

Approximation Algorithms for Orienting Mixed Graphs

Michael Elberfeld^{*,1}, Danny Segev^{*,2}, Colin R. Davidson³, Dana Silverbush⁴,
and Roded Sharan⁴

¹ Institute of Theoretical Computer Science, University of Lübeck, 23538 Lübeck,
Germany, elberfeld@tcs.uni-luebeck.de

² Department of Statistics, University of Haifa, Haifa 31905, Israel,
segevd@stat.haifa.ac.il

³ Faculty of Mathematics, University of Waterloo, Waterloo, Canada, N2L 3G1,
colinrdavidson@gmail.com

⁴ Blavatnik School of Computer Science, Tel Aviv University, Tel Aviv 69978, Israel,
{danasilv,roded}@post.tau.ac.il

Abstract. Graph orientation is a fundamental problem in graph theory that has recently arisen in the study of signaling-regulatory pathways in protein networks. Given a graph and a list of ordered source-target vertex pairs, it calls for assigning directions to the edges of the graph so as to maximize the number of pairs that admit a directed source-to-target path. When the input graph is undirected, a sub-logarithmic approximation is known for the problem. However, the approximability of the biologically-relevant variant, in which the input graph has both directed and undirected edges, was left open. Here we give the first approximation algorithm to this problem. Our algorithm provides a sub-linear guarantee in the general case, and logarithmic guarantees for structured instances.

Key words: protein-protein interaction network, mixed graph, graph orientation, approximation algorithm

1 Introduction

Protein-protein interactions (PPIs) form the skeleton of signal transduction in the cell. While many of these interactions carry directed signaling information, current PPI measurement technologies, such as yeast two hybrid [10] and co-immunoprecipitation [14], cannot reveal the direction in which the signal flows. The problem of inferring this hidden directionality information is fundamental to our understanding of how these networks function. Previous work on this problem has relied on information from perturbation experiments [23], in which a gene is perturbed (cause) and as a result other genes change their expression levels (effects), to guide the orientation inference. Specifically, it is assumed that for an effect to take place, there must be a directed path in the network from the causal gene to the affected gene. The arising combinatorial problem is to

* These authors contributed equally to this work.

orient the edges of the network such that a maximum number of cause-effect pairs admit a directed path from the causal to the affected gene. When studying a PPI network in isolation, the input network is undirected. However, the more biologically relevant variant considers also protein-DNA interactions as these are necessary to explain the expression changes. Moreover, the directionality of some PPIs, like kinase-substrate interactions, is known in advance. Thus, in general, the input network is a mixed graph containing both directed and undirected edges.

The optimization problem that we study draws its recent interest from applications in network biology, but is rooted at practical applications from already a century ago: In 1939, Robbins [20], who was motivated by applications in street network design, showed that an undirected graph has a strongly connected orientation if and only if it has no bridge edge. The corresponding decision problem can be solved in linear time [22]. The characterization of Robbins was extended to mixed graphs by Boesch et al. [5]; linear time algorithms for deciding whether a mixed graph admits a strongly connected orientation were presented by Chung et al. [7]. Hakimi et al. [15] presented a polynomial algorithm for the problem of orienting an undirected graph so as to maximize the number of source-target pairs out of all possible ordered vertex pairs that admit a directed source-to-target path. A recent work by Dorn et al. [9] studies the parameterized complexity of orienting graphs. We refer to the textbook of Bang-Jensen and Gutin [3] for a comprehensive discussion of various graph orientation problems.

More recently, the problem of network orientation has been motivated by applications in network biology. Medvedovsky et al. [18] who formulated the problem that we study here, focused on restricted instances where the input graph is undirected, providing a logarithmic approximation algorithm for the problem. The approximation guarantee was later improved to $\Omega(\log \log n / \log n)$ by Gamzu et al. [13], where n denotes the number of vertices in the input graph. Gamzu et al. also showed that the orientation problem on mixed graphs can be approximated to within a poly-logarithmic ratio of $\Omega(1/\log^l n)$ where l is the maximum number of alternations between undirected and directed edges on a source-to-target path. Silverbush et al. [21] developed an ILP-based algorithm to optimally orient mixed networks, but the approximability of the problem (for non-constant l) was left open.

In this work, we study the approximability of the orientation problem on mixed graphs. We show that the problem is NP-hard to approximate to within a factor of $7/8$. We then reduce the problem to orienting acyclic mixed graphs. We provide logarithmic approximation guarantees for tree-like reduced instances and use those to develop a sub-linear approximation algorithm for general instances.

The paper is organized as follows: In the next section we formally define the orientation problem, discuss its complexity and describe a generic reduction to acyclic mixed graphs. In Section 3 we present logarithmic factor approximation algorithms for tree-like instances. Section 4 presents the sub-linear approximation algorithm for the general case.

