \exp(-c^{\epsilon})$ with probability $1 - o(1)$.

Proof of Theorem 1. First we will show that $\mathbf{G}_a \wedge_{\mu} \mathbf{G}_b$ has a large k -core, so there is some $\hat{\mu} \subseteq \mu$ such that $\hat{\mu}$ is a k -core alignment and $|\hat{\mu}| \geq n(1 - o(1))$. We have $(\mathbf{G}_a \wedge_{\mu} \mathbf{G}_b) \sim G(n, p_{11})$ and $np_{11} \geq \omega(1)$. From the application of Theorem 3 with $\epsilon = \frac{1}{4}$, for $k = np_{11}(1 - (np_{11})^{-\frac{1}{4}}) \geq np_{11}(1 - o(1))$, G_{μ} has a k -core of size

$$n(1 - \exp(-(np_{11})^{\frac{1}{4}})) \geq n(1 - e^{-\omega(1)}) \geq n(1 - o(1))$$

with probability $1 - o(1)$.

Now we will show that $\hat{\mu} \subseteq \mu^*$, i.e. every vertex pair in $\hat{\mu}$ is correct. From Lemma 1, the probability of error is at most $\exp(n^2\xi) - 1$ and from Lemma 4 we have $\xi \leq n^{-\omega(1)}$, so the probability that $\hat{\mu} \not\subseteq \mu^*$ is $o(1)$. \square

3 Converse

Recall from Section 1.3 that when we are trying to estimate a subset of μ , the quality of a partial matching μ' depends on the list of bijections that extend it.

Definition 10. Let \mathbf{X} and \mathbf{Y} be random variables on \mathcal{X} and \mathcal{Y} respectively. A list estimator for \mathbf{Y} given \mathbf{X} is a function $S : \mathcal{X} \rightarrow 2^{\mathcal{Y}}$. The estimator succeeds when $\mathbf{Y} \in S(\mathbf{X})$.

Lemma 5. Let \mathbf{Y} be a random variable on a finite set \mathcal{Y} with distribution P_Y . Let $S \subseteq \mathcal{Y}$ such that $|S| = \ell$. Then

$$\Pr[\mathbf{Y} \in S] \leq \mathbb{E} \left[\min \left(1, \frac{\ell}{|\{y' \in \mathcal{Y} : P_Y(y') \geq P_Y(\mathbf{Y})\}|} \right) \right] \quad (10)$$

Proof. Let (p_0, p_1, \dots) be the list of distinct probabilities that appear in P_Y , sorted from largest to smallest. Let $\mathcal{Y}_i = \{y \in \mathcal{Y} : P_Y(y) = p_i\}$. Let S^* be a set of size ℓ that maximizes $\Pr[\mathbf{Y} \in S]$.

If $\sum_{i=0}^j |\mathcal{Y}_i| \leq \ell$, then $\mathcal{Y}_j \subseteq S^*$. If $\sum_{i=0}^j |\mathcal{Y}_i| > \ell$, then $|\mathcal{Y}_j \cap S^*| = \max(0, \ell - \sum_{i=0}^{j-1} |\mathcal{Y}_i|)$. We have the inequality

$$\left(\ell - \sum_{i=0}^{j-1} |\mathcal{Y}_i| \right) \left(\sum_{i=0}^j |\mathcal{Y}_i| \right) = \ell |\mathcal{Y}_j| + \left(\ell - \sum_{i=0}^j |\mathcal{Y}_i| \right) \left(\sum_{i=0}^{j-1} |\mathcal{Y}_i| \right) \leq \ell |\mathcal{Y}_j|,$$

so the fraction of \mathcal{Y}_j appearing in S^* is

$$\frac{|\mathcal{Y}_j \cap S^*|}{|\mathcal{Y}_j|} = \frac{\max(0, \ell - \sum_{i=0}^{j-1} |\mathcal{Y}_i|)}{|\mathcal{Y}_j|} \leq \frac{\ell}{\sum_{i=0}^j |\mathcal{Y}_i|}.$$

