



الگوریتمی برای پیدا کردن موتیف های

توالی DNA بر مبنای جبر خطی

شیخ عطار علیرضا¹ و²، اصلاح چی چنگیز² و¹، پزشک حمید³، صادقی مهدی⁴
1-دانشکده ریاضی، دانشگاه شهید بهشتی 2-پژوهشکده کامپیوتر، پژوهشگاه
دانشهای بنیادی 3- گروه آمار دانشکده ریاضی آمار و علوم کامپیوتر
دانشگاه تهران 4-پژوهشگاه ملی مهندسی ژنتیک و زیست شناسی

چکیده

بررسی الگوهای تقریباً مشابه روی توالی های DNA، یا همان موتیف ها (Motifs)، یک مسئله بسیار مهم در زیست شناسی محاسباتی است که در سالهای اخیر، تحقیقات و کوشش های فراوانی به منظور دست یابی به روش مناسب تر برای یافتن آن صورت گرفته است. در این راستا ده ها روش محاسباتی متفاوت ارائه شده است که اخیراً تلاشهایی برای مقایسه نه همه، بلکه تعداد اندکی از این روش ها به عمل آمده است اما هنوز به سختی می توان با استفاده از این کوشش ها یکی از روش های مذکور را برجسته تلقی کرد. ما با دیدگاه جبر خطی، بر آنیم تا روش جدیدی برای یافتن موتیف ها ارائه دهیم.

کلمات کلیدی: موتیف، ردیف سازی توالی، پروفایل، نظریه گراف

مقدمه

یکی از مسئله های بسیار مهمی که زیست شناسان را با چالشی بزرگ مواجه کرده است، دانستن ساز و کاری است که ناظر بر قاعده مندی بیان ژن می باشد. گام اساسی برای مو شکافی دقیق تر چنین نظام حاکم بر ژن ها، روشن ساختن مناطق قاعده مندی روی DNA است. این مطلب منجر به مسئله یافتن موتیف ها می شود: "با در نظر گرفتن تعدادی از توالی های DNA داده شده، مناطقی از آنها را شناسایی کنید که نامزد خوبی برای موتیف ها باشند." انتظار می رود این مناطق با توجه به قاعده مندی بیان ژن، همگی الگوی تقریباً مشابهی داشته باشند. این تشابه باعث می شود مسئله یافتن موتیف ها که شالوده ای زیستی دارد به مسئله ای کاملاً محاسباتی تبدیل شود.

تئوری و روش

در اکثر روش های موجود برای یافتن موتیف ها، هدف پیدا کردن زیررشته هایی هم اندازه به طول l از توالی های DNA است که حداکثر d نوکلئوتید غیرمنطبق داشته باشند. در روش ارائه شده خود را محدود به یافتن موتیف هایی به طول خاص نکرده ایم و علاوه بر آن به جای حداکثر d نوکلئوتید غیرمنطبق، از مفهومی به نام همردیفی محلی بهره برده ایم. در این روش ابتدا هر توالی DNA طبق نظام خاصی قطعه قطعه می شود. در واقع به هر نوکلئوتید بر اساس میزان پیش بینی پذیری اش نسبت به چند نوکلئوتید قبلی، یکی از اعداد 1، 2، 3 و یا 4 را نسبت می دهیم. فرض کنید f نگاشتی است که نوکلئوتیدهای C, T, A و G را به ترتیب به چهار پایه استاندارد مربوط به فضای برداری R^4 ، یعنی e_1, e_2, e_3, e_4 می نگارد. اگر چهار نوکلئوتید پشت سرهم روی یک توالی DNA باشد،

منظور از یک ترکیب خطی برای آن عبارت است از $\sum_{n=0}^r c_n f(N_{a-n})$ که در

آن c_n ها چهار عدد حقیقی در فاصله $[-1, 1]$ هستند. هدف یافتن c_n هایی

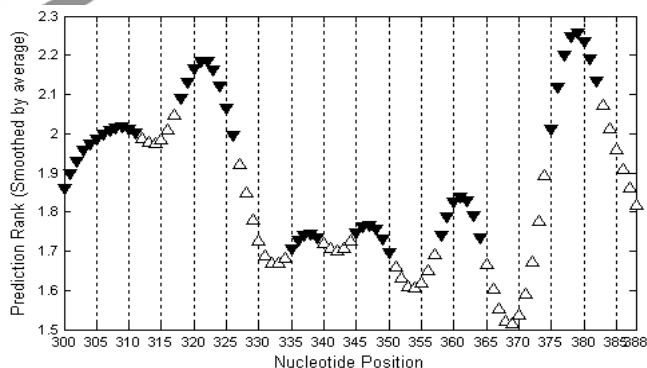


است که فاصله اقلیدسی $\sum_{n=0}^r c_n f(N_{a-n})$ تا $f(N_{a+1})$ کمتر یا مساوی با فاصله اقلیدسی $\sum_{n=0}^r c_n f(N_{a-n})$ تا هر پایه استاندارد دیگری شود. در واقع ضرایبی را خواهیم یافت که بتوانیم نوکلئوتید پنجم را پیش بینی کنیم. اگر منظور از $\|V_1, V_2\|$ فاصله اقلیدسی دو بردار V_1 و V_2 باشد، نقطه ای که در آن تابع زیر کمینه می شود، گزینه خوبی برای c_n های مورد انتظار خواهند بود:

$$P(c_0, c_1, c_2, c_3) = \sum_{t=0}^r \left\| f(N_{a+t}), \sum_{n=0}^r c_n f(N_{a-n+t}) \right\|^2$$