2 Preliminaries

Notation and terminology. We focus on simple graphs with no loops or parallel edges. A *mixed graph* is a triple $G = (V, E_U, E_D)$ that consists of a vertex set V , a set of *undirected edges* $E_U \subseteq \{e \subseteq V \mid |e| = 2\}$, and a set of *directed edges* $E_D \subseteq V \times V$. We assume that every pair of vertices is either connected by a single edge of a specific type (directed or undirected) or not connected at all. We also write $V(G)$, $E_U(G)$, and $E_D(G)$ to refer to the sets V , E_U , and E_D , respectively. When G is clear from the context, we will denote $n = |V|$.

Let G_1 and G_2 be two mixed graphs. The graph G_1 is a *subgraph* of G_2 when the relations $V(G_1) \subseteq V(G_2)$, $E_U(G_1) \subseteq E_U(G_2)$, and $E_D(G_1) \subseteq E_D(G_2)$ hold. A *path* of length ℓ in a mixed graph G is a sequence $p = \langle v_1, v_2, \dots, v_\ell, v_{\ell+1} \rangle$ of distinct vertices such that for every $1 \leq i \leq \ell$, we have $\{v_i, v_{i+1}\} \in E_U(G)$ or $(v_i, v_{i+1}) \in E_D(G)$. It *crosses* a vertex $v \in V(G)$ when $v = v_i$ for some $i \in \{1, \dots, \ell + 1\}$. It is a *cycle* when $v_1 = v_{\ell+1}$. Given $s \in V(G)$ and $t \in V(G)$, we say that t is *reachable from* s when there exists a path in G that goes from s to t . In this case we also say that G *satisfies* the pair (s, t) . A mixed graph with no cycles is called a *mixed acyclic graph* (MAG).

Let G be a mixed graph. An *orientation* of G is a directed graph over the same vertex set, whose edge set contains all the directed edges of G and a single directed instance of every undirected edge, but nothing more. When only a subset of the undirected edges have been oriented, we obtain a *partial orientation*.

Problem statement. The MAXIMUM-MIXED-GRAPH-ORIENTATION problem is defined as follows:

Input: A mixed graph G , and a collection of source-target vertex pairs $P \subseteq V(G) \times V(G)$.

Output: An orientation of G that satisfies a maximum number of pairs from P .

Hardness result. Arkin and Hassin [1] showed that it is NP-complete to decide whether, for a given mixed graph G and a collection of source-target pairs P , the graph G can be oriented to satisfy all pairs in P . Their reduction, based on the 3-SATISFIABILITY problem, guarantees that for every $k \in \mathbb{N}$ there exists an assignment with k satisfied clauses if and only if there exists an orientation with k satisfied pairs. Thus, the inapproximability of MAXIMUM-3-SATISFIABILITY [16] directly transfers to MAXIMUM-MIXED-GRAPH-ORIENTATION, implying that it is NP-hard to approximate it to within a factor of $7/8$. We note that this bound is slightly lower than the $12/13$ -bound known for the special case where the input graph is undirected [18].

Reduction to mixed acyclic graphs. Given an orientation instance (G, P) , we can orient the undirected edges of any mixed cycle in a consistent direction without affecting the maximum number of source-target pairs that are satisfied by an optimal orientation. This observation gives rise to a polynomial-time reduction from mixed graphs to MAGs: First, we iteratively orient mixed cycles in the input

graph. Then, we contract strongly connected components into single vertices, and connect two components by an undirected (directed) edge when some vertex in the first component is connected by such an edge to a vertex in the second component (note that there cannot be more than one edge type as otherwise the two strongly connected components would have been merged). The pairs from P are adjusted accordingly from vertices of G to component vertices. A formal correctness proof of this reduction is given by Silverbush et al. [21].

Given a MAG G , the components of the undirected graph $(V(G), E_U(G))$ are called the *undirected components* of G ; they must be trees that are connected by directed edges from $E_D(G)$ without producing cycles. The *graph of undirected components* of G is the directed acyclic graph G_{UCC} with $V(G_{\text{UCC}}) = \{G_i \mid G_i \text{ is an undirected component of } G\}$, and there is a directed edge from a node G_i to a node G_j when there is an edge from some vertex $v \in G_i$ to some vertex $w \in G_j$.

By the reduction above, we may focus our attention on treating MAGs. In addition, we may assume that each of the input pairs can be satisfied by some orientation; otherwise, it can be eliminated without affecting the optimum solution. Thus, throughout the paper, all instances considered will be assumed to satisfy these two properties.