Observe that for $y \in \mathcal{Y}_j$, $\{y' \in \mathcal{Y} : P_Y(y') \geq P_Y(y)\} = \bigcup_{i=0}^j \mathcal{Y}_i$. Thus for all j , the fraction of \mathcal{Y}_j appearing in S^* is at most as large as the contribution of \mathcal{Y}_j to the right side of (10). \square

Let $\mu \subseteq \mathcal{U}_a \times \mathcal{U}_b$ be a matching and let $\pi \subseteq \mu \times \mu$ be a bijection. Then we can extract another matching $\gamma(\pi) \subseteq \mathcal{U}_a \times \mathcal{U}_b$ as follows. Observe that

$$\pi \subseteq (\mu \times \mu) \subseteq (\mathcal{U}_a \times \mathcal{U}_b) \times (\mathcal{U}_a \times \mathcal{U}_b)$$

and define $\gamma(\pi) = \{(u_a, v_b) : ((u_a, u_b), (v_a, v_b)) \in \pi\}$.

Cullina and Kiyavash proved the following fact.

Lemma 6. Let G_a and G_b be graphs, let μ be a matching between their vertex sets, and let $\frac{p_{01}p_{10}}{p_{11}p_{00}} < 1$. For all $\pi \in \text{Aut}(G_a \wedge_{\mu} G_b)$,

$$\Pr[\mu = \gamma(\pi) | (\mathbf{G}_a, \mathbf{G}_b) = (G_a, G_b)] \geq \Pr[\mu = \mu | (\mathbf{G}_a, \mathbf{G}_b) = (G_a, G_b)]$$

Proof. This follows immediately by combining Lemma II.2 and Lemma V.1 of [3]. \square

Theorem 4. Let $\mathbf{S} = S(\mathbf{G}_a, \mathbf{G}_b)$ be a list estimator for μ . Then

$$\Pr[\mu \in \mathbf{S} \wedge |\mathbf{S}| \leq \ell] \leq \mathbb{E} \left[\min(1, \ell |\text{Aut}(\mathbf{G}_a \wedge_{\mu} \mathbf{G}_b)|^{-1}) \right]$$

Proof. Suppose that $|S(G_a, G_b)| > \ell$ for some (G_a, G_b) . Then by instead using a shorter list in these cases, we can create some S' such that $|S'(G_a, G_b)| \leq \ell$ for all (G_a, G_b) and this can only increase $\Pr[\mu \in \mathbf{S} \wedge |\mathbf{S}| \leq \ell]$.

Applying Lemmas 6 and 5 for all (G_a, G_b) and then averaging over $(\mathbf{G}_a, \mathbf{G}_b)$ gives the claim. \square

Lemma 7. *If $\mathbf{G} \sim \text{ER}(n, p)$ and $p \leq \mathcal{O}(\frac{1}{n})$, then \mathbf{G} has $\Omega(n)$ isolated vertices with probability $1 - o(1)$.*

Proof. The expected number of isolated vertices is $n(1-p)^{n-1} \geq ne^{\Omega(1)}$. By a standard use of the second moment method, with probability $1 - o(1)$, the number of isolated vertices is within a factor of $1 - o(1)$ of the mean. \square

Proof of Theorem 2. The estimated partial matching $\hat{\mu}$ can be interpreted as a list estimator. There are $(n - |\hat{\mu}|)!$ bijections μ that extend $\hat{\mu}$, i.e. $|\mu| = n$ and $\hat{\mu} \subseteq \mu$. Then

$$\begin{aligned} \Pr[\boldsymbol{\mu} \in \mathbf{S} \wedge |\mathbf{S}| \leq \epsilon(n)!] &\stackrel{(a)}{\leq} \mathbb{E} [\min(1, \epsilon(n)! |\text{Aut}(\mathbf{G}_a \wedge_{\boldsymbol{\mu}} \mathbf{G}_b)|^{-1})] \\ &\leq \Pr[|\text{Aut}(\mathbf{G}_a \wedge_{\boldsymbol{\mu}} \mathbf{G}_b)| \leq \epsilon(n)!] + \frac{\epsilon(n)!}{|\text{Aut}(\mathbf{G}_a \wedge_{\boldsymbol{\mu}} \mathbf{G}_b)|} \\ &\stackrel{(b)}{\leq} o(1) + \frac{\epsilon(n)!}{\Omega(n)!} \\ &\leq o(1). \end{aligned}$$

