

基于多策略人工蜂群的多序列比对算法

匡芳君^{1,2†}, 张思扬¹, 刘传才²

(1. 温州商学院 信息工程学院, 浙江 温州 325035; 2. 南京理工大学 计算机科学与工程学院, 南京 210094)

摘要: 多序列比对是生物信息学中最重要和最具挑战性的任务之一. 基于多序列比对是NP完全组合优化问题, 引入Tent混沌初始化种群策略、不同蜂种的邻域搜索策略和锦标赛选择策略等, 提出一种基于多策略人工蜂群的多序列比对算法. 该算法应用Tent混沌初始化种群策略以使初始个体多样化并获取较好初始解; 针对不同蜂种的特性设计不同的邻域搜索策略以平衡算法的全局探索和局部开发能力. 同时引入序列比对的蜜源编码方法以适应多序列比对的离散性. 实验结果表明, 所提出算法的鲁棒性较强, 能获得较好的比对性能和生物特性.

关键词: 人工蜂群算法; 多策略; Tent混沌初始化; 邻域搜索; 多序列比对

中图分类号: TP18

文献标志码: A

Multiple sequence alignment algorithm based on multi-strategy artificial bee colony

KUANG Fang-jun^{1,2†}, ZHANG Si-yang¹, LIU Chuan-cai²

(1. School of Information Engineering, Wenzhou Business College, Wenzhou 325035, China; 2. School of Computer Science and Engineering, Nanjing University of Science and Technology, Nanjing 210094, China)

Abstract: Multiple sequence alignment(MSA), known as NP-complete combinatorial optimization problem, is one of the most important and challenging tasks in bioinformatics. A multi-strategy artificial bee colony(MS-ABC) algorithm is proposed for MSA, which is composed of multiple strategies, such as the Tent chaotic initialization population strategy, different neighborhood search strategies and tournament selection strategy. In the MSA-ABC algorithm, the Tent chaotic initialization population strategy is presented to diversify the initial individuals and to obtain good initial solutions. Then, the different neighborhood search strategies for different bee species are designed to balance the global exploration and the local exploitation. Moreover, the food source encoding method is used to adapt discreteness of MSA. The experimental results demonstrate that the proposed algorithm is more robust and can obtain better alignment quality and biological characteristics.

Keywords: artificial bee colony; multi-strategy; Tent chaotic initialization; neighborhood search; multiple sequence alignment

0 引言

多序列比对(Multiple sequence alignment, MSA)是生物信息学研究的热点之一, 通过多序列比对可以挖掘生物序列中的结构、进化和功能等信息, MSA也是一个NP完全组合优化难题^[1]. 自然启发式算法作为一类适用于求解NP难题的优化算法, 具有精度高、对度量标准不敏感等优势, 吸引着大量研究人员, 其研究成果颇丰. 如Gupta等^[2]提出了一种基于遗传算法的多序列比对方法(MSA-GA); Gao^[3]提出了一种基于惯性权重粒子群优化算法(MS-PSO); Tsvetanov

等^[4]提出了一种基于改进蚁群算法的多序列比对方法(MS-ACO); Öztürk等^[5]提出了一种新的人工蜂群算法的多序列比对方法(ABC-Aligner); Zhu等^[6]提出了一种基于分解的多目标进化算法(MOMSA)以解决多序列比对问题, 初始种群通过使用空位插入操作生成, 进化操作利用遗传操作的交叉和变异算子实现; Liu等^[7]提出了一种基于隐马尔可夫模型和后验概率分配函数的多序列比对算法(MSAProbs); Rani等^[8]提出了基于遗传和人工蜂群的混合算法(GA-ABC)和多目标细菌觅食优化算法(MO-BFO), 并应

收稿日期: 2017-06-20; 修回日期: 2017-08-22.

基金项目: 国家自然科学基金项目(61373063, 61233011, 61402227).

责任编委: 巩敦卫.

作者简介: 匡芳君(1976-), 女, 教授, 博士, 从事群智能与多目标优化、模式识别、生物信息学及其应用等研究; 刘传才(1963-), 男, 教授, 博士生导师, 从事智能计算、计算机视觉、图像处理及应用等研究.