3 Logarithmic approximations for tree-like instances

In this section we provide logarithmic approximations that apply to orientation instances where the graph is “similar” to a tree, as formally defined in the sequel. In the remainder of this section, we make use of the following result about orienting undirected trees due to Medvedovsky et al. [18]:

Lemma 3.1. *Let (G, P) be an orientation instance where G is an undirected tree. There is a polynomial-time algorithm that computes an orientation satisfying at least $|P|/(4\lceil \log n \rceil)$ pairs.*

3.1 Orienting mixed trees

The above lemma guarantees that a logarithmic fraction of the input pairs can always be satisfied, and since it is constructive, we immediately derive an $\Omega(1/\log n)$ approximation algorithm for undirected trees. The following sequences of claims are of similar nature: We prove the existence of orientations satisfying a certain fraction of all input pairs, and this leads to approximation algorithms with the corresponding ratio. We start with orientations for mixed trees.

Lemma 3.2. *Let (G, P) be an orientation instance with a mixed tree G . There is a polynomial-time algorithm that computes an orientation satisfying at least $|P|/(4\lceil \log n \rceil)$ of the pairs.*

Proof. First contract all directed edges of the tree into single vertices and update the end vertices of the input pairs accordingly. The resulting graph is an undirected tree with source-to-target paths for every pair in P . By Lemma 3.1, there exists a polynomial-time algorithm that finds an orientation of the resulting tree satisfying at least $1/(4\lceil \log n \rceil)$ of the pairs. Carrying over the edge directions to the initial mixed tree produces an orientation that satisfies exactly the same collection of pairs in the original orientation instance. \square

3.2 Crossings through a junction component

Let (G, P) be an orientation instance and let T_1, T_2, \dots be the undirected components of G . We construct a subgraph of G , called the *skeleton* $S = S(G)$ of G by deleting all but one directed edge between any pair of trees T_i and T_j . Note that the exact structure of a skeleton graph depends on the (polynomial-time) procedure used for its construction; we choose any fixed procedure to define S unambiguously. It is not difficult to verify that the skeleton S contains source-to-target paths for the pairs P , and that every orientation of S satisfying certain pairs directly translates into an orientation for G satisfying at least the same pairs. Figure 1 shows an example of a graph G and a skeleton for it. Note that

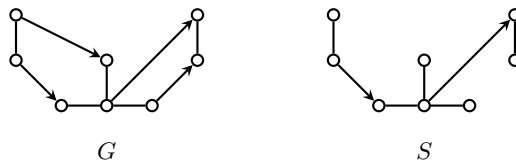


Fig. 1. An example MAG G and a skeleton S of it.

for any MAG G and its skeleton $S = S(G)$, we have $G_{\text{UCC}} = S_{\text{UCC}}$.

The next lemma is crucial to establish the remaining results of this section, as well as the sub-linear approximation algorithm described in Section 4.

Lemma 3.3. *Let (G, P) be an orientation instance and T be an undirected component of G . If each pair in P admits a source-to-target path that crosses a vertex from T then there is a polynomial-time algorithm that computes an orientation satisfying at least $|P|/(4\lceil \log n \rceil)$ of the pairs.*

Proof. Since the skeleton $S = S(G)$ is a MAG and, therefore, S_{UCC} is a directed acyclic graph, for every undirected component $T' \neq T$ of S_{UCC} , exactly one of the following options holds: (1) T' is reachable from T in S_{UCC} ; (2) T is reachable from T' in S_{UCC} ; or (3) there is no path between T and T' in either direction. Consequently, we can consider two subtrees that are rooted at T : The first subtree spans the vertices of S_{UCC} that are reachable from T , and the second subtree spans the vertices of S_{UCC} that can reach T . We merge both subtrees at T and call the resulting directed tree T_{UCC} . To compute an orientation for G , we

consider the subtree of S that is constructed by taking all undirected components from T_{UCC} , and connect two vertices in different components by a directed edge if this edge is already present in S and the components are connected in T_{UCC} . This subtree of S contains a source-to-target path for each pair in P . Therefore, by Lemma 3.2, we can construct (in polynomial time) an orientation satisfying at least $|P|/(4\lceil\log n\rceil)$ pairs in S and, thus, also in the original graph G . \square

Lemma 3.3 implies an $\Omega(1/\log n)$ approximation for a special case of the orientation problem, which we call MAXIMUM-JUNCTION-TREE-ORIENTATION. In Section 4 we shall apply the algorithm to instances where all pairs have source-to-target paths crossing a distinguished vertex r .

3.3 Orientations for small feedback vertex sets or treewidth

We end this section by providing logarithmic approximations to the orientation problem on tree-like instances. Precisely, we consider two graph parameters: *feedback vertex number* and *treewidth*, showing that whenever either one of these is bounded by a constant, it is possible to compute an orientation that satisfies a poly-logarithmic fraction of the input pairs.