where (a) uses Theorem 4, and (b) uses the following argument. If a graph G has j isolated vertices, then $|\text{Aut}(G)| \geq j!$. The true intersection graph $\mathbf{G}_a \wedge_{\boldsymbol{\mu}} \mathbf{G}_b$ has distribution $\text{ER}(n, p_{11})$, so from Lemma 7 and (4), $|\text{Aut}(\mathbf{G}_a \wedge_{\boldsymbol{\mu}} \mathbf{G}_b)| \geq (\Omega(n))!$ with probability $1 - o(1)$. \square

A Proof of Lemma 2

A.1 Generating function combinatorics

Definition 11. *Define the following matrices indexed by $\{0, 1\} \times \{0, 1\}$:*

$$P = \begin{pmatrix} p_{00} & p_{01} \\ p_{10} & p_{11} \end{pmatrix} \quad Z = \begin{pmatrix} 1 & 1 \\ 1 & z \end{pmatrix}.$$

For $\ell \geq 1$, define the generating functions

$$\begin{aligned} a_{\ell}(z) &= \vec{1}^T (PZ)^{\ell} P \vec{1} \\ a_{\ell}^{\circ}(z) &= \text{tr}((PZ)^{\ell}). \end{aligned}$$

Lemma 8.

$$A_{\mu, \mu^*}(z) = \prod_{\ell \geq 1} a_{\ell}(z)^{t_{\ell,1}} a_{\ell}^{\circ}(z)^{t_{\ell,1}^{\circ}} a_{\ell}(z^2)^{t_{\ell,2}} a_{\ell}^{\circ}(z^2)^{t_{\ell,2}^{\circ}}$$

Proof. From (7), we have

$$\mathbf{M}_{\mu, \mu^*} = \sum_{w \in \nu_1} (\mathbf{G}_a \wedge_{\mu} \mathbf{G}_b)(w) + 2 \sum_{w \in \nu_2} (\mathbf{G}_a \wedge_{\mu} \mathbf{G}_b)(w).$$

and these two terms are independent. From this we obtain

$$\begin{aligned} \Pr[\mathbf{G}_a = G_a, \mathbf{G}_b = G_b] z^{M(\mu, \mu^*; G_a, G_b)} &= \\ \prod_{(w_a, w_b) \in \nu_1^*} p_{G_a(w_a), G_b(w_b)} &\prod_{(w_a, w_b) \in \nu_2^*} p_{G_a(w_a), G_b(w_b)} \prod_{(w_a, w_b) \in \nu_1} z^{G_a(w_a) \wedge G_b(w_b)} \prod_{(w_a, w_b) \in \nu_2} z^{2(G_a(w_a) \wedge G_b(w_b))} \end{aligned}$$

which is the contribution of a particular graph pair to the generating function for \mathbf{M}_{μ, μ^*} . Each w_a or w_b appears at most twice in this expression: once in a factor of $p_{G_a(w_a), G_b(w_b)}$ and up to once in a factor of $z^{(G_a(w_a) \wedge G_b(w_b))}$ or $z^{2(G_a(w_a) \wedge G_b(w_b))}$. Thus

$$A_{\mu, \mu^*}(z) = \sum_{G_a, G_b} \Pr[\mathbf{G}_a = G_a, \mathbf{G}_b = G_b] z^{M(\mu, \mu^*; G_a, G_b)}$$

factorizes based on the cycle and path decomposition of $\ell(\mu) + \ell(\mu^*)$. Because the value of $G_a(w_a)$ is used in at most places, the sum over $G_a(w_a) \in [2]$ can be interpreted as a matrix multiplication. The matrix that contributes the factor from $\ell(\mu)$ is Z and the matrix that contributes the factor from $\ell(\mu^*)$ is P . \square