†通讯作者. E-mail: kfjztb@126.com

用于多序列比对,但在工作中主要集中于MO-BFO算法;Sun等^[9]提出了一种基于量子粒子群优化和隐马尔可夫模型的多序列比对算法;Rubio-Largo等^[10]结合多目标和进化启发式算法,提出了一种基于混合蛙跳算法的多序列比对方法;Zambrano-Vega等^[11]进行了基于多目标启发式算法的多序列比对方法的比较研究,并对这些启发式方法进行了综述和比较。

人工蜂群算法(Artificial bee colony, ABC)是由Karaboga于2005年提出的群智能算法^[12],具有全局探索能力强、控制参数少、收敛速度快和鲁棒性强等优势,被广泛应用于求解复杂优化问题^[13],但当其接近全局最优解时,种群多样性减少,搜索速度变慢,甚至易陷入局部最优,因此也衍生了很多改进算法和应用^[14-16]。本文在研究多序列比对算法作为多目标优化问题的基础上,在人工蜂群算法中引入Tent混沌初始化种群策略、不同蜂种的邻域搜索策略和锦标赛选择策略等,提出了一种基于多策略人工蜂群的多序列比对算法。

1 多序列比对优化模型

多序列比对主要反映给定序列之间的进化关系。目的是使参与比对的序列尽可能相似(或距离最小),即相同残基的位点位于同一列。多序列比对问题描述如下^[6]:假设参与比对的序列组 $S = (S_1, S_2, \dots, S_i, S_N)$ 含有 N 条序列,每条序列长度分别为 l_1, l_2, \dots, l_N ,序列中字符由4种核苷酸或20种氨基酸组成。通过插入或删除空位找到一个多序列比对 S' ,使对齐的字符数最多。多序列比对 S' 的约束条件为:1) S' 中各条比对后的序列长度相等;2) 若删除 S' 中的第 i 行所有空位,则 $S'_i = S_i$;3) 使多序列比对 S' 的评分 $\text{Score}(S')$ 最大。

多序列比对常用目标函数主要有:WSP权重计分函数,隐Markov模型,COFFEE计分函数等^[6]。本文采用WSP函数,其评分函数 $\text{Score}(S')$ 表示为

$$\text{Score}(S') = \sum_{l=1}^L S'_l \tag{1}$$

$$S'_l = \sum_{i=1}^{N-1} \sum_{k=i+1}^N w_{ij} p(A_i, A_j) \tag{2}$$

$$p(A_i, A_j) = \begin{cases} 0, & A_i = A_j = '-' \\ -g_{\text{open}}, & (A_i = '-' \text{ and } A_j \neq '-') \text{ or } \\ & (A_i \neq '-' \text{ and } A_j = '-') \text{ and is opening;} \\ -g_{\text{extend}}, & (A_i = '-' \text{ and } A_j \neq '-') \text{ or } \\ & (A_i \neq '-' \text{ and } A_j = '-') \text{ and is extension;} \\ \text{Blosum}(A_i, A_j), & A_i \neq '-' \text{ and } A_j \neq '-' \end{cases} \tag{3}$$

其中: L 为多序列比对 S' 的长度; S'_l 为多序列比对 S' 第 l 列的代价; N 为序列比对组中序列的条数; w_{ij} 为序列 s_i 和 s_j 的权重,取决于序列间的距离,本文利用编辑距离(Levenshtein distance, LD)来计算 w_{ij} ,令 $w_{ij} = 1 - \frac{\text{LD}(s_i, s_j)}{\max(|s_i|, |s_j|)}$,其中 $\max(|s_i|, |s_j|)$ 为两条序列 s_i 和 s_j 长度的最大值; $p(A_i, A_j)$ 为序列中残基 A_i 和 A_j 的评分; $g_{\text{open}}, g_{\text{extend}}$ 分别为单位空位开放罚分 and 单位空位扩展罚分; $\text{Blosum}(A_i, A_j)$ 为根据残基片段替换矩阵(Blocks substitution matrix)计算残基匹配评分^[6]。使序列间匹配的残基对个数最大,而空位罚分最小是多序列比对追求的目标,故多序列比对多目标优化函数表示如下:

$$\begin{aligned} \max \text{imize } F(S') &= (f_1(S'), f_2(S'))^T, \\ f_1(S') &= \text{Score}(S'), \\ f_2(S') &= -g(S'), \\ S' &\in \Omega. \end{aligned} \tag{4}$$

其中: $g(S') = n_1 \times g_{\text{open}} + n_2 \times g_{\text{extend}}$ 为对齐序列组 S' 的空位罚分, n_1, n_2 分别为空位开放的数量和空位扩展的数量; Ω 为多序列比对空间。本文选用Blosum 62替换矩阵, $g_{\text{open}} = 6, g_{\text{extend}} = 0.85$ 。

