

Non-approximability of Weighted Multiple Sequence Alignment

Bodo Manthey¹

*Universität zu Lübeck
Institut für Theoretische Informatik
Wallstraße 40, 23560 Lübeck, Germany*

Abstract

We consider a weighted generalization of multiple sequence alignment with sum-of-pair score. Multiple sequence alignment without weights is known to be \mathcal{NP} -complete and can be approximated within a constant factor, but it is unknown whether it has a polynomial time approximation scheme. Weighted multiple sequence alignment can be approximated within a factor of $O(\log^2 n)$ where n is the number of sequences.

We prove that weighted multiple sequence alignment is MAX \mathcal{SNP} -hard and give a numerical lower bound on its approximability, namely $\frac{324}{323} - \epsilon$. This lower bound is obtained already for the simple binary weighted case where the weights are restricted to 0 and 1. Furthermore, we show that weighted multiple sequence alignment and its restriction to binary weights can be approximated to the same degree.

Key words: Multiple sequence alignment; Non-approximability; Computational biology; SP-score

1 Introduction

Multiple sequence alignment (MSA) is an important problem in computational biology (see e.g. Karp [13]). The alignment of a group of protein or nucleotide sequences yields information about the relationships between these sequences and it is also used to detect similarities (so called “homologous regions”)

Email address: manthey@tcs.mu-luebeck.de (Bodo Manthey).

URL: <http://www.tcs.mu-luebeck.de/pages/manthey/> (Bodo Manthey).

¹ Birth name: Bodo Siebert. Supported by DFG research grant Re 672/3

between them. This information is applied in constructing evolutionary trees and finding coherences between the function and structure of proteins and their sequences. For a general survey of this topic, see for instance Gusfield [9].

Many objective functions have been suggested to measure the quality of a multiple sequence alignment. One of the most widely used is the so called sum-of-pair score (SP-score, see Carrillo et al. [7]).

The problem of finding a multiple sequence alignment with minimum SP-score is \mathcal{NP} -hard (see Wang and Jiang [16] and Bonizzoni and Della Vedova [6]). For the case that the scoring function does not have to be a metric, Just has shown that MSA with SP-score is MAX \mathcal{SNP} -hard [12]. Akutsu et al. have investigated the multiple sequence alignment problem under several scoring functions, namely $\#LOG\#$ -score and IC -score [1]. They have shown that a variant of the multiple sequence alignment problem called local multiple alignment is MAX \mathcal{SNP} -hard under these scoring schemes.

However, if the scoring function fulfils the triangle inequality, no lower bound for the approximability of this problem is known so far. The complexity of MSA over an alphabet of fixed size with metric SP-scoring functions is of main interest. According to Jiang et al. the approximability of MSA with metric SP-score is an important open problem in computational biology [11].

To represent existing knowledge about the relationships of the sequences considered, a weighted variant of MSA was introduced by Wu et al. [17]. Each pair of sequences is assigned a nonnegative value reflecting their degree of relationship. This means that a pair that is assumed to be closely related will be assigned a high weight while a less related pair will be assigned a smaller weight. This generalization of MSA is called weighted multiple sequence alignment, or WMSA for short.

In this paper we also examine a restricted version of WMSA called binary weighted multiple sequence alignment (BMSA), where the weights are restricted to 0 and 1. The binary weights can be used to represent an arbitrary graph over which multiple sequence alignments can be determined. We prove that BMSA is equivalent to WMSA with respect to their approximability. Thus, an approximation algorithm for BMSA directly yields an approximation algorithm for the general case with the same performance ratio. Moreover, we prove the MAX \mathcal{SNP} -hardness and a numerical lower bound for the approximability of BMSA. These results are obtained even if the sequences are of fixed length and the alphabet is of fixed size. Thus, the difficulty of computing an optimal alignment is caused by the number of sequences, not by their length.

In the next section we give a formal definition of the problems considered. The reduction from WMSA to BMSA is presented in Section 3. In Section 4

we prove a lower bound for the approximability of a problem called Max-E2-neg-Lin2. This result will be used in Section 5 to prove a lower bound for the approximability of BMSA.

2 Definitions and Notations

Let Σ be an alphabet and $\Sigma' := \Sigma \cup \{-\}$, where “-” denotes a *gap symbol*. $S[k]$ denotes the k -th symbol of a sequence S . Let $\mathcal{S} = \{S_1, \dots, S_n\}$ be a family (a multiset) of sequences over Σ . An *alignment* of \mathcal{S} is a family $\mathcal{A} = \{\tilde{S}_1, \dots, \tilde{S}_n\}$ of sequences over Σ' such that all \tilde{S}_i have equal length and \tilde{S}_i is obtained from S_i by inserting gaps. The following is an example of an alignment of the three sequences ALIGNMENT, ALGORITHM, and APPROXIMATION.

```

ALIGN--MENT---
AL-GORI---THM-
APPROXIMA-TION

```

A function $d : \Sigma'^2 \rightarrow \mathbb{N}$ will be called *scoring function* if it is a metric, i.e. for any $x, y, z \in \Sigma'$ we have $d(x, y) = 0$ if and only if $x = y$, $d(x, y) = d(y, x)$, and $d(x, z) \leq d(x, y) + d(y, z)$. We define the distance of two sequences \tilde{S}_i and \tilde{S}_j of length ℓ as

$$D(\tilde{S}_i, \tilde{S}_j) := \sum_{k=1}^{\ell} d(\tilde{S}_i[k], \tilde{S}_j[k]).$$

Carrillo and Lipman [7] introduced a scoring scheme for alignments called *sum-of-pair score* (SP-score). The SP-score of an alignment $\mathcal{A} = \{\tilde{S}_1, \dots, \tilde{S}_n\}$ is defined by

$$D(\mathcal{A}) := \sum_{1 \leq i < j \leq n} D(\tilde{S}_i, \tilde{S}_j).$$

Multiple sequence alignment (MSA) is the problem of finding an alignment with minimum SP-score.