Definition 12. For a sequence $f \in \{0, 1\}^\ell$, let $k_1(f)$ be the number of ones in the sequence, $k_{11}(f)$ be the number of pairs of consecutive ones, and $k_{11}^\circ(f)$ be the number of pairs of consecutive ones counting the first and last position as consecutive. Define the generating functions

$$b_\ell(x, y) = \sum_{f \in \{0, 1\}^\ell} x^{k_1(f)} y^{k_{11}(f)}$$

$$b_\ell^\circ(x, y) = \sum_{f \in \{0, 1\}^\ell} x^{k_1(f)} y^{k_{11}^\circ(f)}$$

Lemma 9. For all $\ell \geq 1$,

$$a_\ell(z) = b_\ell\left(p_{1*} p_{*1}(z-1), \frac{p_{11}}{p_{1*} p_{*1}}\right)$$

$$a_\ell^\circ(z) = b_\ell^\circ\left(p_{1*} p_{*1}(z-1), \frac{p_{11}}{p_{1*} p_{*1}}\right)$$

Proof. Let $Z_0 = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$ and $Z_1 = \begin{pmatrix} 0 & 0 \\ 0 & z-1 \end{pmatrix}$. Then

$$\begin{aligned} a_\ell(z) &= \vec{1}^T (PZ)^\ell P \vec{1} \\ &= \vec{1}^T (P(Z_0 + Z_1))^\ell P \vec{1} \\ &= \sum_{f \in \{0, 1\}^\ell} \vec{1}^T \left(\prod_{i \in [\ell]} PZ_{f(i)} \right) P \vec{1} \\ &\stackrel{(a)}{=} \sum_{f \in \{0, 1\}^\ell} (p_{1*} p_{*1}(z-1))^{k_1(f)} \left(\frac{p_{11}}{p_{1*} p_{*1}} \right)^{k_{11}(f)}. \end{aligned}$$

Because $Z_0 = \vec{1} \vec{1}^T$, $\vec{1}^T \left(\prod_{i \in [\ell]} PZ_{f(i)} \right) P \vec{1}$ is the product of terms of the form $\vec{1}^T (PZ_1)^j P \vec{1} = p_{01}(z-1)^j p_{11}^{j-1} p_{10}$. Each run of j consecutive ones in f contributes j to $k_1(f)$ and $j-1$ to $k_{11}(f)$, which gives us (a).

The identity for a_ℓ° follows analogously. \square

A.2 Inequalities

Lemma 10. For all $\ell \geq 2$, $x \in \mathbb{R}$ and $y \geq 1$, $b_\ell^\circ(x, y) \leq b_2^\circ(x, y)^{\ell/2}$.

Proof. Let $B = \begin{pmatrix} 1 & 1 \\ x & xy \end{pmatrix}$ and observe that

$$b_\ell^\circ(x, y) = \sum_{f \in [2]^\ell} \prod_{i \in [\ell]} B_{i, (i+1) \bmod \ell} = \text{tr}(B^\ell).$$

Let λ_0 and λ_1 be the eigenvalues of B . When

$$(\lambda_0 - \lambda_1)^2 = \text{tr}(B)^2 - 4 \det(B) = y^2 x^2 + (4 - 2y)x + 1 \geq 0,$$

the eigenvalues are real. The discriminant of this quadratic is $(4 - 2y)^2 - 4y^2 = 16 - 16y$, which is negative when $y \geq 1$. Thus when $y \geq 1$, the eigenvalues are real for all x . We have the Jordan decomposition $B = C^{-1} \Lambda C$ where Λ is upper triangular and has diagonal entries λ_0 and λ_1 . Then

$$\text{tr}(B^\ell) = \text{tr}(\Lambda^\ell) = \lambda_0^\ell + \lambda_1^\ell \leq |\lambda_0|^\ell + |\lambda_1|^\ell \leq (\lambda_0^2 + \lambda_1^2)^{\ell/2}$$

from the standard inequality between p -norms. \square

Lemma 11. *For all $\ell \geq 1$, all $x \geq 0$, and all $y \geq 1$, $b_\ell(x, y)^2 \leq b_2^\circ(x, y)^\ell$.*