例 1 假设参与比对的序列组为 S ,通过序列比对后得到的对齐序列比对为 S' ,结果如图 1 所示。

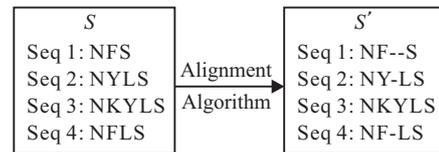


图 1 多序列比对实例

1) 计算序列组 S 中每一对序列间的权重 w_{ij} ,即 $w_{12} = 1 - \frac{\text{LD}(s_1, s_2)}{\max(|s_1|, |s_2|)} = 1 - \frac{2}{4} = 0.5$ 。同理, $w_{13} = 0.4, w_{14} = 0.75, w_{23} = 0.8, w_{24} = 0.75, w_{34} = 0.6$ 。

2) 按式(3)和(4)计算序列 S' 每一列代价值,即 $S'_1 = w_{12} \times p(S'_{11}, S'_{21}) + w_{13} \times p(S'_{11}, S'_{31}) + w_{14} \times p(S'_{11}, S'_{41}) + w_{23} \times p(S'_{21}, S'_{31}) + w_{24} \times p(S'_{21}, S'_{41}) + w_{34} \times p(S'_{31}, S'_{41}) = (0.5 \times \text{Blosum}(N, N)) + (0.4 \times \text{Blosum}(N, N)) + (0.75 \times \text{Blosum}(N, N)) + (0.8 \times \text{Blosum}(N, N)) + (0.75 \times \text{Blosum}(N, N)) + (0.6 \times \text{Blosum}(N, N)) = 22.8$ 。

同理, $S'_2 = 3.45, S'_3 = -10.8, S'_4 = 7.1975, S'_5 = 15.2$ 。

3) 按式(1)计算对齐序列 S' 的评分 $\text{Score}(S') = 37.8475$ 。

4) 计算对齐序列组 S' 的空位罚分 $g(S') = n_1 \times g_{open} + n_2 \times g_{extend} = 18.85$. 因此, 此对齐序列组 S' 的目标函数值为 $f_1 = 37.8475, f_2 = -18.85$.

2 多策略人工蜂群算法

为更好地均衡 ABC 算法的全局探索能力与局部开采能力之间的关系, 本文利用 Tent 混沌初始化种群, 多种邻域搜索和锦标赛选择等策略, 提出多策略人工蜂群算法 (Multi-strategy ABC, MS-ABC).

2.1 ABC 算法

ABC 算法^[11]的蜂群分为引领蜂 (Employed bee)、跟随蜂 (Oonlooked bee) 和侦察蜂 (Scout bee). 可行解的质量或适应度用蜜源的丰富程度表示, 问题的最优解用最大适应度的蜜源表示. 优化问题的可行解通过式 (5) 随机产生蜜源的位置来表示. 引领蜂与蜜源一一对应, 引领蜂的数量与跟随蜂数量相等, 且等于蜂群总数的一半.

$$X_i^j(0) = X_{min}^j + R \times (X_{max}^j - X_{min}^j). \quad (5)$$

其中: $X_i^j(0)$ 为初始化时第 i 个蜜源的第 j 维向量; $i = 1, 2, \dots, N, j = 1, 2, \dots, D, N$ 为蜜源数量, D 为问题空间的维数; X_{max}^j, X_{min}^j 分别为第 j 维向量的最大值和最小值; R 为 $[0,1]$ 间的随机数. 蜜源产生后就分派给引领蜂. 为了能根据旧蜜源位置 X_i^j 产生新蜜源位置 V_i^j , ABC 算法的基本搜索策略^[11]为

$$V_i^j(t) = X_i^j(t) + \Phi(X_i^j(t) - X_k^j(t)). \quad (6)$$

其中: $i = 1, 2, \dots, N, j = 1, 2, \dots, D, k = 1, 2, \dots, N$ 为随机选择的下标, 且 $k \neq i; \Phi$ 为 $[-1, 1]$ 之间的随机数.

如果一个蜜源位置循环蜜源搜索限定次数 Limit 后仍不能更新, 则该蜜源处的引领蜂变成侦察蜂, 并且该侦察蜂按式 (5) 搜索策略在解空间中产生新蜜源位置来替换原蜜源位置.

2.2 Tent 混沌初始化种群策略

为了使产生的初始个体尽可能均匀分布, 并且让种群保持多样性, 考虑到 Tent 映射比 Logistic 映射具有更好的收敛速度和遍历均匀性, 采用 Tent 混沌初始化种群^[15]产生蜜源位置, 既能保留算法初始化的随机性, 又能提高种群多样性, 从而择优产生蜜源位置. Tent 映射通过贝努利移位变换后可以表示为 $x_{t+1} = (2x_t) \bmod 1$, 计算机进行 Tent 映射运算时, 小数部分的二进制数进行无符号左移^[15].

