

# 双対分解によるマルチプルアラインメント

穴水 拓郎<sup>1,a)</sup> 榎原 康文<sup>1,b)</sup> 佐藤 健吾<sup>1,c)</sup>

概要：バイオインフォマティクスにおいて塩基配列やアミノ酸配列を解析する手法として、配列アラインメントが広く用いられている。ペアワイズアラインメントは動的計画法で効率的に解けることが知られている。一方、マルチプルアラインメントは、その計算量の問題から累進法と呼ばれる近似手法で解くことが多いが、最適解が保証されない。本研究では、マルチプルアラインメントを整数計画問題として定式化し、双対分解によって総当たりのペアワイズアラインメントに相当する部分問題に分割することによって最適解を効率よく求める手法を提案する。

TAKURO ANAMIZU<sup>1,a)</sup> YASUBUMI SAKAKIBARA<sup>1,b)</sup> KENGO SATO<sup>1,c)</sup>

## 1. はじめに

配列アラインメントとは、二本以上の配列を入力として類似する塩基どうしが縦に揃うように並べる操作である。特に、配列本数が二本のときをペアワイズアラインメント、三本以上のときをマルチプルアラインメントと呼ぶ。ペアワイズアラインメントは動的計画法で効率的に解ける一方、マルチプルアラインメントは累進法と呼ばれる近似手法で解くことが多い。しかし累進法では最適解が保証されず、繰り返し最適化法による精度改善は、計算時間とトレードオフになってしまう。本研究では、マルチプルアラインメントを整数計画問題として定式化し、双対分解 [4] によって総当たりのペアワイズアラインメントに相当する部分問題に分割することによって最適解を効率よく求める手法 DMSA (Dual decomposition for Multiple Sequence Alignment) を提案する。

## 2. 手法

配列  $a$  の長さを  $|a|$  と書くこととする。二本の配列  $a, b$  が与えられたとき、 $\mathcal{A}(a, b)$  を  $a$  と  $b$  がとりうる全ての配列アラインメントの集合とする。アラインメント  $z \in \mathcal{A}(a, b)$  は  $|a| \times |b|$  次元の二値行列  $z = (z_{ij})$  で表わす。ここで  $z_{ij} = 1$  は、塩基  $a_i$  と  $b_j$  がアラインされていることを表す。 $n$  本の配

列  $s_1, \dots, s_n$  のマルチプルアラインメント  $\theta \in \mathcal{A}(s_1, \dots, s_n)$  が総当たりのペアワイズアラインメントからなるとき、 $\theta = (z^{(1,2)}, \dots, z^{(n-1,n)})$  と書くこととする。

マルチプルアラインメントの利益関数を、総当たりのペアワイズアラインメントの利益関数の総和として定義する：

$$G(\theta, \hat{\theta}) = G(z^{(1,2)}, \hat{z}^{(1,2)}) + \dots + G(z^{(n-1,n)}, \hat{z}^{(n-1,n)}) \quad (1)$$

ペアワイズアラインメント  $z$  に関する利益関数  $G(z, \hat{z})$  は次のように定義される：

$$G(z, \hat{z}) = (1 - \sigma)TP(z, \hat{z}) + \sigma TN(z, \hat{z}) \quad (2)$$

ここで  $TP(z, \hat{z}) = \sum_{i,j} I(z_{ij} = 1)I(\hat{z}_{ij} = 1)$  は true positive の数、 $TN(z, \hat{z}) = \sum_{i,j} I(z_{ij} = 0)I(\hat{z}_{ij} = 0)$  は true negative の数、 $\sigma$  は true positive と true negative のバランスを制御するパラメータである。期待精度最大化原理 [2] に基づき、アラインメント空間  $\mathcal{A}(s_1, \dots, s_n)$  上に与えられる確率分布  $P(\theta | s_1, \dots, s_n)$  の下で、利益関数の期待値を最大化するマルチプルアラインメント  $\hat{\theta}$  を求める：

$$\mathbb{E}[G(\theta, \hat{\theta})] = \sum_{\theta \in \mathcal{A}(s_1, \dots, s_n)} P(\theta | s_1, \dots, s_n) G(\theta, \hat{\theta}) \quad (3)$$

ここで、マルチプルアラインメントの確率分布を次のように積近似する：

$$P(\theta | s_1, \dots, s_n) \approx \prod_{1 \leq k < l \leq n} P(z^{(k,l)} | s_k, s_l) \quad (4)$$

その結果、期待利益関数は次のように近似することができる：

<sup>1</sup> 慶應義塾大学理工学部生命情報学科  
Department of Biosciences and Informatics, Keio University