Proof. First, observe that for all $f \in \{0, 1\}^\ell$, $k_{11}(f) \leq k_{11}^\circ(f)$. Thus for all $\ell \geq 2$, $j \in \mathbb{N}$, and $y \geq 1$, we have the stronger inequality $[x^k]b_\ell(x, y) \leq [x^j]b_\ell^\circ(x, y)$. Combining this with Lemma 10 gives the claim for $\ell \geq 2$. Finally, $b_1(x, y)^2 = (1 + x)^2 \leq 1 + 2x + x^2 y^2 = b_2^\circ(x, y)$. \square

Lemma 12. *If $\frac{p_{11}p_{00}}{p_{10}p_{01}} \geq 1$, then for all $\ell \geq 1$, $a_\ell(z)^2 \leq a_2^\circ(z)^\ell$ and for all $\ell \geq 2$, $a_\ell^\circ(z)^2 \leq a_2^\circ(z)^\ell$.*

Proof. We have $y = \frac{p_{11}}{p_{1*}p_{*1}} \geq 1$ if and only if $\frac{p_{11}p_{00}}{p_{10}p_{01}} \geq 1$. Then the claim follows from Lemmas 9, 10, and 11. \square

Proof of Lemma 2. Let $n' = |\mu|$. We have

$$\begin{aligned} A_{\mu, \mu^*}(z) &\stackrel{(a)}{=} \prod_{\ell \geq 1} a_\ell(z)^{t_{1,\ell}} a_\ell^\circ(z)^{t_{1,\ell}^\circ} a_\ell(z^2)^{t_{2,\ell}} a_\ell^\circ(z^2)^{t_{2,\ell}^\circ} \\ &\stackrel{(b)}{\leq} a_2^\circ(z)^{t_{1,1}/2} a_1^\circ(z)^{t_{1,1}^\circ} a_2^\circ(z^2)^{t_{2,1}/2} a_1^\circ(z^2)^{t_{2,1}^\circ} \\ &\quad \cdot \prod_{\ell \geq 2} (a_2^\circ(z)^{\ell/2})^{t_{1,\ell} + t_{1,\ell}^\circ} (a_2^\circ(z^2)^{\ell/2})^{t_{2,\ell} + t_{2,\ell}^\circ} \\ &\stackrel{(c)}{=} a_1^\circ(z^2)^{t_{2,1}^\circ} a_2^\circ(z)^{d(n'-d)/2} a_2^\circ(z^2)^{\binom{d}{2} - t_{2,1}^\circ} \end{aligned} \tag{11}$$

where (a) follows from Lemma 8, (b) follows from Lemma 12, and (c) uses the facts $t_{1,1}^\circ = 0$, $\sum_{\ell \geq 1} (t_{1,\ell} + t_{1,\ell}^\circ) = d(n' - d)$, and $\sum_{\ell \geq 1} (t_{2,\ell} + t_{2,\ell}^\circ) = \binom{d}{2}$.

We can easily compute each of the factors in (11):

$$\begin{aligned} a_1^\circ(z) &= \text{tr}(PZ_0) + \text{tr}(PZ_1) \\ &= 1 + p_{11}(z - 1) \\ a_2^\circ(z) &= \text{tr}(PZ_0 PZ_0) + 2 \text{tr}(PZ_1 PZ_0) + \text{tr}(PZ_1 PZ_1) \\ &= 1 + 2p_{1*}p_{*1}(z - 1) + p_{11}^2(z - 1)^2. \end{aligned}$$

Now we will bound each factor:

$$\begin{aligned}
a_1^\circ(z^2) &\leq \exp(p_{11}(z^2 - 1)) \\
a_2^\circ(z^2) &\leq \exp(2p_{1*}p_{*1}(z^2 - 1) + p_{11}^2(z^2 - 1)^2) \\
a_2^\circ(z)^2 &\leq \exp(4p_{1*}p_{*1}(z - 1) + 2p_{11}^2(z - 1)^2) \\
&= \exp(4(p_{1*}p_{*1} - p_{11}^2)(z - 1) + 2p_{11}^2(z^2 - 1)) \\
&\stackrel{(a)}{\leq} \exp(2(p_{1*}p_{*1} - p_{11}^2)(z^2 - 1) + p_{11}^2(z^4 - 1)) \\
&= \exp(2p_{1*}p_{*1}(z^2 - 1) + p_{11}^2(z^2 - 1)^2).
\end{aligned}$$

