



الگوریتمی برای بلاک بندی هاپلوتیپ ها بر مبنای درخت های فیلوژنتیک کامل

اصلاحی چنگیز^{1,3,*}، بزشک حمید²، صادقی مهدی⁴، افضلی نرجس¹

- 1- دانشکده ریاضی، دانشگاه شهید بهشتی
- 2- گروه بیوانفورماتیک، دانشکده ریاضی، آمار و علوم کامپیوتر، پردیس علوم، دانشگاه تهران
- 3- گروه بیوانفورماتیک، پژوهشکده علوم کامپیوتر، پژوهشگاه دانش های بنیادین
- 4- پژوهشگاه ملی مهندسی ژنتیک و زیست فناوری

چکیده فارسی:

در این مقاله ما روشی برای بلوک بندی هاپلوتایپ های یک همیت ارائه داده ایم که بر اساس دو پارامتر این کار را انجام می دهد. پارامتر اول میزان نزدیکی بلوک را به بلوک کامل نشان می دهد و پارامتر دوم میزان هاپلوتایپ های تکراری را معین می کند. همچنین با گرفتن پارامتر های سومی را که درصد هاپلوتایپ های را که کاربر می خواهد توسط Tag SNPs شناسایی شوند میگیریم و روشی برای یافتن Tag SNPs ارائه می دهیم.

کلمات کلیدی:

بازتولید هاپلوتایپ، مساله NP-سخت، الگوریتم خوشه بندی

مقدمه:

تفاوت های ژنتیکی بین انسان ها در نواحی کوچکی از نوکلئوتیدها به نام اسنیپ (SNP) رخ می دهد. در واقع SNP ها مهم ترین عامل بروز تفاوت های ژنتیکی در انسان های مختلف هستند و در تشخیص بیماری ها و ساخت دارو های جدید از اهمیت بالا برخوردار هستند. به دلایلی که هنوز معلوم نیست در هر موقعیت اسنیپ به طور معمول دو نوع نوکلئوتید ظاهر می شود. در هر موقعیت اسنیپ می توان علل با فراوانی بیشتر را با 0 و دیگری را با 1 نشان داد. به نظر می آید که انتقال کروموزوم ها از پدر و مادر به بچه در قطعات معینی صورت می گیرد که به آن بلوک گفته می شود یک بلوک را Perfect می نامیم اگر در آن در هر موقعیت SNP بیش از یک بار جهش در طول تاریخ رخ نداده باشد. از طرفی بلوک بندی را می توان بر اساس پارامتر های دیگر بنا کرد از جمله میزان تکرار هاپلوتایپ ها (Common) و تحلیل پیوست اسنیپ ها (Linkage Analysis). بلوک ها دقیقاً Prefect نیستند ولی از میزان Prefect بودن بالایی برخوردار هستند. همچنین میزان Common بودن در هر بلوک هم می تواند متغیر باشد. در این مقاله امکان بلوک بندی با مقادیر مختلف برای این پارامترها وجود دارد. نکته دیگر در مورد tag SNP هاست یعنی زیر مجموعه ای از SNP های یک بلوک که به کمک آن ها تمام اسنیپ های دیگر در هاپلوتایپ ها قابل شناسایی باشد

مواد و روش ها:



بخش اول: برای مجموعه داده‌های معین که شامل چندین هاپلوتایپ باشد. ما ابتدا هاپلوتایپ‌های تکراری را حذف کرده و سپس روی مجموعه‌ی جدید بلوک بندی می‌کنیم. برای مجموعه داده‌هایی که در آن یک سری از داده‌ها مخدوش هستند (missing data) مفهوم جدیدی به نام Prefect تعریف کرده‌ایم که با یکی کردن تمام هاپلوتایپ‌هایی که در داده‌های غیر مخدوش مشترک هستند بدست می‌آید. می‌توان ثابت کرد اگر مجموعه Prefect های حاصل Prefect باشد مجموعه‌ی هاپلوتایپ‌های اولیه نیز Prefect است. ما دو پارامتر (α, σ) را که اعدادی بین 0 و 1 هستند می‌گیریم که α در صد Common بودن σ میزان نزدیکی بلوک به بلوک Prefect را نشان می‌دهد. به طوری که با حداکثر η تغییر در ورودی‌ها بلوک Prefect شود. η را به صورت زیر تعریف می‌کنیم:

$$\eta = \begin{cases} \sigma & \sigma = 1 \\ (a - \sigma)mn + 1 & \sigma \neq 1 \end{cases}$$