2.3 锦标赛选择策略

ABC 算法的跟随蜂如果利用比例选择策略进行蜜源的选择, 则可能导致较差蜜源不能得到较好地更

新, 使算法种群多样性减少, 从而导致算法后期个体适应度在种群进化时趋于一致, 进而使算法不易跳出局部最优. 因此, 本文采用锦标赛选择策略^[15]选择蜜源. 该策略是基于局部竞争机制的选择过程, 只根据适应度值的相对值来选择蜜源, 从而避免超级个体对算法的影响. 锦标赛选择概率为

$$P_i(t) = \frac{c_i(t)}{\sum_{i=1}^N c_i(t)}, \quad (7)$$

其中 c_i 为每个个体的得分.

2.4 多邻域搜索策略

在求解复杂优化问题时, 如何平衡全局探索与局部开发能力是群智能算法提升性能的关键^[16]. ABC 算法前期全局探索能力强, 收敛速度快, 后期局部开发能力弱, 收敛速度慢. 同时考虑到多序列比对是一个特殊的离散问题, 式 (5)、(6) 不能直接用于多序列比对问题. 因此, 本文将根据蜂群的种类设计不同的邻域搜索策略, 以克服基本人工蜂群算法易于陷入局部最优解的缺陷.

2.4.1 引领蜂邻域搜索

在引领蜂邻域搜索时, 既考虑个体 E_i , 又考虑整个种群迄今为止搜索到的最优引领蜂 (蜜源) E_{best} , 对任一引领蜂 E_i 和当前最优引领蜂 E_{best} 通过单点交叉操作产生新的引领蜂 V_i , 即在蜜源 E_{best} 选取较好交叉位块, 以替代 E_i 中相应的序列位块, 若不足序列长度, 则在中间补足空位. 假设当前最优蜜源 E_{best} 为

S'_1 : GARFIE- LDTHE- - - - FAST- - - CAT
 S'_2 : GARFI-ELDTHEVERYFAST- - - CAT
 S'_3 : GARFI- ELDTHE- - - -LASTFAT CAT
 S'_4 : - - - - - THE - - - - FA- T - - - CAT

引领蜂 E_i 为

S'_1 : GA-RFIELD- -THE- - - FAST- CAT - -
 S'_2 : GA-RFIELD- -THEVERYFAST- CAT
 S'_3 : G- -ARFIELDT-HEL- -AST-FAT CAT
 S'_4 : - - - -THEFAT- - - - - CAT- - -

则对 E_i 和 E_{best} 通过单点交叉操作后产生新的引领蜂 V_i 为

S'_1 : GA-RFIELD- -THE- - - FAST- - - CAT
 S'_2 : GA-RFIELD- -THEVERYFAST- CAT
 S'_3 : G- -ARFIELDT-HEL- -AST- FATCAT
 S'_4 : - - - -THEFAT- - - - - - - - - CAT

2.4.2 跟随蜂邻域搜索

跟随蜂根据锦标赛选择操作选中引领蜂 E_i 作为蜜源时, 跟随蜂在蜜源邻域位置通过变异操作产生一个新的蜜源位置 V_i , 为了提高算法精度, 引入一个新

$$S_2 = Y - D - EILYQ - SKRF.$$

4) 混沌映射产生混沌序列.

利用 Tent 映射产生的混沌序列表示插入空位的位置, 故应对映射后的序列做取整操作, 并判断欲插入的空位位置是否已存在. 若已存在, 则在剩余位置中任选一位置插入, 否则直接插入到该位置.

3.2 算法流程

基于 MS-ABC 多序列比对算法步骤描述如下:

Step 1: 设置多序列比对算法的参数包括种群规模 SN、蜜源(引领蜂)的个数为 $M = SN/2$ 、最大迭代次数 G_{max} 、蜜源开采限制次数 Limit 等参数, 设置当前迭代次数 $iter = 0$, 设置同一蜜源的开采标志向量 $trial(i) = 0$.

Step 2: 根据“2.2 Tent 混沌初始化策略”和“3.1 编码设计”产生初始蜜源插入空位的位置, 初始比对的长度 L 在 l_{max} 与 $1.2l_{max}$ 之间, 从而形成初始化种群, 并保证其对应的序列比对结果中不存在只由空位“-”组成的列.

Step 3: 对于每个个体 S_i , 根据式(4)计算目标函数 $f_1(S_i)$ 和 $f_2(S_i)$. 对适应度排序, 取前 M 个个体作为蜜源位置, 令 $F_1 = (f_1(S_1), f_1(S_2), \dots, f_1(S_L))$, $F_2 = (f_2(S_1), f_2(S_2), \dots, f_2(S_N))$, 则可以得到 $F = (\max(F_1), \max(F_2))$ 为目标理想值, 并记录种群当前最优位置 S_{best} .