Wu et al. [17] generalized MSA to weighted sum-of-pair score. The weights are given by $W := (w_{S_i, S_j})_{S_i, S_j \in \mathcal{S}}$, a symmetric matrix of nonnegative integers. The weighted SP-score of an alignment \mathcal{A} is

$$D_W(\mathcal{A}) := \sum_{1 \leq i < j \leq n} w_{S_i, S_j} \cdot D(\tilde{S}_i, \tilde{S}_j).$$

This generalization is called *weighted multiple sequence alignment* (WMSA). The aim is to find an alignment with minimum weighted SP-score.

An instance of WMSA is a 4-tuple $(\Sigma, \mathcal{S}, d, W)$. We consider the case of a fixed alphabet Σ and a fixed scoring function d . Thus, a problem instance of WMSA is given by a pair (\mathcal{S}, W) . It is easy to see that any lower bound for this case also holds if we allow arbitrary scoring functions and alphabets.

In this paper we also consider a special case of WMSA called *binary weighted multiple sequence alignment* (BMSA), where the weights are restricted to 0 and 1.

For the case that the scoring function d does not have to fulfil $d(x, x) = 0$ for any $x \in \Sigma'$, Wang and Jiang [16] showed that MSA with SP-score is \mathcal{NP} -complete. This result was extended by Bonizzoni and Della Vedova [6]. They showed that MSA with metric SP-score is \mathcal{NP} -complete even if $|\Sigma| = 2$. Gusfield [8] presented an algorithm which achieves a performance ratio of $2 - \frac{2}{n}$ where n is the number of sequences. This result was improved by Bafna et al. [4]. For an arbitrary fixed constant r and $n \geq r$, their algorithm computes an alignment whose score is at most a factor $2 - \frac{r}{n}$ greater than the score of an optimal alignment. It is unknown whether MSA admits a polynomial time approximation scheme (PTAS, see e.g. Ausiello et al. [3]). WMSA with arbitrary weights can be approximated within a factor of $O(\log^2 n)$ (see Wu et al. [17]). Using a technique of Bartal [5] one can obtain a randomized $O(\log n \cdot \log \log n)$ approximation.

Papadimitriou and Yannakakis [14] introduced a complexity class of optimization problems called $\text{MAX } \mathcal{SNP}$. They showed that there exist problems that are $\text{MAX } \mathcal{SNP}$ -complete with respect to L-reductions. In the following, $\text{opt}(I)$ denotes the optimal score of an instance I of an optimization problem. For example, $\text{opt}(\mathcal{S})$ denotes the score of an optimal (weighted) alignment of \mathcal{S} .

Definition 1 *Let Π and Π' be two optimization problems. Then Π L-reduces to Π' if there exist polynomial time computable functions f_1, f_2 and constants $\gamma_1, \gamma_2 > 0$ such that for every instance I of Π the following properties hold:*

(1) *Function f_1 produces an instance $I' = f_1(I)$ of Π' such that*

$$\text{opt}(I') \leq \gamma_1 \cdot \text{opt}(I).$$

(2) *Given a solution S' for I' with cost $c'(S')$, function f_2 produces a solution $S = f_2(I, S')$ for I with cost $c(S)$ such that*

$$|c(S) - \text{opt}(I)| \leq \gamma_2 \cdot |c'(S') - \text{opt}(I')|.$$

No $\text{MAX } \mathcal{SNP}$ -hard problem has a polynomial time approximation scheme, unless $\mathcal{NP} = \mathcal{P}$ (see Arora et al. [2]).

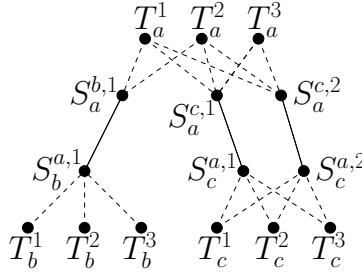


Fig. 1. Connections between the sequences obtained from $\mathcal{S} = \{S_a, S_b, S_c\}$, $w_{S_a, S_b} = 1$, $w_{S_a, S_c} = 2$, $w_{S_b, S_c} = 0$, and $K = 3$.

3 Reduction from WMSA to BMSA

Let $\mathcal{S} = \{S_1, \dots, S_n\}$ be a family of sequences over an alphabet Σ and $W = (w_{S_i, S_j})_{S_i, S_j \in \mathcal{S}}$ be a weight matrix. Let ℓ_{\max} be the maximal length of the sequences in \mathcal{S} and d_{\max} be the maximum value of the scoring function d . We assume that the weights and the scoring function are unary coded. In practice the weights are very small and the scoring function is fixed. Therefore, we may restrict to the case of unary coded weights.