Lemma 3.4. *Let (G, P) be an orientation instance where the underlying undirected subgraph of G_{UCC} can be turned into a tree by deleting at most k vertices. There is a polynomial-time algorithm that computes an orientation satisfying at least $|P|/(4(2k+1)\lceil\log n\rceil)$ pairs.*

Proof. We begin by detecting a small-sized feedback vertex set $F = \{T_1, \dots, T_\ell\}$, consisting of ℓ vertices whose removal turns the underlying undirected subgraph of G_{UCC} into a tree. Even though finding a minimum cardinality vertex set of this type is NP-hard [17], this problem can be approximated to within a factor of 2 in undirected graphs [2], implying that we can assume $\ell \leq 2k$. We now partition P into two subsets, the collection of pairs P^+ for which we can find source-to-target paths in G that cross undirected components from F , and the collection $P^- = P \setminus P^+$. We further partition P^+ into ℓ subsets P_1^+, \dots, P_ℓ^+ , where a pair $(s, t) \in P^+$ lies in P_i^+ if i is the minimal index for which there exists a source-to-target path for this pair that crosses the undirected component T_i . With these definitions at hand, note that by deleting the undirected components F from G , we can use Lemma 3.2 to efficiently compute an orientation of G satisfying at least $|P^-|/(4\lceil\log n\rceil)$ pairs; after deleting F the skeleton of the resulting graph is a tree and all pairs in P^- remain connected since they are only connected through paths that not visit vertices from F . On the other hand, for each collection P_i^+ we can satisfy at least $|P_i^+|/(4\lceil\log n\rceil)$ pairs by applying Lemma 3.3. Picking the option that generates the highest number of satisfied pairs results in an orientation satisfying at least $|P|/(4(2k+1)\lceil\log n\rceil)$ of the pairs in P . \square

We note that the above approximation result can be improved by a factor of 2 if the feedback vertex set has bounded size. For such instances we can invoke

an exact fixed parameter algorithm [6] to find an optimal feedback set, rather than using the 2-approximation algorithm.

Lemma 3.5. *Let (G, P) be an orientation instance where the underlying undirected subgraph of G_{UCC} has treewidth k . There is a polynomial-time algorithm that computes an orientation satisfying at least $|P|/(4(k+1)\lceil\log n\rceil^2)$ pairs.*

Proof. We first compute a tree decomposition of width k for the undirected underlying graph of G_{UCC} . A tree decomposition $(\mathcal{T}, \{B_t\}_{t \in V(\mathcal{T})})$ consists of a tree \mathcal{T} whose nodes are labeled with possibly-overlapping subsets B_t of vertices, called *bags*, such that: (1) the incident vertices of every edge are both contained in some bag; and (2) for every original vertex, the nodes of the bags that contain it make up a connected subtree. Its width is defined as the maximum number of vertices in a bag minus 1. For a comprehensive discussion on tree decompositions and their polynomial-time computability in the case of bounded tree width, we refer to the book of Flum and Grohe [11].

Based on the tree decomposition $(\mathcal{T}, \{B_t\}_{t \in V(\mathcal{T})})$, we partition P into subsets P_1, P_2, \dots, P_L with $L \leq \lceil\log n\rceil$ such that for every subset we can efficiently find an orientation that satisfies a fraction of at least $1/(4(k+1)\lceil\log n\rceil)$ of its pairs. By picking the largest subset of pairs and its corresponding orientation, we obtain an orientation satisfying at least $|P|/(4(k+1)\lceil\log n\rceil^2)$ pairs.

For the purpose of constructing P_1 , consider a *centroid* node t of \mathcal{T} whose removal breaks this tree into subtrees of cardinality at most $|V(\mathcal{T})|/2$, noting that any tree necessarily contains a centroid (see, for instance, [12]). Let P_1 be the pairs in P with source-to-target paths that cross vertices from undirected components of $B_t = \{T_1, \dots, T_l\}$, where $l \leq k+1$. We further partition P_1 into l collections P_1^1, \dots, P_1^l such that a pair $(s, t) \in P_1$ lies in P_1^i if there exists an s - t path that crosses vertices from T_i but no s - t paths that cross vertices from components T_j with $j < i$. By Lemma 3.3, we can compute an orientation that satisfies at least $|P_1^i|/(4\lceil\log n\rceil)$ of the pairs in P_1^i , for every $1 \leq i \leq l$. By taking the largest collection, we can satisfy at least $|P_1|/(4(k+1)\lceil\log n\rceil)$ pairs in P_1 .

To construct P_2 , we proceed with the pair collection $P \setminus P_1$ that contains exactly the pairs from P with no source-to-target paths that cross vertices from the components of P_1 . We delete the node t from \mathcal{T} , as well as the components in B_t from G . This results in a graph that contains source-to-target paths for all pairs from $P \setminus P_1$ and a forest of tree decompositions for the graph. For each tree decomposition we compute a centroid bag and, in the same way as above, the collection P_2 of pairs in $P \setminus P_1$ with source-to-target paths that cross components from these centroid bags. Using the same arguments as above, we can compute an orientation that satisfies at least $|P_2|/(4(k+1)\lceil\log n\rceil)$ of the pairs in P_2 . We proceed recursively in the same way to construct P_3, P_4, \dots, P_L as long as each tree decomposition (and the corresponding subgraph of G) is not empty. Since the maximal size of a subtree decreases by a factor of at least 2 in each level of the recursion, this process terminates within $\lceil\log n\rceil$ steps. \square

4 Sub-linear approximations for general instances

In what follows, we focus our attention on approximating the orientation problem in its utmost generality, that is, without making simplifying structural assumptions on the underlying (mixed-acyclic) graph G and on the collection of input pairs P . The main result of this section can be briefly stated as follows.