where (a) uses the inequality $x^2 - 1 = (x - 1)^2 + 2(x - 1) \geq 2(x - 1)$. Combining these with (11) and

$$\frac{d(n' - d)}{4} + \frac{d(d - 1) - 2t_{2,1}^\circ}{4} = \frac{d(n' - 1) - 2t_{2,1}^\circ}{4}$$

gives the claimed bound. \square

B Proof of Lemma 3

Proof. The optimal choice of z satisfies

$$\begin{aligned}
0 &= (2q_2z + q_1) \exp(q_2(z^2 - 1) + q_1(z - 1))z^{-\tau} \\
&\quad - \tau \exp(q_2(z^2 - 1) + q_1(z - 1))z^{-\tau-1} \\
0 &= 2q_2z^2 + q_1z - \tau
\end{aligned} \tag{12}$$

The equation (12) has one positive root and one negative root. The positive root is

$$z^* = \frac{-q_1 + \sqrt{q_1^2 + 8\tau q_2}}{4q_2} = \frac{2\tau}{q_1 + \sqrt{q_1^2 + 8\tau q_2}}.$$

Because $q_1 \leq \sqrt{q_1^2 + 8\tau q_2}$, we have the bounds

$$\frac{\tau}{\sqrt{q_1^2 + 8\tau q_2}} \leq z^* \leq \frac{\tau}{q_1}. \tag{13}$$

Starting with one of the factors from the left side of (8), we have

$$\exp(q_2(z^2 - 1) + q_1(z - 1)) = \exp\left(\frac{q_1}{2}z + \frac{\tau}{2} - q_2 - q_1\right) \leq e^{\tau - q_2 - q_1} \leq e^\tau$$

where we used (12) to eliminate the q_2z^2 term, applied the upper bound from (13), and used $q_1 \geq 0$ and $q_2 \geq 0$. From the lower bound in (13),

$$z^{-2} \leq \frac{q_1^2}{\tau^2} + \frac{8q_2}{\tau} \leq \max\left(\frac{2q_1^2}{\tau^2}, \frac{16q_2}{\tau}\right)$$

so $\exp(q_2(z^2 - 1) + q_1(z - 1))z^{-\tau} \leq \zeta^\tau$. \square

C Proof of Lemma 4

Proof. For $\mu \in \mathcal{M}(\mu^*, d)$ and $z > 0$,

$$\begin{aligned} \Pr[\mathbf{M}_{\mu, \mu^*} \geq kd] &\leq z^{-kd} A_{\mu, \mu^*}(z) \\ &\leq (z^2)^{-\tau} \exp(q_1(z^2 - 1) + q_2(z^4 - 1)) \end{aligned}$$

where we have used Lemma 2 and

$$q_2 = \frac{\tilde{t}}{4} p_{11}^2 \quad q_1 = t_{2,1}^\circ p_{11} + \frac{\tilde{t}}{2} (p_{1^*} p_{*1} - p_{11}^2) \quad \tau = \frac{dk}{2}.$$

Applying Lemma 3 we obtain $\Pr[\mathbf{M}_{\mu, \mu^*} \geq kd] \leq \zeta^\tau$, where

$$\zeta = \max \left(\sqrt{2} e \frac{q_1}{\tau}, 4e \left(\frac{q_2}{\tau} \right)^{\frac{1}{2}} \right).$$

Thus $\frac{1}{d} \log \Pr[\mathbf{M}_{\mu, \mu^*} \geq kd] \leq \frac{k}{2} \log \zeta$. We have

$$\begin{aligned} q_2 &\leq \frac{dn}{4} p_{11}^2 \\ q_1 &\stackrel{(a)}{\leq} \frac{d}{2} p_{11} + \frac{dn}{2} (p_{1^*} p_{*1} - p_{11}^2) \\ \tau &\geq \Omega(dn p_{11}) \\ \frac{q_2}{\tau} &\leq \mathcal{O}(p_{11}) \\ \frac{q_1}{\tau} &\leq \mathcal{O} \left(\frac{1}{n} + \frac{p_{1^*} p_{*1} - p_{11}^2}{p_{11}} \right) \end{aligned}$$

where (a) uses the fact that $t_{2,1}^\circ$ is equal to the number of cycles of length four in $\mu + \mu^*$, so it is at most $d/2$. From condition (3) and $p_{00} \geq \Omega(1)$ we have

$$\frac{p_{1^*} p_{*1} - p_{11}^2}{p_{11}} = \frac{p_{01} p_{10}}{p_{11}} + p_{01} + p_{10} \leq n^{-\Omega(1)}.$$