We construct a family of sequences \mathcal{S}' as an instance of BMSA as follows. Let $K := 2 \cdot d_{\max} \cdot \ell_{\max}$. For a sequence $S_i \in \mathcal{S}$ generate K copies $T_i^k \in \mathcal{S}'$ ($1 \leq k \leq K$) of this sequence. Furthermore, for each $1 \leq j \leq n$ construct w_{S_i, S_j} copies $S_i^{j, \mu} \in \mathcal{S}'$ ($1 \leq \mu \leq w_{S_i, S_j}$) of S_i . The weight matrix $W' = (w'_{I, J})_{I, J \in \mathcal{S}'}$ is given by

$$w'_{I, J} := \begin{cases} 1 & \text{if } I \equiv S_i^{j, \mu} \text{ and } J \equiv S_j^{i, \mu}, \\ 1 & \text{if } I \equiv S_i^{j, \mu} \text{ and } J \equiv T_i^k \text{ or vice versa,} \\ 0 & \text{otherwise,} \end{cases}$$

where $A \equiv B$ means that A and B are not only equal but denote the same sequence. An example is shown in Figure 1. There is an edge between two sequences I and J if and only if $w'_{I, J} = 1$.

Since the weights and the scoring function are unary coded, the input size N of the instance of WMSA fulfils the bound

$$N \in \Omega\left(n + \ell_{\max} + \sum_{i, j=1}^n w_{S_i, S_j}\right).$$

On the other hand, the input size N' of the constructed instance of BMSA satisfies

$$N' \in O\left(\underbrace{n \cdot K \cdot \ell_{\max}}_{T_j^k} + \underbrace{\ell_{\max} \cdot \sum_{i, j=1}^n w_{S_i, S_j}}_{S_j^{i, \mu}} + \underbrace{\left(n \cdot K + \sum_{i, j=1}^n w_{S_i, S_j}\right)^2}_{W'}\right).$$

Note that N' is polynomially bounded by N .

Lemma 2 *If \mathcal{S} has an alignment \mathcal{A} with weighted score $D_W(\mathcal{A})$ then \mathcal{S}' has an alignment \mathcal{A}' with weighted score $D_{W'}(\mathcal{A}') = D_W(\mathcal{A})$.*

PROOF. Let $\mathcal{A} = \{\tilde{S}_1, \dots, \tilde{S}_n\}$ be an alignment of \mathcal{S} with weighted score D . We obtain an alignment $\mathcal{A}' = \{\tilde{A} \mid A \in \mathcal{S}'\}$ of \mathcal{S}' by setting $\tilde{T}_i^k = \tilde{S}_i$ and $\tilde{S}_i^{j,\mu} = \tilde{S}_i$ for all i, k, j, μ . The score of \mathcal{A}' with respect to the weight matrix W' is

$$\begin{aligned} D_{W'}(\mathcal{A}') &= \sum_{i,j=1}^n \sum_{\mu=1}^{w_{S_i, S_j}} \sum_{k=1}^K \underbrace{D(\tilde{S}_i^{j,\mu}, \tilde{T}_i^k)}_{=0} \\ &+ \sum_{1 \leq i < j \leq n} \sum_{\mu=1}^{w_{S_i, S_j}} \underbrace{D(\tilde{S}_i^{j,\mu}, \tilde{S}_j^{i,\mu})}_{=D(\tilde{S}_i, \tilde{S}_j)} = D_W(\mathcal{A}). \end{aligned}$$

□

Lemma 3 *Given an alignment \mathcal{A}' of \mathcal{S}' with weighted score $D_{W'}(\mathcal{A}')$ we can construct an alignment \mathcal{A} of \mathcal{S} with less or equal score in polynomial time.*

PROOF. Let $\mathcal{A}' = \{\tilde{A} \mid A \in \mathcal{S}'\}$ be an arbitrary alignment of \mathcal{S}' with score $D_{W'}(\mathcal{A}')$. We call the copies of a sequence $S_i \in \mathcal{S}$ *consistent* if there exists a sequence B_i with $\tilde{T}_i^k = B_i$ and $\tilde{S}_i^{j,\mu} = B_i$ for all k, j, μ . The sequence B_i is called *block*.

We consider the case that for some i_0 the copies of S_{i_0} are not consistent and distinguish two cases. On the one hand, if not all $\tilde{T}_{i_0}^k$ are equal, let

$$D_k := \sum_{j=1}^n \sum_{\mu=1}^{w_{S_{i_0}, S_j}} D(\tilde{T}_{i_0}^k, \tilde{S}_{i_0}^{j,\mu})$$

be the score of $\tilde{T}_{i_0}^k$ with the sequences $\tilde{S}_{i_0}^{j,\mu}$ ($1 \leq j \leq n, 1 \leq \mu \leq w_{S_{i_0}, S_j}$). We choose k_0 such that D_{k_0} is minimal among all D_k and set $\tilde{T}_{i_0}^k = \tilde{T}_{i_0}^{k_0}$ for all $k \neq k_0$. This way we obtain a new alignment with less or equal score.

On the other hand, we consider the case that there exists a B_{i_0} such that $\tilde{T}_{i_0}^k = B_{i_0}$ for all k . Then there exists a sequence $\tilde{S}_{i_0}^{j_0, \mu_0} \neq B_{i_0}$. This sequence yields at least score K with the sequences $\tilde{T}_{i_0}^k$ ($1 \leq k \leq K$), because it yields a score of at least 1 with every $\tilde{T}_{i_0}^k$. Set $\tilde{S}_{i_0}^{j_0, \mu_0} = B_{i_0}$. Then $\tilde{S}_{i_0}^{j_0, \mu_0}$ yields score 0

with every $\tilde{T}_{i_0}^k$ and at most score K with $\tilde{S}_{j_0}^{i_0, \mu_0}$. Thus, the new alignment has less or equal score.