<sup>a)</sup> anamizu@dna.bio.keio.ac.jp

<sup>b)</sup> yasu@bio.keio.ac.jp

<sup>c)</sup> satoken@bio.keio.ac.jp

$$\mathbb{E}[G(\theta, \hat{\theta})] \approx \sum_{1 \leq k < l \leq n} \sum_{i=1}^{|s_k|} \sum_{j=1}^{|s_l|} [p_{ij}^{(k,l)} - \sigma] z_{ij}^{(k,l)} + C \quad (5)$$

ここで

$$p_{ij}^{(a,b)} = \sum_{z \in \mathcal{A}(a,b)} P(z | a, b) I(z_{ij} = 1) \quad (6)$$

はアラインメント事後確率であり、 $C$  は  $\theta$  に依存しない定数である。

本手法の目的は、マルチプルアラインメントが満たすべき制約の下、期待利益関数を最大化するマルチプルアラインメントを計算することである。この最適化は次のような整数計画問題として定式化することができる：

maximize:

$$S(z; s) = \sum_{1 \leq k < l \leq n} \sum_{i=1}^{|s_k|} \sum_{j=1}^{|s_l|} [p_{ij}^{(k,l)} - \sigma] z_{ij}^{(k,l)} \quad (7)$$

subject to:

$$\sum_{j=1}^{|s_l|} z_{ij}^{(k,l)} \leq 1 \quad (1 \leq \forall k \leq n, 1 \leq \forall l \leq n, 1 \leq \forall i \leq |s_k|) \quad (8)$$

$$z_{iq}^{(k,l)} + z_{jp}^{(l,m)} \leq 1 \quad (9)$$

$$(1 \leq \forall k < \forall l \leq n, 1 \leq \forall i < \forall j \leq |s_k|, 1 \leq \forall p < \forall q \leq |s_l|)$$

$$z_{ij}^{(k,l)} - z_{jh}^{(l,m)} + z_{hi}^{(m,k)} \leq 1 \quad (10)$$

$$z_{ij}^{(k,l)} + z_{jh}^{(l,m)} - z_{hi}^{(m,k)} \leq 1 \quad (11)$$

$$-z_{ij}^{(k,l)} + z_{jh}^{(l,m)} + z_{hi}^{(m,k)} \leq 1 \quad (12)$$

$$(1 \leq \forall k < \forall l < \forall m \leq n, 1 \leq \forall i \leq |s_k|, 1 \leq \forall j \leq |s_l|, 1 \leq \forall h \leq |s_m|)$$

$$z_{i_1 i_2}^{(\rho_1, \rho_2)} + \dots + z_{i_{k-1} i_k}^{(\rho_{k-1}, \rho_k)} + z_{i_k i_1}^{(\rho_k, \rho_1)} \leq k - 1 \quad (13)$$

$$(\forall \rho \in P_k, 1 \leq \forall i_p \leq |s_{\rho_p}| \text{ for } 1 \leq \forall p \leq k)$$

ただし、 $P_k$  は  $\{1, \dots, n\}$  の大きさ  $k$  の巡回順列の集合である。制約 (8) は、二本の配列の各々の塩基は高々 1 つの塩基とのみアラインされることを意味する。制約 (9) は交差するアラインメントを許さないことを表す。制約 (10)-(12) は、consistency transformation [3] に相当する。制約 (13) は、アラインメント列の順番に矛盾がないことを表す。

制約 (10)-(12) および制約制約 (13) を目的関数に移すことによって、ラグランジュ双対関数を定義する。

$$\begin{aligned} L(\lambda, \mu, \nu, \xi) = & \max S(z; s) \quad (14) \\ & + \sum_{1 \leq k < l < m \leq n} \sum_{i=1}^{|s_k|} \sum_{j=1}^{|s_l|} \sum_{h=1}^{|s_m|} \lambda_{ijh}^{(k,l,m)} (1 - z_{ij}^{(k,l)} - z_{jh}^{(l,m)} + z_{hi}^{(m,k)}) \\ & + \sum_{1 \leq k < l < m \leq n} \sum_{i=1}^{|s_k|} \sum_{j=1}^{|s_l|} \sum_{h=1}^{|s_m|} \mu_{ijh}^{(k,l,m)} (1 - z_{ij}^{(k,l)} + z_{jh}^{(l,m)} - z_{hi}^{(m,k)}) \\ & + \sum_{1 \leq k < l < m \leq n} \sum_{i=1}^{|s_k|} \sum_{j=1}^{|s_l|} \sum_{h=1}^{|s_m|} \nu_{ijh}^{(k,l,m)} (1 + z_{ij}^{(k,l)} - z_{jh}^{(l,m)} - z_{hi}^{(m,k)}) \\ & + \sum_{3 \leq k \leq n} \sum_{\rho \in P_k} \sum_{i_1=1}^{|s_{\rho_1}|} \dots \sum_{i_k=1}^{|s_{\rho_k}|} \xi_{i_1 \dots i_k}^{(\rho)} (k - 1 - (z_{i_1 i_2}^{(\rho_1, \rho_2)} + \dots + z_{i_{k-1} i_k}^{(\rho_{k-1}, \rho_k)} + z_{i_k i_1}^{(\rho_k, \rho_1)})) \end{aligned}$$