که در آن m تعداد هاپلوتایپ‌ها و n تعداد اسنیپ‌هاست. همچنین بلوک بندی دارای این خاصیت است که $\alpha\%$ از هاپلوتایپ‌ها Common باشد. که این کار را به صورت حریصانه با شروع از 2 ستون اول انجام می‌دهیم. اگر دو ستون در این شرایط صدق کند که به سراغ ستون بعدی می‌رویم و گرنه ستون اول به عنوان بلوک در نظر گرفته می‌شود. در مرحله‌ی k ام، $k+1$ ستون اول را در نظر می‌گیریم. به کمک دو روش حریصانه و برنامه‌سازی پویا، تقریب مناسبی برای کمترین تعداد تغییرات لازم برای Prefect کردن بلوک بدست می‌آوریم و سپس چک می‌کنیم اگر در شرایط صدق کند به مرحله‌ی $k+1$ می‌رویم و در غیر اینصورت k ستون اول را به عنوان بلوک در نظر می‌گیریم و بلوک‌های بعدی را نیز به همین شیوه بلوک بندی می‌کنیم.

نتایج و بحث:

بعد از بلوک بندی برای هر بلوک تقریب مناسبی برای کمترین تعداد tag SNP ها بدست می‌آوریم. ثابت کرده‌ایم که اگر بلوک با m هاپلوتایپ و n اسنیپ، Prefect و بدون داده‌های مخدوش باشد، آنگاه کمترین تعداد tag SNP ها دقیقاً برابر $m-1$ است که الگوریتم‌ها آن را در زمان $O(mn)$ می‌یابد. در جدول 1 الگوریتم خود را با پارامترهای $(\alpha, \sigma) = (0, 1)$ و $(\alpha, \sigma) = (0, 0.98)$ با الگوریتم حریصانه‌ی پاتیل (Patil et al.) و الگوریتم برنامه‌سازی پویای ژانگ (Zhang et al) که هر دو بر پایه‌ی Common بودن هستند مقایسه کرده‌ایم. همانطور که مشاهده می‌شود. تمام این نتایج مربوط به کروموزم 21 انسان است که از 20 هاپلوتایپ و 24407 اسنیپ تشکیل شده است که در 4



بخش مختلف قرار دارند. پارامترهای $(\alpha, \sigma) = (0, 1)$ بیان می‌دارد که بلوک باید کاملاً Perfect باشد ولی روی میزان Common بودن قیدی نمی‌گذارد. پارامترهای $(\alpha, \sigma) = (0, 0.98)$ حداکثر 2% در ایده‌ها مجاز به تغییر می‌داند تا بلوک Perfect گردد و دو بار روی میزان Common بودن قیدی نمی‌گذارد. در حالت $(\alpha, \sigma) = (0, 0.98)$ نسبت به پارامترهای $(1, 0)$ تعداد بلوک‌ها به طور قابل ملاحظه‌ای از 3737 به 1860 در حالت $(0, 0.98)$ کاهش می‌یابد که این مقدار در الگوریتم پاتیل برابر 4135 و در الگوریتم ژانگ 2575 می‌باشد. در حالت $(\alpha, \sigma) = (0, 1)$ 657 بلوک بیشتر از 10 اسنیپ دارند که در حالت $Co, 0.98$ به 744 می‌رسد. در پاتیل این مقدار برابر 589 و ژانگ 742 می‌باشد. همچنین تعداد بلوک‌های با کمتر از 3 اسنیت در حالت (0.1) برابر 1116 و در حالت $(0.0/98)$ برابر 97، در پاتیل برابر 2138 و در ژانگ برابر 924 می‌باشد. متوسط تعداد هاپلوتایپ‌های پترن در حالت (0.1) برابر $3/49$ در حالت $(0.0/98)$ برابر $4/11$ و در پاتیل برابر 2072 و در ژانگ برابر $3/05$ می‌باشد. همچنین تعداد کل tag SNP ها در حالت (0.1) برابر 11271 و در حالت $(0.0/98)$ برابر 7873 می‌باشد.