Step 4: 引领蜂 i 按 2.3.1 节引领蜂邻域搜索策略在蜜源附近搜索产生新的序列比对 V_i , 并按式(4)计算适应度值 $F(V_i)$.

Step 5: 如果 $F(V_i) > F(S_i)$, 则 $S_i = V_i$, 否则保持原始序列比对 S_i 不变.

Step 6: 根据式(7)锦标赛选择策略计算选择概率 $P_i(t)$.

Step 7: 跟随蜂根据 $P_i(t)$ 选择蜜源, 按 2.3.2 节跟随蜂邻域搜索策略在蜜源附近搜索产生新的蜜源 V_i , 并按式(4)计算其适应度值 $F(V_i)$.

Step 8: 如果 $F(V_i) > F(S_i)$, 则 $S_i = V_i$, $trial(i) = 0$; 否则保持原始序列 S_i 不变, $trial(i) = trial(i) + 1$.

Step 9: 若 $trial(i) \geq Limit$, 则引领蜂 i 放弃当前蜜源变为侦察蜂, 然后侦察蜂按照 2.3.3 节侦察蜂邻域搜索策略在解空间搜索产生新蜜源 V_i 替代原蜜源.

Step 10: 记录当前所有最优蜜源, 保存迄今为止的最优序列比对 S_{best} .

Step 11: 更新迭代次数 $iter = iter + 1$, 如果 $iter \geq G_{max}$, 则输出最优序列比对结果, 否则返回 Step 4.

MS-ABC 多序列比对算法的流程如图 2 所示.

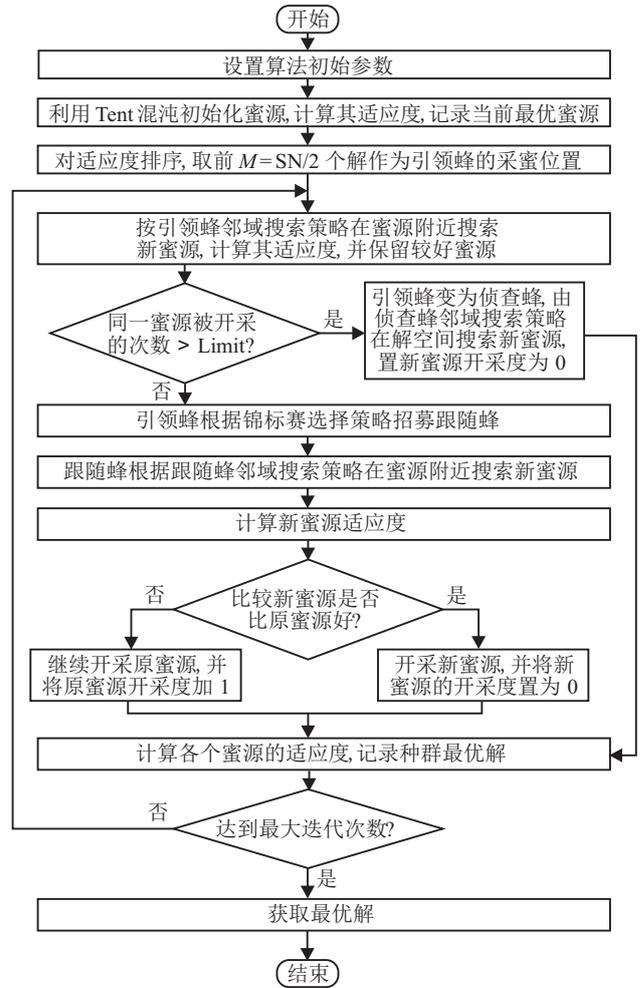


图 2 基于多策略人工蜂群的多序列比对算法流程

4 实验结果与分析

4.1 实验数据

为验证 MS-ABC 算法解决多序列比对的性能和有效性, 采用 BAliBASE 3.0 数据库^[17], 它是基于三维重叠的标准比对库, 包含 6 255 条序列, 218 个测试实例, 被划分为 6 个序列集合: BB11, BB12, BB20, BB30, BB40, BB50, 每组集合代表不同生物特性, 且包含的序列相似度也不同. BB11 集合中包含 38 个测试实例, 序列间相似度小于 20%; BB12 集合包含 44 个测试实例, 序列间相似度在 20% 到 40% 之间; BB20 集合包含 41 个测试实例, 序列间相似度大于 40%; BB30 集合包含 30 个测试实例, 序列间相似度大于 40%; BB40 集合包含 49 个测试实例, 序列间相似度大于 20%; BB50 集合包含 16 个测试实例, 序列间相似度大于 20%.