To handle the case where ζ is equal to the first entry of the maximum, we have

$$\begin{aligned} \frac{k}{2} \log \left(\frac{\tau}{\sqrt{2} e q_1} \right) &\geq \Omega(n p_{11}) \log(n^{\Omega(1)}) \\ &\geq \omega(\log n). \end{aligned}$$

In the second case, when $\omega(\frac{1}{n}) \leq p_{11} \leq n^{-\Omega(1)}$, we have

$$\begin{aligned} \frac{k}{2} \log \left(\left(\frac{\tau}{16e^2 q_2} \right)^{\frac{1}{2}} \right) &\geq \Omega(n p_{11}) \log \left(\frac{1 - o(1)}{8e^2 p_{11}} \right) \\ &\geq \omega(1) \log(n^{\Omega(1)}) \\ &\geq \omega(\log n) \end{aligned}$$

The function $f(x) = -x \log(8e^2 x)$ is increasing on the interval $0 \leq x \leq \frac{1}{8e^3}$, so condition (2) can replace $p_{11} \leq n^{-\Omega(1)}$. \square

References

- [1] E. M. Wright, “Graphs on unlabelled nodes with a given number of edges,” *Acta Mathematica*, vol. 126, no. 1, pp. 1–9, 1971.
- [2] P. Pedarsani and M. Grossglauser, “On the privacy of anonymized networks,” in *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2011, pp. 1235–1243.
- [3] D. Cullina and N. Kiyavash, “Exact alignment recovery for correlated Erdős Rényi graphs,” *arXiv:1711.06783 [cs, math]*, Nov. 2017, arXiv: 1711.06783. [Online]. Available: <http://arxiv.org/abs/1711.06783>
- [4] —, “Improved achievability and converse bounds for Erdős-Rényi graph matching,” in *Proceedings of the 2016 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Science*. ACM, 2016, pp. 63–72.
- [5] E. Kazemi, L. Yartseva, and M. Grossglauser, “When Can Two Unlabeled Networks Be Aligned Under Partial Overlap?” in *Proceedings of the 53rd Annual Allerton Conference on Communication, Control, and Computing*, 2015.
- [6] F. Shirani, S. Garg, and E. Erkip, “Typicality Matching for Pairs of Correlated Graphs,” *arXiv preprint arXiv:1802.00918*, 2018.
- [7] D. Cullina, P. Mittal, and N. Kiyavash, “Fundamental Limits of Database Alignment,” *arXiv:1805.03829 [cs, math]*, May 2018, arXiv: 1805.03829. [Online]. Available: <http://arxiv.org/abs/1805.03829>
- [8] S. Ji, W. Li, N. Z. Gong, P. Mittal, and R. Beyah, “On Your Social Network De-anonymizability: Quantification and Large Scale Evaluation with Seed Knowledge,” 2015.
- [9] W.-H. Lee, C. Liu, S. Ji, P. Mittal, and R. B. Lee, “Blind de-anonymization attacks using social networks,” in *Proceedings of the 2017 on Workshop on Privacy in the Electronic Society*. ACM, 2017, pp. 1–4.
- [10] L. Backstrom, C. Dwork, and J. Kleinberg, “Wherefore art thou r3579x?: anonymized social networks, hidden patterns, and structural steganography,” in *Proceedings of the 16th international conference on World Wide Web*. ACM, 2007, pp. 181–190.
- [11] N. Korula and S. Lattanzi, “An efficient reconciliation algorithm for social networks,” *Proceedings of the VLDB Endowment*, vol. 7, no. 5, pp. 377–388, 2014.
- [12] N. Malod-Dognin and N. Prulj, “L-GRAAL: Lagrangian graphlet-based network aligner,” *Bioinformatics*, vol. 31, no. 13, pp. 2182–2189, 2015.
- [13] O. Kuchaiev, T. Milenkovi, V. Memievi, W. Hayes, and N. Prulj, “Topological network alignment uncovers biological function and phylogeny,” *Journal of the Royal Society Interface*, p. rsif20100063, 2010.
- [14] R. Singh, J. Xu, and B. Berger, “Global alignment of multiple protein interaction networks with application to functional orthology detection,” *Proceedings of the National Academy of Sciences*, 2008.

