



## یافتن موقعیت های دایاد با الگوریتم ژنتیک

هاشمی فر، سمیه<sup>1\*</sup> - صادقی، مهدی<sup>2,3,4\*</sup> - بهرام، گلپایانی<sup>2</sup> - نونذری، عباس<sup>1,2</sup> - اهرابیان، هایده<sup>1,2</sup> - زارع، فاطمه

<sup>1</sup>گروه علوم کامپیوتر دانشکده ریاضی، آمار و علوم کامپیوتر دانشگاه تهران

<sup>2</sup>مرکز تحقیقات بیوشیمی و بیوفیزیک دانشگاه تهران

<sup>3</sup>پژوهشکده مهندسی ژنتیک و زیست فناوری

<sup>4</sup>پژوهشکده کامپیوتر و پژوهشگاه دانش های بنیادی

### چکیده

در این مقاله یک الگوریتم ژنتیک جدید برای یافتن موتیف های دایاد ناشناخته در یک مجموعه از دنباله های داده شده ارائه شده است. الگوریتم بر روی یک مجموعه داده ای تست گردیده و نتایج آن با الگوریتم های دایاد دیگر مقایسه شده است. نتایج مقایسه، کارایی بسیار خوبی را برای الگوریتم نشان می دهد.

کلمات کلیدی: عوامل الگوبرداری، موتیف، الگوریتم های یافتن موتیف، الگوریتم ژنتیک.

### مقدمه

یکی از مسائل مهم در بیولوژی مولکولی، یافتن چگونگی ارتباط بین ژن هایی است که یک فعالیت معین را انجام می دهند و یا با یکدیگر بیان می شوند [4]. آزمایشات تجربی بر روی ناحیه پروموتور تعدادی از ژن ها نشان داده است که برای شروع پروسه نسخه برداری، یک مولکول خاص که (TF) فاکتور رونویسی نامیده می شود به زیر رشته های کوتاهی در ناحیه پروموتور ژن های بیان شده با هم پیوند می خورد. به هر کدام از این زیر رشته ها یک مکان اتصال (bs) از فاکتور رونویسی مزبور می گویند. این مکان های اتصال دارای یک الگوی مشترک هستند که اغلب موتیف یا سیگنال نامیده می شود و به مساله یافتن این الگوها یافتن موتیف می گویند [4]. ترکیبی از چندین الگوی مرتبط با هم که همزمان رخ می دهند یک الگوی ترکیبی نام دارد. مشکل اساسی در یافتن سیگنال های composite این است که بعضی از الگوهای موجود در ترکیب ممکن است خیلی ضعیف باشند که بر اساس دیدگاه های قدیمی به سختی پیدا می شوند [5]. یک نمونه از الگوهای ترکیبی سیگنال های دایاد هستند که ترکیبی از یک جفت الگوی تکی (موناد) می باشند که در فاصله ثابتی از یکدیگر رخ می دهند [2]. برخی از الگوریتم های معروف که تا کنون برای یافتن سیگنالها ارائه شده اند عبارتند از: MEME، AlignACE و MEME بر اساس تکنیک expectation-maximization، عمل می کند [1]، در حالی که AlignACE [3] با بهره گیری از روش Gibbs sampling موتیف ها را می یابد. در این مقاله یک الگوریتم ژنتیک برای یافتن سیگنال های دایاد ارائه می شود. یکی از خصوصیات بهینه این روش یافتن دو بخش سیگنال دایاد به طور همزمان می باشد که منجر به یافتن الگوی موتیف به نحو موثرتری می شود. تعمیم این الگوریتم برای حالت سیگنال های ترکیبی نیز براحتی قابل انجام است.

### مواد و روشها

### مجموعه داده ها

مجموعه داده هایی که در این مقاله مورد آزمایش قرار گرفته است، 5 تنظیم کننده مخمر است که در پایگاه داده scpd گزارش شده است [4]. در این پایگاه همراه با هر موتیف ژن های تنظیم شده بوسیله آن و موقعیت مصداق هایش ذکر شده است. برای هر تنظیم کننده ناحیه ای به طول 850 bp از موقعیت 800- تا +50 از ژن های بیان شده بوسیله آن آزمایش شده است.