Theorem 4.1. *The MAXIMUM-MIXED-GRAPH-ORIENTATION problem can be approximated within a factor of $\Omega(1/(M^{1/\sqrt{2}} \log n))$, where $M = \max\{n, |P|\}$.*

In addition, we provide an improved approximation guarantee for input instances with bounded-distance pairs. This result is described in Section 4.3.

4.1 The algorithm

For each pair $(s_i, t_i) \in P$, let p_i be a shortest path from s_i to t_i in G , and let \mathcal{P} be the set of all shortest paths, i.e., $\mathcal{P} = \{p_i : (s_i, t_i) \in P\}$. Our algorithm is based on a greedy framework where paths in \mathcal{P} are oriented (from source to target) one after the other, trying not to interfere with future orientations of too many other paths by picking the shortest path in each step. Somewhat informally, this process concludes as soon as one of the following termination conditions is met:

The greedy step. At any point in time, we will be holding a partial orientation G_ℓ of G and a subset $\mathcal{P}_\ell \subseteq \mathcal{P}$ of shortest paths, where these sets are indexed according to the step number that has just been completed. In other words, at the conclusion of step ℓ we have G_ℓ and \mathcal{P}_ℓ , where initially $G_0 = G$ and $\mathcal{P}_0 = \mathcal{P}$. Now, as long as none of the termination conditions described below is met, we proceed as follows:

1. Let $\hat{p} = \langle s, \dots, t \rangle$ be a shortest path in \mathcal{P}_ℓ .
2. Orient \hat{p} in the direction from s to t to obtain $G_{\ell+1}$.
3. Discard from \mathcal{P}_ℓ the path \hat{p} as well as any path that has a non-empty edge intersection with \hat{p} . This way, we obtain $\mathcal{P}_{\ell+1}$.

Termination conditions. There are two conditions that will cause the greedy iterations to terminate. For now, we state both conditions in terms of two parameters $\alpha \geq 0$ and $\beta \geq 0$, whose values will be optimized later on.

Condition 1: $|\mathcal{P}_\ell| \leq n^\alpha$. In this case, we will orient an arbitrary path from \mathcal{P}_ℓ , and update the current orientation to G_ℓ , as in the preceding greedy iterations. We then complete the orientation by arbitrarily orienting all yet-unoriented edges.

Condition 2: There exists a vertex r such that at least $|\mathcal{P}_\ell|^\beta$ paths in \mathcal{P}_ℓ go through r . We construct a MAXIMUM-JUNCTION-TREE-ORIENTATION instance with input graph G_ℓ , junction vertex r , and pairs $\{(s_i, t_i) : p_i \in \mathcal{P}_\ell \text{ goes through } r\}$. We then apply the algorithm described in Section 3.2 for this special case, and return its output as our final orientation.

4.2 Analysis

To establish a lower bound on the number of satisfied pairs, we break the analysis into two cases, depending on the condition that caused the greedy iterations to terminate. In the remainder of this section, we assume that L greedy iterations have been completed prior to satisfying one of the termination conditions.

Connections due to condition 1: In this case we satisfy a single pair out of $\{(s_i, t_i) : p_i \in \mathcal{P}_L\}$, noting that $|\mathcal{P}_L| \leq n^\alpha$.

Connections due to condition 2: Following Lemma 3.3, the number of pairs satisfied out of $\{(s_i, t_i) : p_i \in \mathcal{P}_L\}$ is $\Omega(1/\log n) \cdot |\mathcal{P}_L|^\beta$.

We proceed by arguing that an $\Omega(1/n^{1-\alpha(1-2\beta)})$ fraction of the pairs in $\{(s_i, t_i) : p_i \notin \mathcal{P}_L\}$ are already satisfied by the partial orientation G_L . To this end, note that in each iteration $1 \leq \ell \leq L$ we satisfy a single pair by orienting the shortest path $\hat{p} \in \mathcal{P}_{\ell-1}$, and eliminating several others to obtain \mathcal{P}_ℓ . To prove the claim above, it is sufficient to show that the number of eliminated paths satisfies $|\mathcal{P}_{\ell-1} \setminus \mathcal{P}_\ell| \leq n^{1-\alpha(1-2\beta)}$. Denote by $E(p)$ the set of edges of a path p , so that $|E(p)|$ is its length. We begin by observing that, since condition 2 has not been met in iteration ℓ , each edge can have at most $|\mathcal{P}_{\ell-1}|^\beta$ paths from $\mathcal{P}_{\ell-1}$ going through it, implying that $|\mathcal{P}_{\ell-1} \setminus \mathcal{P}_\ell| \leq |E(\hat{p})| \cdot |\mathcal{P}_{\ell-1}|^\beta$. Since $|E(\hat{p})|$ is upper bounded by the average length of the paths in $\mathcal{P}_{\ell-1}$, we have