ここで、 $\lambda, \mu, \nu, \xi$  はラグランジュ未定乗数である。上の式

は次のように書きかえることができる。

$$L(\lambda, \mu, \nu, \xi) = \sum_{1 \leq k < l \leq n} \sum_{i=1}^{|s_k|} \sum_{j=1}^{|s_l|} [p_{ij}^{(k,l)} - \sigma + \phi_{ij}^{(k,l)}] z_{ij}^{(k,l)} \quad (15)$$

ただし、

$$\begin{aligned} \phi_{ij}^{(k,l)} = & \sum_{m>l} \sum_{h=1}^{|s_m|} (-\lambda_{ijh}^{(k,l,m)} - \mu_{ijh}^{(k,l,m)} + \nu_{ijh}^{(k,l,m)}) \\ & + \sum_{k<m<l} \sum_{h=1}^{|s_m|} (\lambda_{ijh}^{(k,m,l)} - \mu_{ijh}^{(k,m,l)} - \nu_{ijh}^{(k,m,l)}) \\ & + \sum_{m<k} \sum_{h=1}^{|s_m|} (-\lambda_{ijh}^{(m,k,l)} + \mu_{ijh}^{(m,k,l)} - \nu_{ijh}^{(m,k,l)}) \\ & - \sum_{3 \leq m \leq n} \sum_{\rho \in P_m \text{ s.t.}} \sum_{i_1=1}^{|s_{\rho_1}|} \dots \sum_{i_{l-1}=1}^{|s_{\rho_{l-1}}|} \sum_{i_{l+1}=1}^{|s_{\rho_{l+1}}|} \dots \sum_{i_k=1}^{|s_{\rho_k}|} \xi_{i_1 \dots i_k}^{(\rho)} \end{aligned}$$

劣勾配法によって式 (15) を  $\lambda, \mu, \nu, \xi$  に関して最小化することによって元々の問題 (7) の解を得ることができる。そのアルゴリズムは以下のように動作する：(1) 入力配列の総当たりのペアワイズアラインメントに分解して、それぞれ独立に最適化する。各ペアワイズアラインメントは Needleman-Wunsch アルゴリズムによって高速に解くことができる。(2) マルチプルアラインメントが満たすべき制約 (10)-(13) と矛盾するアラインメントカラムに対応するスコアにペナルティを与える。以上の手順を繰り返すことによって、マルチプルアラインメントが満たすべき制約の下、目的関数すなわち期待精度を最大化するマルチプルアラインメントを得る。

### 3. 結果

前節で述べたアルゴリズムに基づき DMSA を実装した。アラインメント確率行列 (6) を計算するために、ProbCons [1] を用いた。実験結果は発表時に示す。

#### 参考文献

- [1] Do, C. B., Mahabhashyam, M. S., Brudno, M. and Batzoglou, S.: ProbCons: Probabilistic consistency-based multiple sequence alignment, *Genome Res.*, Vol. 15, No. 2, pp. 330–340 (2005).
- [2] Hamada, M. and Asai, K.: A classification of bioinformatics algorithms from the viewpoint of maximizing expected accuracy (MEA), *J. Comput. Biol.*, Vol. 19, No. 5, pp. 532–549 (2012).
- [3] Notredame, C., Higgins, D. G. and Heringa, J.: T-Coffee: A novel method for fast and accurate multiple sequence alignment, *J. Mol. Biol.*, Vol. 302, No. 1, pp. 205–217 (2000).
- [4] Sato, K., Kato, Y., Akutsu, T., Asai, K. and Sakakibara, Y.: DAFS: simultaneous aligning and folding of RNA sequences via dual decomposition, *Bioinformatics*, Vol. 28, No. 24, pp. 3218–3224 (2012).