4.2 实验结果与分析

为了验证本文所提算法 MS-ABC 的性能, 分别与 MSA-GA^[2]、MS-PSO^[3]、MS-ACO^[4]、ABC-Aligner^[5]、MOMSA^[6]、MSAProbs^[7]、MO-BFO^[8]、MUSCLE^[16]、Clustal Ω ^[18] 等多序列比对算法进行比较实验. 采用

表 1 各算法解决多序列比对的 SP 平均值及相应标准差

算法	SP						
	BB11(38)	BB12(44)	BB20(41)	BB30(30)	BB40(49)	BB50(16)	Overall(218)
MS-ABC	0.728 4±0.108 1	0.951 6±0.011 2	0.934 5±0.003 8	0.887 3±0.013 3	0.930 8±0.007 3	0.899 6±0.011 5	0.890 7±0.010 8
MSA-GA	0.462 8±0.201 2	0.838 9±0.120 5	0.814 7±0.019 6	0.788 6±0.082 4	0.804 5±0.105 2	0.774 6±0.062 4	0.761 5±0.093 8
MSA-PSO	0.617 5±0.153 4	0.926 8±0.101 8	0.906 9±0.098 6	0.839 8±0.121 6	0.906 3±0.100 8	0.834 5±0.078 2	0.846 7±0.158 5
MSA-ACO	0.605 1±0.167 8	0.913 4±0.093 6	0.900 9±0.128 4	0.824 6±0.098 9	0.881 3±0.105 2	0.823 4±0.120 5	0.824 7±0.094 8
ABC-Aligner	0.695 0±0.094 5	0.942 5±0.065 6	0.924 7±0.030 2	0.853 2±0.107 9	0.922 6±0.086 4	0.888 6±0.015 8	0.875 9±0.100 8
MOMSA	0.567 8±0.124 7	0.866 5±0.098 4	0.892 5±0.008 3	0.759 5±0.183 0	0.824 1±0.009 8	0.765 8±0.104 3	0.754 6±0.158 4
MSAProbs	0.682 1±0.100 5	0.956 4±0.068 7	0.928 7±0.009 8	0.865 7±0.031 4	0.923 6±0.080 5	0.901 6±0.019 8	0.875 4±0.024 5
MO-BFO	0.635 4±0.105 8	0.934 8±0.089 6	0.902 4±0.127 8	0.843 5±0.086 7	0.894 5±0.100 2	0.846 7±0.084 6	0.834 9±0.088 7
MUSCLE	0.682 6	0.944 7	0.928 4	0.875 5	0.925 4	0.894 4	0.875 2
Clustal Ω	0.590 1	0.906 0	0.911 6	0.862 4	0.901 0	0.862 0	0.838 9

表 2 各算法解决多序列比对的 TC 平均值及相应标准差

算法	TC						
	BB11(38)	BB12(44)	BB20(41)	BB30(30)	BB40(49)	BB50(16)	Overall(218)
MS-ABC	0.591 2±0.125 4	0.876 9±0.030 2	0.515 7±0.070 5	0.640 5±0.098 6	0.658 4±0.018 9	0.622 5±0.034 7	0.650 6±0.086 5
MSA-GA	0.316 7±0.302 3	0.635 8±0.195 1	0.245 8±0.115 8	0.608 1±0.190 4	0.358 7±0.208 7	0.298 9±0.097 4	0.356 4±0.129 3
MSA-PSO	0.384 6±0.193 5	0.836 8±0.090 6	0.356 4±0.187 6	0.489 2±0.107 5	0.583 4±0.097 3	0.509 5±0.100 3	0.548 9±0.182 1
MSA-ACO	0.368 5±0.173 4	0.803 1±0.078 3	0.398 7±0.163 2	0.486 4±0.077 3	0.508 1±0.085 6	0.439 8±0.072 3	0.516 4±0.098 4
ABC-Aligner	0.486 4±0.112 3	0.870 2±0.092 7	0.446 9±0.127 6	0.568 7±0.128 5	0.625 8±0.112 4	0.558 4±0.099 8	0.584 6±0.124 0
MOMSA	0.389 5±0.140 5	0.814 3±0.112 9	0.363 8±0.119 2	0.384 1±0.098 2	0.534 6±0.009 8	0.514 6±0.187 0	0.561 8±0.101 4
MSAProbs	0.484 2±0.127 1	0.871 5±0.098 5	0.469 8±0.094 8	0.611 8±0.100 9	0.610 9±0.010 9	0.611 2±0.069 5	0.605 8±0.095 3
MO-BFO	0.432 1±0.208 2	0.814 8±0.100 9	0.401 3±0.131 1	0.501 3±0.210 6	0.545 9±0.094 5	0.486 7±0.185 4	0.578 4±0.152 8
MUSCLE	0.440 9	0.861 9	0.479 4	0.622 7	0.603 9	0.593 1	0.600 3
Clustal Ω	0.362 2	0.793 8	0.452 9	0.579 1	0.582 6	0.537 4	0.551 3

多序列比对评估指标: 残基对评分和 (Sum of pairs score, SP), TC 度量 (Total column, TC). SP 指的是算法测试得到的多序列比对与标准参考比对的残基对的比例, 而 TC 是指所有序列中正确对齐列的比例.