### الگوریتم یافتن موتیف



مجموعه  $\{s_1, \dots, s_t\}$  متشکل از  $t$  رشته به طول  $n$  داده شده است و هدف یافتن موتیف ناشناخته در بین آنها است. یک موتیف  $E$  دایاد  $E$  به طول  $L$  از دو بخش  $E_1$  به طول  $l_1$  و  $E_2$  به طول  $l_2$  ( $l_1 + l_2 = L$ ) تشکیل شده است که در فاصله  $k$  از  $s_1$  آغاز شده و تمام  $(n - L + 1) \times (l_{\max} - l_{\min} + 1)$  زیر رشته  $L$  تایید دایاد موجود در  $s_1$  با  $E$  مقایسه می گردد و در نهایت زیر رشته  $E$  با بیشترین تعداد تطابق به عنوان اولین مصداق ( $B_1$ ) در نظر گرفته می شود. تعداد تطابق های بین  $E$  با هر کدام از زیر رشته های  $[s_1[i] \dots s_1[i + l_1 + l_2 - 1]]$  برابر با جمع تعداد تطابق های  $E_1$  و  $E_2$  با زیر رشته های  $[s_1[i] \dots s_1[i + l_1 - 1]]$  و  $[s_1[i + l_1 + k] \dots s_1[i + l_1 + l_2 - 1]]$  به طول  $l_2$  می باشد که به صورت زیر محاسبه می گردد:

$$M(s_1[i] \dots s_1[i + l_1 + l_2 - 1], E[1] \dots E[l_1 + l_2]) = \sum_{j=1}^{l_1+l_2} d(s_1[i + j - 1], E[j]),$$

$$d(a, b) = \begin{cases} 1 & \text{if } a = b \\ 0 & \text{otherwise.} \end{cases}$$

به ازای دیگر رشته های  $s_i$  ( $2 \leq i \leq t$ ) این عمل تکرار می شود تا در نهایت مصداق های  $\{B_1 \dots B_t\}$  تشکیل گردد. با توجه به مطالب فوق الگوریتم ژنتیک یافتن موتیف های دایاد به صورت زیر می باشد:

1. تشکیل جمعیت اولیه  $p_0$  به سایز  $Psize$ .
2. ارزش دهی به اعضای  $p_0$  بر اساس تابع برازندگی  $F$  و مرتب کردن آن ها.
3. مقدار دهی  $k = 1$  و انجام عملیات زیر به ازای  $N$  نسل:
  - انتقال  $(1 - pc) \times Psize$  اعضای بهتر جمعیت  $P_{k-1}$  به  $P_k$ .
  - انجام عمل crossover روی  $pc \times Psize / 2$  عضو از  $P_{k-1}$  و انتقال آنها به  $P_k$ .
  - مرتب کردن  $P_k$  بر اساس تابع برازندگی.
  - انجام عمل mutation روی  $pm \times Psize$  عضو از  $P_k$  و افزایش 1 واحد به  $k$ .
4. انتخاب 10 موتیف با رتبه بالاتر به عنوان نتیجه نهایی.