$$\begin{aligned} |E(\hat{p})| &\leq \frac{1}{|\mathcal{P}_{\ell-1}|} \sum_{p_i \in \mathcal{P}_{\ell-1}} |E(p_i)| \leq \frac{1}{|\mathcal{P}_{\ell-1}|} \sum_{p_i \in \mathcal{P}_{\ell-1}} |V(p_i)| \\ &= \frac{1}{|\mathcal{P}_{\ell-1}|} \sum_{v \in V} |\{p_i \in \mathcal{P}_{\ell-1} : v \in V(p_i)\}| \\ &\leq \frac{1}{|\mathcal{P}_{\ell-1}|} \cdot n \cdot |\mathcal{P}_{\ell-1}|^\beta = \frac{n}{|\mathcal{P}_{\ell-1}|^{1-\beta}}, \end{aligned}$$

where the third inequality holds since condition 2 has not been met. Hence,

$$|\mathcal{P}_{\ell-1} \setminus \mathcal{P}_\ell| \leq \frac{n}{|\mathcal{P}_{\ell-1}|^{1-2\beta}} \leq \frac{n}{n^{\alpha(1-2\beta)}} = n^{1-\alpha(1-2\beta)},$$

where the second inequality follows from $|\mathcal{P}_{\ell-1}| > n^\alpha$, as condition 1 has not been met.

Putting it all together. Based on the above discussion, it follows that the number of satisfied pairs when we terminate the algorithm due to condition 1 is

$$\begin{aligned} \Omega\left(\frac{1}{n^{1-\alpha(1-2\beta)}}\right) (|P| - |\mathcal{P}_L|) + 1 &= \Omega\left(\frac{1}{n^{1-\alpha(1-2\beta)}}\right) (|P| - n^\alpha) + \frac{1}{n^\alpha} n^\alpha \\ &= \Omega\left(\frac{1}{\max\{n^{1-\alpha(1-2\beta)}, n^\alpha\}}\right) |P| \\ &= \Omega\left(\frac{1}{n^{\max\{1-\alpha(1-2\beta), \alpha\}}}\right) |P|. \end{aligned}$$

Similarly, the number of satisfied pairs when the algorithm is terminated due to condition 2 is

$$\begin{aligned}
& \Omega\left(\frac{1}{n^{1-\alpha(1-2\beta)}}\right)(|P| - |\mathcal{P}_L|) + \Omega\left(\frac{1}{\log n}\right)|\mathcal{P}_L|^\beta \\
&= \Omega\left(\frac{1}{n^{1-\alpha(1-2\beta)}}\right)(|P| - |\mathcal{P}_L|) + \Omega\left(\frac{1}{|\mathcal{P}_L|^{1-\beta} \log n}\right)|\mathcal{P}_L| \\
&= \Omega\left(\frac{1}{\max\{n^{1-\alpha(1-2\beta)}, |P|^{1-\beta} \log n\}}\right)|P| \\
&= \Omega\left(\frac{1}{M^{\max\{1-\alpha(1-2\beta), 1-\beta\}}}\right)\frac{1}{\log n}|P|.
\end{aligned}$$

To obtain the best-possible performance guarantee, we pick values for α and β so as to minimize $\max\{\alpha, 1 - \beta, 1 - \alpha(1 - 2\beta)\}$. As explained below, the last term is optimized for $\alpha^* = \sqrt{1/2}$ and $\beta^* = \sqrt{1/2}/(2\sqrt{1/2} + 1) = 1 - \sqrt{1/2}$, in which case its value is $\sqrt{1/2} \approx 0.707$.

Optimizing α and β . Suppose we know the value of α^* . In this case, β^* should be picked so as to minimize $\max\{1 - \beta, 1 - \alpha^*(1 - 2\beta)\}$. Since $1 - \beta$ is a decreasing linear function of β and $1 - \alpha^*(1 - 2\beta)$ is an increasing linear function, this minimum is attained when $1 - \beta = 1 - \alpha^*(1 - 2\beta)$, that is, $\beta^* = \alpha^*/(2\alpha^* + 1)$. For this value, we have $\min \max\{1 - \beta, 1 - \alpha^*(1 - 2\beta)\} = (\alpha^* + 1)/(2\alpha^* + 1)$. It remains to find a value of α that minimizes $\max\{\alpha, (\alpha + 1)/(2\alpha + 1)\}$. Using similar arguments, it is not difficult to verify that the right value to pick is $\alpha^* = \sqrt{1/2}$.