- [15] Y. Tian and J. M. Patel, “TALE: A tool for approximate large graph matching,” in *Data Engineering, 2008. ICDE 2008. IEEE 24th International Conference on*. IEEE, 2008, pp. 963–972.
- [16] S. Zhang, J. Yang, and W. Jin, “SAPPER: subgraph indexing and approximate matching in large graphs,” *Proceedings of the VLDB Endowment*, vol. 3, no. 1-2, pp. 1185–1194, 2010.
- [17] L. Yartseva and M. Grossglauser, “On the performance of percolation graph matching,” in *Proceedings of the first ACM conference on Online social networks*. ACM, 2013, pp. 119–130.
- [18] E. Kazemi, S. H. Hassani, and M. Grossglauser, “Growing a Graph Matching from a Handful of Seeds,” in *Proceedings of the Vldb Endowment International Conference on Very Large Data Bases*, vol. 8, 2015.
- [19] P. Pedarsani, D. R. Figueiredo, and M. Grossglauser, “A bayesian method for matching two similar graphs without seeds,” in *Communication, Control, and Computing (Allerton), 2013 51st Annual Allerton Conference on*. IEEE, 2013, pp. 1598–1607.
- [20] V. Lyzinski, D. E. Fishkind, and C. E. Priebe, “Seeded graph matching for correlated Erdős-Rényi graphs.” *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 3513–3540, 2014.
- [21] S. Feizi, G. Quon, M. Recamonde-Mendoza, M. Medard, M. Kellis, and A. Jadbabaie, “Spectral Alignment of Graphs,” *arXiv:1602.04181 [cs, math]*, Feb. 2016, arXiv: 1602.04181. [Online]. Available: <http://arxiv.org/abs/1602.04181>
- [22] V. Lyzinski, D. Fishkind, M. Fiori, J. Vogelstein, C. Priebe, and G. Sapiro, “Graph matching: Relax at your own risk,” *IEEE Transactions on Pattern Analysis & Machine Intelligence*, no. 1, pp. 1–1, 2016.
- [23] B. Barak, C.-N. Chou, Z. Lei, T. Schramm, and Y. Sheng, “(Nearly) Efficient Algorithms for the Graph Matching Problem on Correlated Random Graphs,” *arXiv preprint arXiv:1805.02349*, 2018.
- [24] E. Mossel and J. Xu, “Seeded Graph Matching via Large Neighborhood Statistics,” *arXiv preprint arXiv:1807.10262*, 2018.
- [25] C. Moore, “The Computer Science and Physics of Community Detection: Landscapes, Phase Transitions, and Hardness,” Feb. 2017. [Online]. Available: <https://arxiv.org/abs/1702.00467>
- [26] E. Abbe, “Community detection and stochastic block models: recent developments,” *arXiv:1703.10146 [cs, math, stat]*, Mar. 2017, arXiv: 1703.10146. [Online]. Available: <http://arxiv.org/abs/1703.10146>
- [27] Y. Wu and J. Xu, “Statistical Problems with Planted Structures: Information-Theoretical and Computational Limits,” *arXiv:1806.00118 [cs, math, stat]*, May 2018, arXiv: 1806.00118. [Online]. Available: <http://arxiv.org/abs/1806.00118>
- [28] B. Bollobás, *The evolution of sparse graphs, Graph Theory and Combinatorics (Cambridge 1983)*, 35-57. Academic Press, London, 1984.
- [29] T. Łuczak, “Size and connectivity of the k-core of a random graph,” *Discrete Mathematics*, vol. 91, no. 1, pp. 61–68, 1991.