یعنی در طول هشت نوکلئوتید، مجموع مربعات فاصله هر چهار نوکلئوتید تا نوکلئوتید پنجم کمینه شده است. حال رتبه $\left\| f(N_{a+1+t}), \sum_{n=0}^r c_n f(N_{a-n+t}) \right\|$ در بین مجموعه $\left\{ \left\| e_x, \sum_{n=0}^r c_n f(N_{a-n+t}) \right\| \right\}_{x=1,2,3,4}$ را به

N_{a+1+t} نسبت می دهیم و از نو همین کار را تکرار می کنیم تا نموداری بدست آوریم که به هر نوکلئوتید عددی از مجموعه $\{1,2,3,4\}$ نسبت می دهد. البته این نمودار را با روش میانگین گیری هموار می کنیم (شکل 1). دامنه این نمودار می تواند به چهار ناحیه متفاوت تقسیم شود: 1) تقعر های بالا 2) تقعر های پایین 3) شیب های مثبت 4) شیب های منفی. در واقع می توان از روی هر توالی DNA چهار نسخه قطعه شده بدست آورد که هر نسخه شامل قطعه هایی است که از تقسیمات ناحیه ای 1 تا 4 بدست آمده اند (شکل 1). به طور شهودی هر قطعه گویای میزان پیچیدگی ترتیب ظاهر شدن نوکلئوتید ها در آن ناحیه روی توالی DNA است. چون $Motif$ ها دارای الگوی تقریباً مشابهی هستند، انتظار می رود فراوانی آنها در یکی از چهار ناحیه مذکور بیشتر باشد. حال یک قطعه از DNA_1 به نام seg_1 را با قطعه ای از DNA_2 به نام seg_2 به طور محلی همردیف دوگانه می کنیم. اگر هر قطعه بدست آمده روی هر DNA را به صورت یک رأس از یک گراف نمایش دهیم، seg_1 را به seg_2 متصل می کنیم اگر و تنها اگر امتیاز همردیفی بیش از 16 باشد که این امتیاز از آزمایش های فراوان بدست آمده است. این کار را برای هر دو قطعه از دو توالی DNA متمایز انجام می دهیم تا گرافی چند بخشی بدست آوریم (شکل 2). در واقع قطعه های ایجاد شده از یک DNA همگی در یک بخش قرار دارند. زیر گراف های کامل ماکسیمال در این گراف را پیدا می کنیم. به این زیرگراف های کامل، میانگین امتیاز همردیفی دوگانه رأس هایشان را نسبت می دهیم و آنها را بر حسب این امتیاز در یک لیست رتبه بندی می کنیم. برای یافتن یک موتیف اولین زیر گراف کامل را از لیست انتخاب می کنیم. هر قطعه را به صورت یک



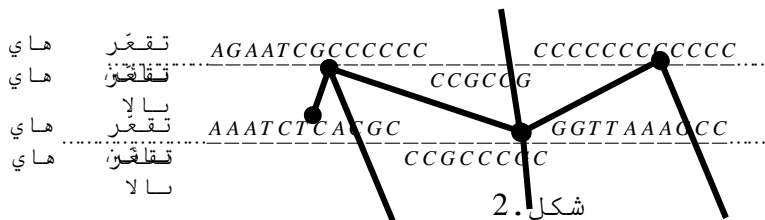
پروفایل (profile)

در نظر می گیریم که تعداد ظاهر شدن هر نوکلئوتید در مکان ها یش را نشان می دهد. حال با روش ردیف سازی پروفایل دو قطعه ای را بدست می آوریم که بیشترین امتیاز همردیفی پروفایل را داشته باشند. منظور از ادغام دو