By these modifications we iteratively obtain a new alignment of \mathcal{S}' such that for any $i \in \{1, \dots, n\}$ the copies of S_i are consistent with block B_i . The blocks of \mathcal{S}' induce an alignment $\mathcal{A} = \{B_1, \dots, B_n\}$ of \mathcal{S} with score

$$D_W(\mathcal{A}) = \sum_{1 \leq i < j \leq n} w_{S_i, S_j} \cdot D(B_i, B_j) \leq D_{W'}(\mathcal{A}').$$

□

With these results we have shown that a λ -approximation for BMSA can be used as a λ -approximation for WMSA. Thus, the following theorem is proved.

Theorem 4 *If BMSA can be approximated within a constant factor λ in polynomial time, then WMSA can also be approximated within λ in polynomial time.* □

4 An approximability gap for Max-E2-neg-Lin2

We consider the multiplicative group $\{1, -1\}$. Let $\mathcal{G} = \{G_1, \dots, G_t\}$ be a multiset of linear equations over the variables $U = \{x_1, \dots, x_r\}$, where G_i is the equation $x_{\alpha_{i,1}} \cdot \dots \cdot x_{\alpha_{i,k}} = a_i$, $a_i \in \{1, -1\}$ is a constant, and $\alpha_{i,q} \in \{1, \dots, r\}$. Max-Ek-Lin2 is the optimization problem of finding the maximum number of simultaneously satisfiable equations. A restriction of Max-Ek-Lin2 is Max-Ek-neg-Lin2, where $a_i = -1$ for all $1 \leq i \leq t$.

Max-E2-neg-Lin2 is exactly the problem Max-Cut (see e.g. Ausiello et al. [3]) where the equations correspond to the edges, the variables correspond to the vertices, and multiple edges are allowed. Therefore, Max-E2-neg-Lin2 is MAX \mathcal{NP} -complete [14]. We use Max-E2-neg-Lin2 here due to the simpler notation.

An instance of Max-Ek-Lin2 or Max-Ek-neg-Lin2 consisting of t equations will be called η -satisfiable if and only if $\eta \cdot t$ is the maximum number of simultaneously satisfiable equations. Håstad [10] proved that it is \mathcal{NP} -hard to distinguish $(1 - \epsilon)$ -satisfiable and $(\frac{1}{2} + \epsilon)$ -satisfiable instances of Max-E3-Lin2 for any $\epsilon > 0$.

Instead of using the known lower bound for the approximability of Max-Cut (see Håstad [10] and Trevisan et al. [15]) we construct a reduction from Max-E3-Lin2 to Max-E2-neg-Lin2 to show that it is \mathcal{NP} -hard to distinguish

$(\frac{18}{22} - \epsilon)$ - and $(\frac{17}{22} + \epsilon)$ -satisfiable instances of Max-E2-neg-Lin2 for any $\epsilon > 0$; the gadget used by Trevisan et al. [15] does not yield such a gap directly. This result will be used in Section 5 to prove the lower bound for the approximability of BMSA.

We will now reduce Max-E3-Lin2 to Max-E2-neg-Lin2. Therefore, let $\mathcal{G} = \{G_1, \dots, G_t\}$ be a multiset of equations over variables U and G_i be the equation $x_{\alpha_{i,1}} \cdot x_{\alpha_{i,2}} \cdot x_{\alpha_{i,3}} = a_i$.

We construct an instance \mathcal{G}' of Max-E2-neg-Lin2 with $22 \cdot t$ equations and $4 \cdot t + 2 \cdot r + 2$ variables. The reduction is similar to the reduction from Max-E3-Lin2 to Max-E2-Lin2 presented by Håstad [10]. The set of variables U' is given by

$$U' = \{x_j^+, x_j^- \mid 1 \leq j \leq r\} \cup \{z^+, z^-\} \cup \{p_{i,1}, p_{i,2}, p_{i,3}, p_{i,z} \mid 1 \leq i \leq t\}.$$

Note that if an assignment satisfies an equation of an instance of Max-E2-neg-Lin2, then the negated assignment also satisfies the equation. So without loss of generality we assume that we always have $z^+ = 1$.

We interpret $x_j^+ = x_j$. We call an assignment *consistent for x_j* if $x_j^+ \neq x_j^-$ and therefore $x_j^+ = x_j = (-x_j^-)$. An assignment that is consistent for every x_j and also fulfils $z^+ \neq z^-$ is called *consistent*.

For an equation G_i we construct the eighteen equations

$$\begin{aligned} x_{\alpha_{i,q}}^+ \cdot p_{i,q'} &= -1 \quad (\text{for } q, q' = 1, 2, 3 \text{ and } q \neq q'), \\ x_{\alpha_{i,q}}^+ \cdot p_{i,z} &= -1 \quad (\text{for } q = 1, 2, 3), \\ x_{\alpha_{i,q}}^- \cdot p_{i,q} &= -1 \quad (\text{for } q = 1, 2, 3), \\ x_{\alpha_{i,q}}^+ \cdot x_{\alpha_{i,q}}^- &= -1 \quad (q = 1, 2, 3), \\ z^+ \cdot z^- &= -1 \quad (\text{three times}). \end{aligned}$$

We add either the four equations

$$z^+ \cdot p_{i,q} = -1 \quad (q = 1, 2, 3) \quad \text{and} \quad z^- \cdot p_{i,z} = -1$$

if $a_i = 1$ or the four equations

$$z^- \cdot p_{i,q} = -1 \quad (q = 1, 2, 3) \quad \text{and} \quad z^+ \cdot p_{i,z} = -1$$

if $a_i = -1$. Note that \mathcal{G}' contains $3 \cdot t$ times the equation $z^+ \cdot z^- = -1$. Let n_j be the number of occurrences of the variable x_j in \mathcal{G} . Then \mathcal{G}' contains n_j times the equation $x_j^+ \cdot x_j^- = -1$.