برای تشکیل جمعیت اولیه  $p_0$  الگوریتم، از هر رشته  $s_i$  ( $2 \leq i \leq t$ )،  $Psize/t$  موتیف به صورت زیر انتخاب می گردد: به ازای موقعیت تصادفی  $j$  از رشته  $s_i$  ( $1 \leq j \leq n - L + 1$ ) تمام  $(l_{\max} - l_{\min} + 1)$  موتیف دایاد ممکن که از دو زیر رشته  $E_1$  در موقعیت  $j$  ام و  $E_2$  در موقعیت  $(j + l_1 + k)$  تشکیل شده است، بررسی می شود. اگر تعداد تطابق ها از یک حد آستانه  $\alpha$  بیشتر باشد و موتیف تصادفی نباشد به عنوان یکی از افراد جمعیت اولیه  $p_0$  در نظر گرفته می شود. این عمل تا زمانی انجام می شود که  $Psize/t$  موتیف مناسب از هر رشته  $s_i$  انتخاب شود. در نهایت  $Psize$  عضو در جمعیت اولیه  $p_0$  به دست خواهد آمد. چون در حالت ایده آل تعداد کل تطابق ها برابر با  $t \times (l_1 + l_2)$  خواهد بود می توان  $\alpha = \lfloor t \times (l_1 + l_2) \times (3/4) \rfloor + 1$  در نظر گرفت. برای بررسی تصادفی بودن موتیف، 20 ناحیه پروموتور از مخمر به طور تصادفی انتخاب می شود. اگر موتیف مورد نظر حداقل با  $\beta$  زیر رشته از این 20 ناحیه به تعداد  $(l_1 + l_2)/2$  تطابق داشته باشد، تصادفی قلمداد می شود. برای محاسبه  $\beta$ ،  $C_x$  برابر با تعداد رخداد های رشته ای به طول  $L$  که از کارا کتر  $x$  تشکیل شده است،



در نظر گرفته می شود که حداقل  $(l_1 + l_2)/2$  تطابق را با نواحی تصادفی انتخاب شده داشته باشد. آنگاه

$$\beta = (C_A + C_C + C_G + C_T)/4$$

برای ارزشگذاری هر موتیف از مجموع سه تابع برآزندگی  $F_{SP}$ ،  $F_{NM}$  و  $F_{IC}$  استفاده می شود:  $F = F_{SP} + F_{NM} + F_{IC}$ . برای محاسبه  $F_{SP}$  که بیانگر *sum of pairs* است، ابتدا تمامی مصداق های الگوی  $E$  یعنی  $\{B_1 \dots B_r\}$  به دست آمده و سپس

مجموع تعداد تطابق های بین هر جفت از مصداق های مطابق فرمول زیر محاسبه می گردد:  $F_{SP} = \sum_{j=1}^l \sum_{r=1}^l M(B_j, B_r)$ .

$F_{NM}$  برابر با مجموع کل تعداد تطابق های بین الگوی  $E$  با همه مصداق های آن می باشد:  $F_{NM} = \sum_{j=1}^l M(E, B_j)$ .

برای محاسبه  $F_{IC}$  که بیانگر Information content از ماتریس وزنی موقعیت مربوط به مصداق های موتیف طبق فرمول

زیر استفاده می شود:  $F_{IC} = \sum_{i=1}^4 \sum_{j=1}^t p_{ij} \times \log(p_{ij} / p_0)$  که  $p_{ij}$  عناصر ماتریس وزنی موقعیت و  $p_0$  احتمال رخداد

نوکلوتید در زمینه است. سپس مقادیر حاصل از این سه تابع با تقسیم بر مقدار ماکزیم هر ارزش، به مقداری در بازه  $[0,1]$  تبدیل می گردد. مقدار ماکزیم برای  $F_{NM}$  برابر با  $t \times (l_1 + l_2)$ ، برای  $F_{IC}$  برابر با  $2 \times (l_1 + l_2)$  و برای  $F_{SP}$  برابر با  $t \times (t-1) / 2 \times (l_1 + l_2)$  می باشد.

در الگوریتم ژنتیک ارائه شده، از crossover یک نقطه ای استفاده می شود. به این منظور دو موتیف  $x_i$  و  $x_j$  از جمعیت را به عنوان دو والد منتخب توسط roulette wheel در نظر گرفته و مصداق های  $\{B_i, \dots, B_i\}$  و  $\{B_j, \dots, B_j\}$  یافته می شود. سپس این دو مجموعه  $t$  تایی در نقطه تصادفی  $r$  ( $2 \times t \leq r \leq t - 2 \times t$ ) با هم crossover شده و دو فرزند جدید به دست می آید. برای mutation ابتدا ماتریس وزنی موقعیت مربوط به مصداق های موتیف محاسبه می شود. هر ستون  $j$  از ماتریس که مقدار Information content آن کمتر از 1 است به عنوان یک ستون غیر conserve در نظر گرفته می شود. سپس روی هر ستون غیر conserve، با توجه به عدد تصادفی  $q_i$  به صورت زیر عمل می شود:

$$\begin{cases} B_{ij} = C & \text{if } q_i \leq .3 \\ B_{ij} = G & \text{if } .3 < q_i \leq .6 \\ B_{ij} = T & \text{if } .6 < q_i \leq .8 \\ B_{ij} = A & \text{otherwise.} \end{cases}$$