4.3 An improved approximation for bounded-distance pairs

In practice, the diameter of biological networks is sub-logarithmic due to their scale-free property [4, 8, 19]. For example, in the yeast physical network described in [21], the maximum source-target distance is 14. This motivates examining the approximation guarantee in terms of the maximum length of a shortest source-target path in the reduced mixed acyclic graph, which we denote by $\Delta = \Delta(G, P)$. In the following we present an $\Omega(1/\sqrt{\Delta}|P| \log n)$ approximation to the orientation problem.

Our algorithm remains essentially unchanged, except for its termination conditions. Unlike the more general procedure, we ignore condition 1, and terminate the greedy iterations as soon as condition 2 is met, i.e., when there exists a vertex $r \in V$ such that at least $|\mathcal{P}_\ell|^\beta$ paths in \mathcal{P}_ℓ go through r . In this case, we construct a MAXIMUM-JUNCTION-TREE-ORIENTATION instance as before, with input graph G_ℓ , junction vertex r , and pairs $\{(s_i, t_i) : p_i \in \mathcal{P}_\ell \text{ goes through } r\}$. Our logarithmic approximation for this particular setting is then applied.

Similarly to the analysis in Section 4.2, we can prove the next two claims:

Connections due to termination condition 2: The number of pairs satisfied out of $\{(s_i, t_i) : p_i \in \mathcal{P}_L\}$ is $\Omega(1/\log n) \cdot |\mathcal{P}_L|^\beta$.

Connections due to greedy iterations: A fraction of $\Omega(1/(\Delta|P|^\beta))$ of the pairs in $\{(s_i, t_i) : p_i \notin \mathcal{P}_L\}$ are already satisfied by the partial orientation \mathcal{G}_L . This follows by observing that the number of paths that are eliminated from $\mathcal{P}_{\ell-1}$ in iteration ℓ is at most $\Delta|\mathcal{P}_{\ell-1}|^\beta \leq \Delta|P|^\beta$.

Consequently, the number of satisfied pairs upon termination is:

$$\begin{aligned} & \Omega\left(\frac{1}{\Delta|P|^\beta}\right) (|P| - |\mathcal{P}_L|) + \Omega\left(\frac{1}{\log n}\right) |\mathcal{P}_L|^\beta \\ &= \Omega\left(\frac{1}{\Delta|P|^\beta}\right) (|P| - |\mathcal{P}_L|) + \Omega\left(\frac{1}{|\mathcal{P}_L|^{1-\beta} \log n}\right) |\mathcal{P}_L| \\ &= \Omega\left(\frac{1}{\max\{\Delta|P|^\beta, |P|^{1-\beta} \log n\}}\right) |P|. \end{aligned}$$

By choosing $\beta = \frac{1}{2}(1 + \log_{|P|}(\frac{\log n}{\Delta}))$, we obtain an approximation ratio of $\Omega(1/\sqrt{\Delta|P| \log n})$.

In this section we used the usual definition of path lengths: the length of a path is the number of its edges. The above analyses work in a similar way if we measure the length of a path by the number of its undirected edges or even by the number of undirected components the path visits. This yields the same asymptotic bounds with respect to the size of the input, but highlights the increasing performance of the algorithm for structured inputs where these path length measures are small.

5 Conclusions

In this paper we presented approximation algorithms for the MAXIMUM-MIXED-GRAPH-ORIENTATION problem, which has recently arisen in the study of biological networks. We first showed that tree-like instances admit orientations (that can be computed in polynomial time) satisfying a poly-logarithmic fraction of the input pairs. Then we extended these algorithms to develop the first approximation algorithm for the problem whose ratio depends only on the size of the input instance, where no structural properties are assumed. The algorithm has a sub-linear approximation ratio, which can be improved when the input pairs are connected by short paths. The known upper and lower bounds for the approximation ratio of MAXIMUM-MIXED-GRAPH-ORIENTATION are far from being tight. Closing this gap, both in the undirected and mixed cases, remains an open problem.

Acknowledgments

M.E. was supported by a research grant from the Dr. Alexander und Rita Besser-Stiftung. C.R.D. would like to thank Gerry Schwartz, Heather Reisman, and the University of Waterloo-Haifa International Experience Program for funding his

visit to the University of Haifa, during which part of this work was done. R.S. was supported by a research grant from the Israel Science Foundation (grant no. 385/06).