For every equation $G_i \in \mathcal{G}$ we have constructed 22 equations for \mathcal{G}' . These 22 equations are called the *representation of G_i* .

Lemma 5 *Let an arbitrary assignment for U be given. Assign $z^- = -1$ and $x_j^+ = x_j$, $x_j^- = (-x_j)$ for $j = 1, \dots, r$. Then for any $i \in \{1, \dots, t\}$ there exists an assignment for $p_{i,1}$, $p_{i,2}$, $p_{i,3}$, and $p_{i,z}$ such that 18 equations of the representation of G_i are satisfied if G_i is satisfied by the given assignment and 16 equations of the representation are satisfied if G_i is not satisfied.*

It is not possible to satisfy more than 18 equations of the representation if G_i is satisfied by the assignment and to satisfy more than 16 equations if G_i is not satisfied by the assignment.

PROOF. The lemma can be proved by testing all possible assignments for the variables occurring in a representation of an equation G_i for the cases $a_i = 1$ and $a_i = -1$. \square

If an assignment for U satisfies g of the t equations of \mathcal{G} , then the corresponding consistent assignment for U' satisfies $16 \cdot t + 2 \cdot g$ equations of \mathcal{G}' . This assignment can be found efficiently by adjusting the assignment for $p_{i,1}$, $p_{i,2}$, $p_{i,3}$, and $p_{i,z}$. On the other hand, a consistent assignment for U' that satisfies $16 \cdot t + 2 \cdot g$ equations of \mathcal{G}' yields an assignment for U that satisfies at least g equations of \mathcal{G} .

Lemma 6 *Given an arbitrary assignment for U' that satisfies $16 \cdot t + 2 \cdot g$ equations of \mathcal{G}' , a consistent assignment that satisfies at least this amount of equations of \mathcal{G} can be computed in polynomial time.*

PROOF. First assume that $z^+ = z^-$ in the given assignment. Then the $3 \cdot t$ equations $z^+ \cdot z^- = -1$ are not satisfied by the assignment. Let $z^- = (-z^+)$. Then these $3 \cdot t$ equations will be satisfied. On the other hand, z^- occurs in only $3 \cdot t$ other equations. Thus, at most $3 \cdot t$ equations are no longer satisfied. Altogether the number of satisfied equations is not decreased by this modification.

If there exists a j with $x_j^+ = x_j^-$, then there are n_j equations $x_j^+ \cdot x_j^- = -1$ that are not satisfied by the assignment. Let $x_j^- = (-x_j^+)$. Then the n_j equations $x_j^+ \cdot x_j^- = -1$ are satisfied by the modified assignment. On the other hand x_j^- occurs in only n_j other equations. Thus, at most n_j equations are no longer satisfied. The number of satisfied equations is thus not decreased by this modification.

This way we iteratively obtain a consistent assignment. The modifications can be computed in polynomial time. \square

Now we can prove the following theorem used in Section 5.

Theorem 7 *For any $\epsilon > 0$ it is \mathcal{NP} -hard to distinguish $(\frac{18}{22} - \epsilon)$ - and $(\frac{17}{22} + \epsilon)$ -satisfiable instances of Max-E2-neg-Lin2.*

PROOF. An instance of Max-E3-Lin2 is η -satisfiable if and only if the corresponding instance of Max-E2-neg-Lin2 is $(\frac{16+2\eta}{22})$ -satisfiable. According to Håstad [10] it is \mathcal{NP} -hard to distinguish $(1 - \xi)$ - and $(\frac{1}{2} + \xi)$ -satisfiable instances of Max-E3-Lin2 for any $\xi > 0$. Thus, it is \mathcal{NP} -hard to distinguish $(\frac{16+2\cdot(1-\xi)}{22})$ - and $(\frac{16+2\cdot(\frac{1}{2}+\xi)}{22})$ -satisfiable instances of Max-E2-neg-Lin2. Choosing $\xi = 11 \cdot \epsilon$ completes the proof. \square

Since Max-Cut and Max-E2-neg-Lin2 are exactly the same problem, we obtain the same approximability gap for Max-Cut.

Corollary 8 *For any $\epsilon > 0$ it is \mathcal{NP} -hard to decide whether the maximum cut of an undirected graph $G = (V, E)$ (where multiple edges are allowed) consists of at most $(\frac{17}{22} + \epsilon) \cdot |E|$ or at least $(\frac{18}{22} - \epsilon) \cdot |E|$ edges. \square*

5 The non-approximability of BMSA

In this section we reduce Max-E2-neg-Lin2 to BMSA. Let $\mathcal{G} = \{G_1, \dots, G_t\}$ be an instance of Max-E2-neg-Lin2 over a set of variables $U = \{x_1, \dots, x_r\}$, where G_i is $x_{\alpha_{i,1}} \cdot x_{\alpha_{i,2}} = -1$, $\alpha_{i,q} \in \{1, \dots, r\}$. We construct a family of sequences

$$\mathcal{S} = \{Z\} \cup \{X_j \mid j = 1, \dots, r\} \cup \{Y_{i,1}, Y_{i,2} \mid i = 1, \dots, t\}$$

over the alphabet $\Sigma = \{\bullet, \circ, \times\}$. Let

$$Z := \circ \circ \circ \circ \circ \circ \circ \circ$$

be a sequence of length 8. Z will be used as a control sequence. For $j \in \{1, \dots, r\}$ let

$$X_j := \bullet \circ \circ \circ \circ \circ \circ \circ \bullet$$

	-	•	○	×
-	0	1	2	5
•	1	0	1	4
○	2	1	0	3
×	5	4	3	0

Fig. 2. The scoring function.

be a sequence of length 9 that represents the variable $x_j \in U$. For each $i \in \{1, \dots, t\}$ create two sequences

$$Y_{i,1} := \bullet \circ \circ \times \circ \times \circ \circ \bullet \text{ and}$$

$$Y_{i,2} := \bullet \circ \circ \circ \times \circ \circ \circ \bullet ,$$

each of length 9. $Y_{i,q}$ represents the variable $x_{\alpha_{i,q}}$ in G_i .