所有算法在 CPU 主频为 4 GHz, 显卡为 NVidia GeForce GTX TITAN X(12 G), 内存为 16 G, 操作系统为 Windows 10 Enterprise 64 位的 PC 机上运行, 在 Matlab R2014 b 和 JAVA 平台编码实现. 所有进化算法的种群规模都取 SN = 100, 其中: ABC 算法种群包含引领蜂 50 个, 跟随蜂 50 个; PSO 算法种群包含粒子 100 个; ACO 中种群包含 100 只蚂蚁; GA 中种群包含 100 个染色体. 所有算法最大迭代次数 $G_{max} = 5000$, 所有 ABC 算法的蜜源开采限制次数 Limit = 25. 所有进化算法在固定迭代次数下, 通过求解同一组序列的比对, 利用算法运行 30 次所得的平均 SP 值和 TC 值来考察算法的性能. 实验独立运行 30 次后, 获得测试结果中的 SP 平均值及其相应的标准差比较如表 1 所示, TC 平均值及其相应的标准差比较如表 2 所示, 较好的算法结果在表中用粗体表示.

表 1 和表 2 分别表示 10 种算法对多序列比对的 SP 平均值及其标准差和 TC 平均值及其标准差的求解. 综合表 1、表 2 分析可以看出, MS-ABC 算法除在 BB12 和 BB50 序列比对组中的性能比 MSAProbs 算

法的 SP 平均值性能稍低外, 其他情况都具有较明显优越性. 同时, 不管序列比对组的序列长度、相似度差别多大, MS-ABC 算法都能得到较好、标准差较小的序列比对, 且对所有序列组比对时也能获得较好比对效果, 这表明 MS-ABC 算法鲁棒性和稳定性较好, 而且对 BB12 和 BB50 序列比对结果比 Clustal Ω 和 MUSCLE 专业软件得到的结果也要好, 特别是对 218 个测试实例中的 6255 条序列比对时, MS-ABC 算法获得的性能也相对而言较好, 其主要原因是 MS-ABC 算法在邻域搜索过程中综合考虑了全局探索能力和局部开采能力的均衡. 对于序列之间的相似性小于 20% 的测试集 BB11, MS-ABC 算法得到 SP 平均值和 TC 平均值虽然比其他非专业比对软件的效果好, 但分别只有 0.728 4 和 0.591 2, 主要是序列间的相似度对序列比对结果的影响较大, 因此, 数据集 BB11 对于 MS-ABC 算法仍然是一个挑战.

表 3 表示的是 10 种比对算法对 BAliBASE 3.0 数据库中 218 个测试实例的运行时间的比较. 由表 3 可知, Clustal Ω 和 MUSCLE 序列比对算法的运行时间比其他启发式算法的运行时间要短, 分别是 9'16" 和 13'29", 主要是因为 Clustal Ω 和 MUSCLE 算法是用于多序列比对的高度专业化的软件, 它在比对序列之前, 根据序列间进化距离创建进化树; MOMSA 算法

表3 各算法对BALiBASE 3.0数据库中218个测试实例的运行时间

h:m:s

算法	MSA-ABC	MSA-GA	MSA-PSO	MSA-ACO	ABC-Aligner	MOMSA	MSAProbs	MO-BFO	MUSCLE	Clustal Ω
Runtime	01:59:23	01:47:34	02:58:52	02:46:25	03:06:42	23:37:56	03:28:21	02:48:17	00:13:29	00:09:16

最慢,主要是因为MOMSA算法没有考虑对参数进行优化,使用了默认参数,而其他启发式序列比对算法对参数进行了优化;MS-ABC算法的运行速度相对而言稍快,下一步工作将采用多CPU或GPU优化算法,以提高算法的运行速度。