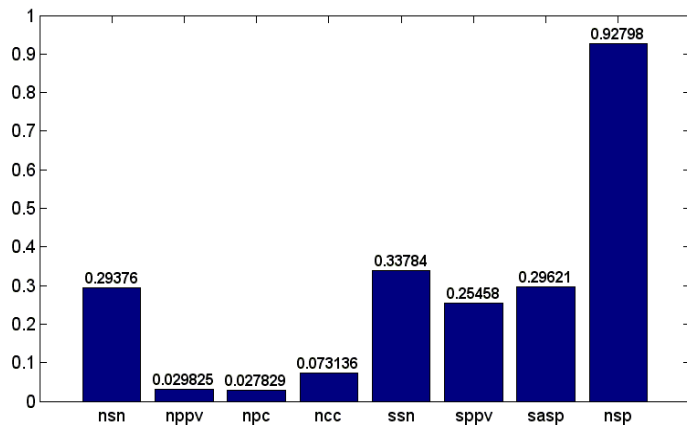
پروفایل، ماتریسی است که از مجموع درایه های فاصله ای که دو پروفایل با هم همریف شده اند، بدست می آید. این دو پروفایل را حذف می کنیم و به جای آن این دو را با هم ادغام می کنیم. سپس این کار را تا زمانی که یک پروفایل بدست آید ادامه می دهیم. این پروفایل را به یک ماتریس احتمال تبدیل می کنیم و آن را روی توالی DNA حرکت می دهیم تا ناحیه ای از توالی DNA را بیابیم که با توجه به پروفایل بیشینه امتیاز را کسب کرده باشد. حاصل یک پیشنهاد برای موتیف است. این کار را روی DNA های دیگر نیز انجام می دهیم تا دیگر پیشنهاد ها بدست آید.

نتایج و بحث

شکل 1 مقدار عددی رتبه پیش بینی پذیری هر نوکلئوتید را نشان می دهد که با روش میانگین گیری هموار شده است. در این نمودار می توان تقریب های بالا و پائین را از هم تشخیص داد.



شکل 2 قسمتی از گرافی را نشان می دهد که در آن تنها رئوس مربوط به تقریب ها دیده می شود.



مقدار عددی آماره های Tompa¹ برای 50 توالی واقعی DNA مربوط به نتایج حاصل از الگوریتم در نمودار روبرو نشان داده شده است. این توالی ها در <http://tare.medisin.ntnu.no/> قابل دسترسی است.

شکل 3



تشکر و قدردانی

بخشی از هزینه این کار از محل طرح شماره (CS 1385-1-02) پژوهشگاه دانشهای بنیادی (IPM) تأمین شده است.

منابع

1. Tompa M, Li N, Bailey TL, Church GM, De Moor B, Eskin E, Favorov AV, Frith MC, Fu Y, Kent WJ, Makeev VJ, Mironov AA, Noble WS, Pavese G, Pesole G, Regnier M, Simonis N, Sinha S, Thijs G, van Helden J, Vandenbogaert M, Weng Z, Workman C, Ye C, Zhu Z (2005) Assessing computational tools for the discovery of transcription factor binding sites. *Nat Biotechnol* 2005, 23:137-44.
2. Geir Kjetil Sandve, Osman Abul, Vegard Walseng and Finn Drabløs (2007) Improved benchmarks for computational motif discovery. *BMC Bioinformatics* 2007, 8:193
3. Elena Zaslavsky and Mona Singh (2006) A combinatorial optimization approach for diverse motif finding Applications. *Algorithms for Molecular Biology* 2006, 1:13



پنجمین همایش ملی بیوتکنولوژی جمهوری اسلامی ایران
3-5 آذر ماه 1386، سالن اجلاس سران
The 5th National Biotechnology Congress of Iran
24-26 Nov, 2007, Summit Meeting Conference Hall, Tehran- Iran



An algorithm based on linear algebra for DNA sequence motif finding

Sheikh Attar Alireza^{1,2*}, Eslahchi Changiz^{1,2**}, Pezeshk Hamid³, Sadeghi Mehdi⁴

1. Faculty of Mathematical Sciences, Shahid Beheshti University, Tehran

2. School of Computer Science, Institute for Studies in Theoretical Physics and Mathematics (IPM), Tehran

3. Center of Excellence in Biomathematics, School of Mathematics, Statistics and Computer Sciences, University College of Science, University of Tehran, Tehran, Iran,

4. National Institute of Genetic Engineering and Biotechnology, Tehran, Iran

Abstract

One of the most important problems in molecular biology is identification of repeated patterns in DNA sequences that are presumed to have a biological function. These patterns are called Motifs. Many computational methods have been proposed and most of them have been designed to find motifs with the fix length (d) and with at most L nucleotides mismatch between each pairs. However, we can get better results by considering sequences that are not restricted on d and L . We introduce a versatile linear algebra approach for the motif finding problem. We assign e_1, e_2, e_3 and e_4 , four standard basis of \mathbb{R}^4 , to A, T, C and G respectively. So we have sequences of vectors instead of nucleotides. By means of linear algebra concepts we define the predictability of each nucleotide in its appearance order. This allows us to cut the DNA to subsequences called segments. Each segment is in fact one of the four kinds of regions in predictability plot. If we show any segments as a vertex of a graph, two vertexes are connected if and only if their pairwise local alignment score is more than 16. So each maximal cliques in this graph derives a profile for a motif.

WWW.IBP.IR

iranian bioinformatics portal