برای یافتن مصداق های نهایی هر موتیف، ماتریس وزنی موقعیت آن با رشته های  $\{S_1, \dots, S_l\}$  متوازن می گردد. برای توازن ماتریس وزنی موقعیت به هر زیر رشته دایاد  $x = x[1], \dots, x[l_1 + l_2]$  که در فاصله  $k$  از یکدیگر قرار دارند رتبه ای به

صورت  $score(x) = \sum_{i=1}^4 \sum_{j=1}^{l_1+l_2} p_{ij} \times S_{ij}$  اختصاص داده می شود که  $S_{ij} = \begin{cases} 1 & \text{if } x[j] = \text{ith base in } \{A, C, G, T\} \\ 0 & \text{otherwise} \end{cases}$ .

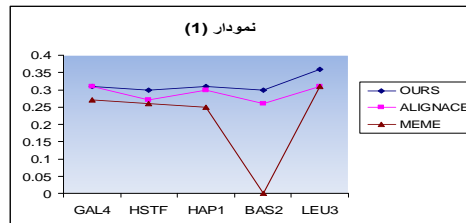
سپس در هر رشته  $S_i$  زیر رشته های با بهترین رتبه به عنوان مصداق انتخاب می شوند.

#### نتایج

برای مقایسه الگوریتم ارائه شده با دو روش AlignACE 3.0 و MEME 3.5.4 از ترکیب پنج معیار زیر استفاده شده است: (1) positive coefficient (2) performance coefficient (3) correlation coefficient (4) specificity (5) sensitivity



prediction . این مقیاس ها در سطح نوکلئوتید موتیف ها ارزیابی می شوند . نمودار 1 نتایج مقایسه سه الگوریتم برای داده های ذکر شده را نشان می دهد.



تشکر و قدردانی

بخشی از هزینه های این پروژه توسط دانشگاه تهران و بخشی از آن توسط پژوهشگاه دانش های بنیادی (CS-1385-1-02) تامین می شود.

منابع

1. Bailey T., Elkan C. (1995) Unsupervised learning of multiple motifs in biopolymers using expectation maximization. *Machine Learning* 21:51-80.
2. Helden J., Rios A.F., Collado-Vides j. (2000) Discovering Regulatory Elements in Non-Coding Sequences by Analysis of Spaced Dyads. *Nucleic Acids Research* 28:1808-18.
3. Hughes J.D., et al. (2000) Computational identification of cisregulatory elements associated with groups of functionally related genes in *Saccharomyces cerevisiae*. *Molecular biology* 296:1205-14.
4. Sinha S., Tompa M. (2003) YMF: A program for discovery of novel transcription factor binding sites by statistical overrepresentation. *Nucleic Acids Research* 31(13):3586-8
5. Tompa M., Li N., et al. (2005) Assessing Computational Tools for the Discovery of Transcription Factor Binding Sites. *Nature Biotechnology* 23:137-144.

#### Abstract

A major challenge in molecular biology is to understand the mechanisms that regulate the expression of genes. An important step in this challenge is the ability to identify regulatory elements, notably the binding sites in DNA for transcription factors. These binding sites are called motifs.

In this paper a new genetic algorithm for finding unknown dyad motifs in a given set of sequences is presented. Not many algorithms for finding dyad motif are presented in the literature and no genetic algorithm is given till now. In our genetic algorithm a novel multi-objective fitness function is employed for the selection of best individuals. The fitness function is designed based on instance consensus motifs and position weight matrix. The algorithm is tested on the different sets of real data. The result of algorithm is compared with the result of two well known algorithms MEME and AlignACE.