The scoring function is shown in Figure 2. Note that it is a metric.

The weight matrix $W = (w_{I,J})_{I,J \in \mathcal{S}}$ is given by

$$w_{I,J} := \begin{cases} 1 & \text{if } I \equiv Y_{i,q} \text{ and } J \equiv Y_{i,q'} , \\ 1 & \text{if } I \equiv Z \text{ and } J \equiv Y_{i,q} \text{ or vice versa } , \\ 1 & \text{if } I \equiv Y_{i,q} \text{ and } J \equiv X_{\alpha_{i,q}} \text{ or vice versa } , \\ 0 & \text{otherwise .} \end{cases}$$

An example how the sequences are connected is shown in Figure 3(a).

The set $\mathcal{S}_i = \{Y_{i,1}, Y_{i,2}, X_{\alpha_{i,1}}, X_{\alpha_{i,2}}\}$ will be called the *representation of G_i* . Note that in general a sequence X_j occurs in more than one representation.

Let $\mathcal{A} = \{\tilde{S} \mid S \in \mathcal{S}\}$ be an alignment of \mathcal{S} . Then $D_i(\mathcal{A})$ denotes the score of the equation G_i ,

$$D_i(\mathcal{A}) = D(\tilde{Y}_{i,1}, \tilde{Y}_{i,2}) + D(\tilde{Y}_{i,1}, \tilde{X}_{\alpha_{i,1}}) + D(\tilde{Y}_{i,2}, \tilde{X}_{\alpha_{i,2}}) \\ + D(\tilde{Y}_{i,1}, \tilde{Z}) + D(\tilde{Y}_{i,2}, \tilde{Z}) .$$

By the construction of the weight matrix, we have $D_W(\mathcal{A}) = \sum_{i=1}^t D_i(\mathcal{A})$.

Definition 9 An alignment $\mathcal{A} = \{\tilde{S} \mid S \in \mathcal{S}\}$ of \mathcal{S} will be called *variable-consistent with respect to an assignment for U* if, after eliminating all columns

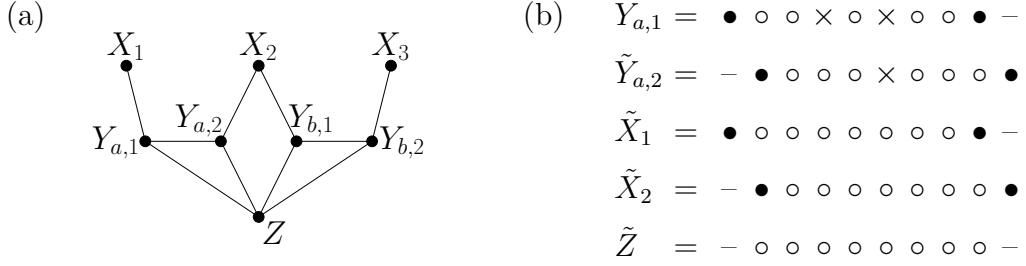


Fig. 3. (a) Connection between the sequences representing the equations G_a , which is $x_1 \cdot x_2 = -1$, and G_b , which is $x_2 \cdot x_3 = -1$. (b) The variable consistent alignment corresponding to G_a with $x_1 = -1$ and $x_2 = 1$.

consisting solely of gaps (which do not affect the score), the following holds for all j, i , and q :

$$(1) \tilde{Z} = -Z-$$

$$(2) \tilde{X}_j = \begin{cases} X_j- & \text{if } x_j = -1 \\ -X_j & \text{if } x_j = 1 \end{cases}$$

$$(3) \tilde{Y}_{i,q} = \begin{cases} Y_{i,q}- & \text{if } x_{\alpha_{i,q}} = -1 \\ -Y_{i,q} & \text{if } x_{\alpha_{i,q}} = 1 \end{cases}$$

The following lemma follows immediately from this definition.

Lemma 10 *An alignment is variable-consistent if and only if for all $i \in \{1, \dots, t\}$ and $q \in \{1, 2\}$ the following properties hold:*

- A. *Either $Y_{i,q}[1]$ or $Y_{i,q}[9]$ matches a gap in Z . No other character of Z or $Y_{i,q}$ matches a gap in the other sequence.*
- B. *No character in either of the two sequences $Y_{i,q}, X_{\alpha_{i,q}}$ matches a gap in the other sequence.*

□

These properties are referred to as Property A and B. An example of a variable-consistent alignment is shown in Figure 3(b).

Note the functional region of a pair $Y_{i,1}, Y_{i,2}$ given by the triples $\times \circ \times$ and $\circ \times \circ$. If $Y_{i,1}$ and $Y_{i,2}$ represent the same value, the functional region yields a weighted score of 9. Otherwise, it yields a weighted score of 3. If an alignment \mathcal{A} is variable-consistent, we have $D_i(\mathcal{A}) = 29$ if G_i is satisfied by the represented assignment and $D_i(\mathcal{A}) = 31$ otherwise.

The next two lemmas have similar proofs. Thus, we only give a proof of the first.