5 结论

本文应用多策略人工蜂群算法优化多序列比对模型目标函数,提出了基于多策略人工蜂群的多序列比对算法.算法特点主要有:1)利用Tent混沌初始化策略提高算法种群的多样性;2)针对不同蜂种特性设计不同的邻域搜索策略以平衡算法的全局探索与局部开发能力,从而使算法有效跳出局部最优,并快速搜索最优解;3)引领蜂使用锦标赛选择策略招募跟随蜂,在一定程度上避免了超级个体对算法的影响;4)引入蜜源编码方法以适应多序列比对的离散性;5)利用多策略人工蜂群算法优化多序列比对多目标优化模型,有利于提供更多有意义的多序列比对供用户选择.实验结果表明,所提出算法不仅能有效地求解多序列比对问题,还能获得良好的比对结果.因此,该方法可为解决多序列比对问题提供新思路,也可为其他相关研究提供新的解决方案.但随着比对序列数量的海量增加,本文算法也会遇到运行速度慢、内存不足和CPU处理效能等瓶颈.另外,算法对于相似度小的数据集比对精度仍需提高.下一步重点考虑如何利用多CPU或GPU和Spark平台搭建并行计算集群来提高算法的并行优化和扩展性,并结合隐马尔可夫模型和群智能混合算法研究多序列比对模型中的多目标优化问题,以提高序列比对质量。

参考文献(References)

- [1] Wang Y X, Wang Z H. Introduction to bioinformatics: Algorithms and applications for high performance computing[M]. Beijing: Tsinghua University Press, 2011.
- [2] Gupta R, Pankaj A, Soni A K. MSA-GA: Multiple sequence alignment tool based on genetic approach[J]. Int J of Soft Computing and Software Engineering, 2013, 3(8): 1-11.
- [3] Gao Y X. A multiple sequence alignment algorithm based on inertia weights particle swarm optimization[J]. J of Bionanoscience, 2014, 8(5): 400-404.
- [4] Tsvetanov S, Ivanova D, Zografov B. Ant colony optimization applied for multiple sequence alignment[J]. Biomath Communications, 2015, 2(1): 800-806.
- [5] Öztürk C, Aslan S. A new artificial bee colony algorithm

- to solve the multiple sequence alignment problem[J]. Int J of Data Mining and Bioinformatics, 2016, 14(4): 332-353.
- [6] Zhu H Z, He Z S, Jia Y Y. A novel approach to multiple sequence alignment using multi-objective evolutionary algorithm based on decomposition[J]. IEEE J of Biomedical and Health Informatics, 2016, 20(2): 717-727.
- [7] Liu Y C, Schmidt B, Maskell D L. MSAProbs: Multiple sequence alignment based on pair hidden Markov models and partition function posterior probabilities[J]. Bioinformatics, 2010, 26(16): 1958-1964.
- [8] Rani R R, Ramyachitra D. Multiple sequence alignment using multi-objective based bacterial foraging optimization algorithm[J]. Biosystems, 2016, 150(10): 177-189.
- [9] Sun J, Wu X J, Fang W, et al. Multiple sequence alignment using the hidden Markov model trained by an improved quantum-behaved particle swarm optimization[J]. Information Sciences, 2012, 182(1): 93-114.
- [10] Rubio-Largo A, Vega-Rodriguez M A, Gonzalez-Alvarez D L. A hybrid multiobjective memetic metaheuristic for multiple sequence alignment[J]. IEEE Trans on Evolutionary Computation, 2016, 20(4): 499-514.
- [11] Zambrano-Vega C, Nebro A J, Durillo J J, et al. Multiple sequence alignment with multiobjective metaheuristics: A comparative study[J]. Int J of Intelligent Systems, 2017, 32(2): 843-861.
- [12] Karaboga D. An idea based on honey bee swarm for numerical optimization[R]. Technical Report-TR06, Erciyes University, Kayseri/Turkey, 2005.
- [13] Öztürk C, Hancer E, Karaboga D. Dynamic clustering with improved binary artificial bee colony algorithm[J]. Applied Soft Computing, 2015, 28(3): 69-80.
- [14] Mao M X, Duan Q C. Modified artificial bee colony algorithm with self-adaptive extended memory[J]. Cybernetics and Systems: An Int J, 2016, 47(7): 585-601.
- [15] Kuang F J, Jin Z, Xu W H, et al. Hybridization algorithm of Tent chaos artificial bee colony and particle swarm optimization[J]. Control and Decision, 2015, 30(5): 839-847.
- [16] Wang Z G, Wang M G. Multi-search strategy of artificial bee colony algorithm based on symbolic function[J]. Control and Decision, 2016, 31(11): 2037-2044.
- [17] Edgar R C. BENCH: A collection of protein sequence alignment benchmarks including BALIBASE v3, PREFAB v4, OXBENCH, and SABRE[EB/OL]. <http://www.drive5.com/bench>, 2016.
- [18] Multiple Sequence Alignment tools[EB/OL]. [2017-02-18]. <http://www.ebi.ac.uk/Tools/msa/>.

(责任编辑: 闫妍)