جدول 1.

روش	تعداد اسنیپ‌ها در بلوک	تعداد بلوک‌ها	تعداد هاپلوتایپ‌ها Pattern	تعداد tag SNP
$(\alpha, \sigma), (0, 1)$	$n > 10$	657	4/45	3070
	$3 \leq n \leq 10$	1964	3/80	6539
	$n > 3$	1116	2/39	1608
	Total	3737	3/49	11217
$(\alpha, \sigma), (0.0/98)$	$n > 10$	744	4/36	4077
	$3 \leq n \leq 10$	1019	3/98	3604
	$n > 3$	97	3/27	192
	Total	1866	4/11	7873
Patil (Greedy)	$n > 10$	589	3/75	
	$3 \leq n \leq 10$	1408	2/92	
	$n > 3$	2138	2/30	
	Total	4135	2/72	
Zhang	$n > 10$	474	4/23	
	$3 \leq n \leq 10$	909	3/03	



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

(Dynamic Programming)	$n > 3$	924	2/12	
	Total	2575	3/05	

تشکر و قدردانی: بخشی از هزینه این کار از محل طرح شماره CS1385-1-02 پژوهشگاه دانشهای بنیادی (IPM) تامین شده است.

منابع:

1. Chou, P.Y., Fasman, G.D. (1974). Prediction of protein conformation. *Biochemistry* 13: 222-245.
2. Garnier J, *et al.*. (1978). Analysis of the Accuracy and Implications of Simple Methods for Predicting the Secondary Structure of Globular Proteins. *J. Mol. Biol.* 120: 97-120.
3. J. Randy Macdonald, *et al.*. (2001). Environmental features are important in determining protein secondary structure. *Protein science* 10: 1172-1177.
4. Zhan-Yang Zhu, *et al.* (1996). The use of amino acid patterns of classified helices and strands in secondary structure prediction. *J. mol. biol.* 260: 261-276.



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 for Haplotype Block Partitioning based on Perfect Phylogeny

Changiz Eslahchi^{1,2,*,**}, Hamid Pezeshk⁴, Mehdi Sadeghi³, Narjess Afzaly¹

¹ Faculty of Mathematics, Shahid-Beheshti University, Tehran, Iran.

² Bioinformatics Group, School of Computer Science, IPM, Tehran, Iran.

³ National Institute of Genetic Engineering and Biotechnology, Tehran-Karaj Highway, Tehran, Iran.

⁴ Center of Excellence in Biomathematics, School of Mathematics, Statistics and Computer Sciences, College of Science, University of Tehran, Tehran, Iran.

Abstract:

Construction of two haplotypes from a set of Single Nucleotide Polymorphism (SNP) fragments is referred to as haplotype reconstruction problem. One of the most popular computational model for this problem is Minimum Error Correction (MEC). Since MEC is an NP-hard problem, here we propose a novel heuristic algorithm based on clustering analysis in data mining for haplotype reconstruction problem. In contrast to MEC model, our algorithm has polynomial time complexity and could be applied to large datasets. Based on maximum normalized hamming distance, our iterative algorithm produces two clusters of fragments; then, in each iteration, the algorithm assigns a fragment to one of the clusters. Our results suggest that the algorithm has less reconstruction error rate in comparison with other algorithms.

WWW.IBP.IR

iranian bioinformatics portal