Lemma 11 *Alignments of the pairs $\{Y_{i,1}, Z\}$ and $\{Y_{i,2}, Z\}$ yield scores of 8 and 5, respectively, if they fulfil Property A. Violating Property A yields scores of at least 10 and 7, respectively.*

PROOF. An alignment of $\{Y_{i,1}, Z\}$ that fulfils Property A yields score 8.

Let us consider an alignment of $\{Y_{i,1}, Z\}$ that does not fulfil Property A. Then at least one of the characters $Y_{i,1}[2], \dots, Y_{i,1}[8], Z[1], \dots, Z[8]$ matches a gap in the other sequence.

We distinguish two cases. If there is an “ \times ” in $Y_{i,1}$ matching a gap in Z , then the alignment yields a score of 5 for this “ \times ” plus 3 for the other “ \times ” plus 1 for each “ \bullet ”. So altogether it yields a score of at least 10.

On the other hand consider the case that no “ \times ” in $Y_{i,1}$ matches a gap in Z . Then there is a “ \circ ” in $Y_{i,1}$ or Z matching a gap in the other sequence. So the alignment yields a score of 3 for each “ \times ” plus 1 for each “ \bullet ” plus 2 for the “ \circ ” matching a gap. So the alignment again yields a score of at least 10.

The statement about $Y_{i,2}$ and Z can be proved in a similar fashion. \square

Lemma 12 *Alignments of the pairs $\{Y_{i,1}, X_{\alpha_{i,1}}\}$ and $\{Y_{i,2}, X_{\alpha_{i,2}}\}$ yield scores of 6 and 3, respectively, if they fulfil Property B. Violating Property B yields scores of at least 8 and 5, respectively.* \square

With the fact that an optimal alignment of a pair $\{Y_{i,1}, Y_{i,2}\}$ has score 7 we can prove the following.

Lemma 13 *Given an arbitrary alignment with score $31 \cdot t - 2 \cdot g$ we can construct a variable-consistent alignment with less or equal score in polynomial time.*

PROOF. Let \mathcal{A} be an arbitrary alignment with $D_W(\mathcal{A}) = 31 \cdot t - 2 \cdot g$.

Let I be the set of all i such that $Y_{i,1}$ and $Y_{i,2}$ fulfil Properties A and B. This implies an assignment for the variables $U_I = \{x_j \in U \mid \exists i \in I : X_j \in \mathcal{S}_i\}$. Let $\bar{I} = \{1, \dots, t\} \setminus I$. Because in every set \mathcal{S}_i for $i \in \bar{I}$ there exists a sequence $Y_{i,q}$ that violates Property A or B, we have $D_i(\mathcal{A}) \geq 31$ for each $i \in \bar{I}$ due to Lemmas 11 and 12.

For $i \in \bar{I}$, if $x_{\alpha_{i,q}} \in U_I$ ($q \in \{1, 2\}$) we realign $Y_{i,q}$ with respect to $x_{\alpha_{i,q}}$. Then we assign an arbitrary value to the variables in $U \setminus U_I$ and realign the corresponding $Y_{i,q}$ and X_j .

By these modifications we obtain an alignment \mathcal{A}' . Then $D_i(\mathcal{A}') = D_i(\mathcal{A})$ for $i \in I$ and $D_i(\mathcal{A}') \leq 31 \leq D_i(\mathcal{A})$ otherwise. Thus, $D_W(\mathcal{A}') \leq D_W(\mathcal{A})$. \mathcal{A}' is variable-consistent due to its construction and can be computed in polynomial time. \square

The alignment obtained yields an assignment that satisfies at least g equations of \mathcal{G} .

Theorem 14 *BMSA is MAX SNP-hard.*

PROOF. We reduce Max-E2-neg-Lin2 to BMSA. The function f_1 is given by the construction of \mathcal{S} from a family \mathcal{G} of t equations. One can see that $\text{opt}(\mathcal{S}) \leq 31 \cdot t$.

An equation of \mathcal{G} will be satisfied by 2 of the 4 possible assignments of its variables. Therefore, for every multiset \mathcal{G} of t equations an assignment exists that satisfies at least $\frac{1}{2} \cdot t$ equations. This yields the inequality $\text{opt}(\mathcal{G}) \geq \frac{1}{2} \cdot t$. Then for $\gamma_1 = 62$ we have $\text{opt}(\mathcal{S}) \leq \gamma_1 \cdot \text{opt}(\mathcal{G})$.

Given an alignment of \mathcal{S} with score $31 \cdot t - 2 \cdot g'$ for some g' we can find an assignment satisfying $g \geq g'$ equations of \mathcal{G} due to Lemma 13. Let $\gamma_2 = \frac{1}{2}$, then we have

$$|g - \text{opt}(\mathcal{G})| \leq \gamma_2 \cdot |(31 \cdot t - 2 \cdot g') - \text{opt}(\mathcal{S})| .$$

\square

Theorem 15 *BMSA has no polynomial time approximation algorithm with approximation ratio $\frac{324}{323} - \epsilon$ for any $\epsilon > 0$, unless $\mathcal{NP} = \mathcal{P}$.*

PROOF. An instance of Max-E2-neg-Lin2 consisting of t equations is η -satisfiable if and only if the corresponding instance of BMSA has an alignment with score $(31 - 2 \cdot \eta) \cdot t$.