References

1. E. M. Arkin and R. Hassin. A note on orientations of mixed graphs. *Discrete Applied Mathematics*, 116(3):271–278, 2002.
2. V. Bafna, P. Berman, and T. Fujito. A 2-approximation algorithm for the undirected feedback vertex set problem. *SIAM Journal on Discrete Mathematics*, 12(3):289–297, 1999.
3. J. Bang-Jensen and G. Gutin. *Digraphs: Theory, Algorithms and Applications*. Springer, 2nd edition, 2008.
4. A.-L. Barabási and Z. N. Oltvai. Network biology: understanding the cell’s functional organization. *Nature Reviews Genetics*, 5(2):101–113, February 2004.
5. F. Boesch and R. Tindell. Robbins’s theorem for mixed multigraphs. *The American Mathematical Monthly*, 87(9):716–719, 1980.
6. J. Chen, F. V. Fomin, Y. Liu, S. Lu, and Y. Villanger. Improved algorithms for feedback vertex set problems. *Journal of Computer and System Sciences*, 74(7):1188–1198, 2008.
7. F. R. K. Chung, M. R. Garey, and R. E. Tarjan. Strongly connected orientations of mixed multigraphs. *Networks*, 15(4):477–484, 1985.
8. R. Cohen, S. Havlin, and D. ben-Avraham. Structural properties of scale-free networks. In *Handbook of Graphs and Networks: From the Genome to the Internet*. Wiley-VCH, 2002.
9. B. Dorn, F. Hüffner, D. Krüger, R. Niedermeier, and J. Uhlmann. Exploiting bounded signal flow for graph orientation based on cause–effect pairs. In *Proceedings of the 1st International ICST Conference on Theory and Practice of Algorithms in (Computer) Systems (TAPAS 2011)*, volume 6595 of *LNCS*. Springer, 2011. To appear.
10. S. Fields. High-throughput two-hybrid analysis. The promise and the peril. *The FEBS Journal*, 272(21):5391–5399, 2005.
11. J. Flum and M. Grohe. *Parameterized Complexity Theory*. Springer, 2006.
12. G. N. Frederickson and D. B. Johnson. Generating and searching sets induced by networks. In *Proceedings 7th International Colloquium on Automata, Languages and Programming (ICALP 1980)*, volume 85 of *LNCS*, pages 221–233. Springer, 1980.
13. I. Gamzu, D. Segev, and R. Sharan. Improved orientations of physical networks. In *Proceedings of the 10th International Workshop on Algorithms in Bioinformatics (WABI 2010)*, volume 6293 of *LNCS*, pages 215–225. Springer, 2010.
14. A. Gavin, M. Börsche, R. Krause, P. Grandi, M. Marzioch, A. Bauer, J. Schultz, J. M. Rick, A. Michon, C. Cruciat, M. Remor, C. Höfert, M. Schelder, M. Brajenovic, H. Ruffner, A. Merino, K. Klein, M. Hudak, D. Dickson, T. Rudi, V. Gnau, A. Bauch, S. Bastuck, B. Huhse, C. Leutwein, M. Heurtier, R. R. Copley, A. Edelmann, E. Querfurth, V. Rybin, G. Drewes, M. Raida, T. Bouwmeester, P. Bork, B. Seraphin, B. Kuster, G. Neubauer, and G. Superti-Furga. Functional organization of the yeast proteome by systematic analysis of protein complexes. *Nature*, 415(6868):141–147, Jan. 2002.

15. S. L. Hakimi, E. F. Schmeichel, and N. E. Young. Orienting graphs to optimize reachability. *Information Processing Letters*, 63(5):229–235, 1997.
16. J. Håstad. Some optimal inapproximability results. *Journal of the ACM*, 48(4):798–859, 2001.
17. R. M. Karp. Reducibility among combinatorial problems. In *Complexity of Computer Computations*, pages 85–103. Plenum Press, 1972.
18. A. Medvedovsky, V. Bafna, U. Zwick, and R. Sharan. An algorithm for orienting graphs based on cause-effect pairs and its applications to orienting protein networks. In *Proceedings of the 8th International Workshop on Algorithms in Bioinformatics (WABI 2008)*, volume 5251 of *LNCS*, pages 222–232. Springer, 2008.
19. M. E. J. Newman. The structure and function of complex networks. *SIAM Review*, 45(2):167–256, 2003.
20. H. E. Robbins. A theorem on graphs, with an application to a problem of traffic control. *The American Mathematical Monthly*, 46(5):281–283, 1939.
21. D. Silverbush, M. Elberfeld, and R. Sharan. Optimally orienting physical networks. In *Proceedings of the 15th Annual International Conference on Research in Computational Molecular Biology (RECOMB 2011)*, volume 6577 of *LNCS*, pages 424–436. Springer, 2011.
22. R. E. Tarjan. A note on finding the bridges of a graph. *Information Processing Letters*, 2(6):160–161, 1974.
23. C. Yeang, T. Ideker, and T. Jaakkola. Physical network models. *Journal of Computational Biology*, 11(2-3):243–262, 2004.