The optimal alignment of a BMSA instance corresponding to a $(\frac{18}{22} - \xi)$ -satisfiable instance of Max-E2-neg-Lin2 has score

$$\left(31 - 2 \cdot \left(\frac{18}{22} - \xi\right)\right) \cdot t = \frac{323 + 22 \cdot \xi}{11} \cdot t .$$

Using the $(\frac{324}{323} - \epsilon)$ -approximation algorithm for BMSA we are able to find an alignment with score at most

$$\left(\frac{324}{323} - \epsilon\right) \cdot \frac{323 + 22 \cdot \xi}{11} \cdot t =: K_1 .$$

The optimal alignment of a BMSA instance corresponding to a $(\frac{17}{22} + \xi)$ -satisfiable instance of Max-E2-neg-Lin2 has score

$$\left(31 - 2 \cdot \left(\frac{17}{22} + \xi\right)\right) \cdot t =: K_2.$$

We have $K_1 < K_2$ if and only if $\xi < \frac{1}{22} \cdot \frac{323^2 \cdot \epsilon}{647 - 323 \cdot \epsilon}$. Choose ξ with $0 < \xi < \frac{1}{22} \cdot \frac{323^2 \cdot \epsilon}{647 - 323 \cdot \epsilon}$. Then the $(\frac{324}{323} - \epsilon)$ -approximation for BMSA can be used to distinguish $(\frac{18}{22} - \xi)$ - and $(\frac{17}{22} + \xi)$ -satisfiable instances of Max-E2-neg-Lin2. This would imply $\mathcal{NP} = \mathcal{P}$ due to Theorem 7. \square

Since WMSA is a generalization of BMSA we obtain the following corollaries.

Corollary 16 *WMSA is MAX \mathcal{SNP} -hard.* \square

Corollary 17 *WMSA has no polynomial time approximation algorithm with approximation ratio $\frac{324}{323} - \epsilon$ for any $\epsilon > 0$, unless $\mathcal{NP} = \mathcal{P}$.* \square

6 Conclusions

We have shown MAX \mathcal{SNP} -hardness and proved a numerical lower bound for the approximability of weighted multiple sequence alignment (WMSA). These results hold even if we restrict the problem to binary weights (BMSA). Furthermore, BMSA and WMSA are equivalent with respect to their approximability. But the distance to the best known upper bound is huge. An obvious goal is to reduce this gap.

Finally, we would like to know how well the unweighted version of the multiple sequence alignment problem with metric SP-score can be approximated.

Acknowledgements

I am grateful to Martin Böhme, Andreas Jakoby, and Rüdiger Reischuk for valuable discussions and comments on this paper.

References

- [1] T. Akutsu, H. Arimura, S. Shimozone, On approximation algorithms for local multiple alignment, in: Proc. 4th Ann. Int. Conf. on Comput. Molec. Biology (RECOMB), ACM, 2000, pp. 1–7.
- [2] S. Arora, C. Lund, R. Motwani, M. Sudan, M. Szegedy, Proof verification and the hardness of approximation problems, *J. ACM* 45 (3) (1998) 501–555.
- [3] G. Ausiello, P. Crescenzi, G. Gambosi, V. Kann, A. Marchetti-Spaccamela, M. Protasi, *Complexity and Approximation: Combinatorial Optimization Problems and Their Approximability Properties*, Springer, 1999.
- [4] V. Bafna, E. L. Lawler, P. A. Pevzner, Approximation algorithms for multiple sequence alignment, *Theoret. Comput. Sci.* 182 (1–2) (1997) 233–244.
- [5] Y. Bartal, On approximating arbitrary metrics by tree metrics, in: Proc. 30th Ann. Symp. on Theory of Comput. (STOC), ACM, 1998, pp. 161–168.
- [6] P. Bonizzoni, G. Della Vedova, The complexity of multiple sequence alignment with SP-score that is a metric, *Theoret. Comput. Sci.* 259 (1) (2001) 63–79.
- [7] H. Carrillo, D. J. Lipman, The multiple sequence alignment problem in biology, *SIAM J. Appl. Math.* 48 (1988) 1073–1082.
- [8] D. M. Gusfield, Efficient methods for multiple sequence alignment with guaranteed error bounds, *Bull. Math. Biology* 55 (1) (1993) 141–154.
- [9] D. M. Gusfield, *Algorithms on Strings, Trees, and Sequences: Computer Science and Computational Biology*, Cambridge University Press, 1997.
- [10] J. Håstad, Some optimal inapproximability results, *J. ACM* 48 (4) (2001) 798–859.
- [11] T. Jiang, P. E. Kearney, M. Li, Some open problems in computational molecular biology, *J. Algorithms* 34 (1) (2000) 194–201.
- [12] W. Just, Computational complexity of multiple sequence alignment with SP-score, *J. Comput. Biology* 8 (6) (2001) 615–623.
- [13] R. M. Karp, Mapping the genome: Some combinatorial problems arising in molecular biology, in: Proc. 25th Ann. Symp. on Theory of Comput. (STOC), ACM, 1993, pp. 278–285.
- [14] C. H. Papadimitriou, M. Yannakakis, Optimization, approximation, and complexity classes, *J. Comput. System Sci.* 43 (3) (1991) 425–440.
- [15] L. Trevisan, G. B. Sorkin, M. Sudan, D. P. Williamson, Gadgets, approximation, and linear programming, *SIAM J. Comput.* 29(6) (2000) 2074–2097.
- [16] L. Wang, T. Jiang, On the complexity of multiple sequence alignment, *J. Comput. Biology* 1 (4) (1994) 337–348.

- [17] B. Y. Wu, G. Lancia, V. Bafna, K.-M. Chao, R. Ravi, C. Y. Tang, A polynomial-time approximation scheme for minimum routing cost spanning trees, *SIAM J. Comput.* 29 (3) (1999) 761–778.