



## ارائه روش‌های خوشه‌بندی مبتنی بر شبکه‌های عصبی در حل مسئله‌ی بازسازی هاپلوتیپ‌ها

با استفاده از اطلاعات ژنوتیپی

محمدزاده\*، جوادی<sup>1</sup> - نوذری دالینی\*\*، عباس<sup>1و2</sup> - اهرابیان، هایده<sup>1و2</sup> - معین زاده، محمد حسین<sup>1</sup>

<sup>1</sup>قطب بیومت، دانشکده ریاضی، آمار و علوم کامپیوتر، دانشگاه تهران  
<sup>2</sup>مرکز تحقیقات بیوشیمی و بیوفیزیک، دانشگاه تهران

### چکیده

مسئله بازسازی هاپلوتیپ، تقسیم کردن قطعه‌های SNP (Single Nucleotide Polymorphism Fragments) به دو قسمت و استنتاج یک جفت هاپلوتیپ از آنهاست. مدل MEC یکی از مدل‌های موجود برای حل این مسئله می‌باشد. استفاده از اطلاعات ژنوتیپ، یک مدل توسعه‌یافته برای MEC می‌باشد که MEC/GI نامیده می‌شود. با توجه به NP-hard بودن این مسئله، به دنبال الگوریتم‌هایی مکاشفه‌ای و روش‌های خوشه‌بندی برای کاهش زمان اجرا می‌باشیم. در این مقاله به ارائه روش‌هایی مبتنی بر شبکه‌های عصبی برای کاهش زمان و افزایش میزان شباهت بین هاپلوتیپ‌های واقعی و بازسازی شده و همچنین برطرف کردن کاستی‌های روش‌های ارائه شده قبلی می‌پردازیم. در نهایت نتایج آزمایشات انجام گرفته بر روی پایگاه داده‌های واقعی مقایسه می‌شود.

**کلمات کلیدی:** هاپلوتیپ، قطعه‌های SNP، اطلاعات ژنوتیپ، خوشه‌بندی، میزان بازسازی، شبکه‌های عصبی

### مقدمه

با توجه به در دسترس بودن توالی ژنوم انسان [1]، امکان بررسی در تفاوت‌های ژنتیکی و ارتباط دادن آنها به بیماری‌های پیچیده، میسر شده است [5]. زیست‌شناسان معتقدند که همگی انسان‌ها عموماً در سطح DNA یکسان و تفاوت تنها در برخی مکان‌های آن وجود دارد که سبب بروز بیماری‌های ژنتیکی می‌شوند [2و3]. این تفاوتها در SNP ظهور می‌کند. SNP، بین افراد مختلف متفاوت بوده و شایع‌ترین عامل بروز تفاوت‌های ژنتیکی است. بیشتر SNPها از دو نوع مختلف که در اینجا 'A' و 'B' نامیده می‌شوند، تشکیل شده‌اند. توالی SNPها در هر یک از دو کروموزوم مربوط به یک ژنوم، هاپلوتیپ (پدری و مادری) نامیده می‌شود دو مکان SNP بر روی رشته ژنوتیپ اگر مقادیر یکسانی داشته باشند، آن مکان هموزیگوت و در غیر اینصورت هتروزیگوت نامیده می‌شود. بدلیل اینکه استخراج هاپلوتیپ‌ها از طریق روش‌های آزمایشگاهی بسیار دشوار و هزینه بر می‌باشد، استفاده از روش‌های محاسباتی مورد توجه قرار گرفته است. متأسفانه بدست آوردن قطعه‌های SNP (داده‌های مساله) همواره با درصدی خطا میسر می‌باشد و حل مساله را با مشکل مواجه کرده است. برای حل این مشکل، چندین مدل از طریق روش‌های محاسباتی ارائه شده است. مدل حداقل تصحیح خطا (MEC) [4] یکی از مدل‌های متداول می‌باشد. MEC به تصحیح تفاوت‌های موجود بین هاپلوتیپ و توالی SNPها می‌پردازد تا قطعات SNP را به دو کلاس افراز نماید. از هر کلاس یک هاپلوتیپ بدست آورده می‌شود. از طرف دیگر با توجه به اینکه MEC یک مسئله NP-Hard است [5] برای حل آن از روش‌های خوشه‌بندی و الگوریتم‌های مکاشفه‌ای استفاده شده است. استفاده از اطلاعات ژنوتیپ (که بسیار آسان‌تر از هاپلوتیپ به دست می‌آیند) یک مدل توسعه یافته از مدل MEC می‌باشد که MEC/GI نامیده می‌شود.

### مواد و روشها



برای حل این مساله از مجموعه داده ACE (Angiotensin Converting Enzyme) [7] استفاده شده است. تغییرات انجام شده بر روی داده ها همانند ژنگ، داده ها را به یک ماتریس  $M_{20 \times 52}$  تبدیل کرده است [5]. این ماتریس با توجه به نرخ گپ که همان از دست دادن اطلاعات میباشد،  $g=0.25, 0.5, 0.75$ ، و نرخ خطا در خواندن  $e=0.2, 0.3, 0.4$  در 9 مجموعه بدست می آیند. فرض کنیم  $m=20$  m قطعه SNP داریم (رشته-هایی از 'A'، 'B' یا '!' (شکاف)) که هر یک با یکی از دو هاپلوتیپ به طول  $n$  ( $n=52$ ) متناظر می باشد. پارتیشن P قطعه ها را به دو کلاس  $C_1$  و  $C_2$  تقسیم می کند. برای بازسازی هاپلوتیپها، قطعه های موجود در هر کلاس، با یکدیگر ترکیب می شوند. این عمل که با فرمول 1 انجام می گردد، رای گیری نامیده می شود.  $N_A^j(C_i)$  تعداد 'A' های ستون j از قطعه های مربوط به کلاس  $C_i$  می باشد. هر ژنوتیپ که از ترکیب دو هاپلوتیپ بدست می آید، رشته ایست که از 'A'، 'B' یا '!' تشکیل شده است (فرمول 2). میزان بازسازی (RR) که برای مقایسه کارایی الگوریتم های طراحی شده به کار رفته و بر اساس فاصله هامینگ بین قطعه ها محاسبه می شود (فرمول 3). در این رابطه،  $HD(h_i, h_k) = \sum_{j=1}^n d(h_{ij}, h_{kj})$  فاصله ی همینگ

و  $d(x, y)$  مقدار تفاوت موجود بر روی دو قطعه در یک مکان خاص می باشد.

$$V_{ij} = \begin{cases} A & N_A^j(C_i) > N_B^j(C_i) \\ B & \text{otherwise} \end{cases}$$

$$i = 1, 2$$

$$0 \leq j < n$$

$$g_{i, (1 \leq i \leq n)} = \begin{cases} A & (h_{1i} = h_{2i} = A) \\ B & (h_{1i} = h_{2i} = B) \\ - & \text{otherwise} \end{cases}$$

(فرمول 2)

$$RR(h, h') = 1 - \frac{\min(r_{11} + r_{22}, r_{12} + r_{21})}{2n}$$

$$r_{ij} = HD(h_i, h'_j), i, j \in \{1, 2\}$$

$$d(m_{ij}, m_{kj}) = \begin{cases} 1 & (m_{ij} \neq m_{kj} \neq -) \\ 0 & \text{otherwise} \end{cases}$$

(فرمول 3)

(فرمول 1)

ژنگ شبکه عصبی نظارتی دو لایه ای (شکل 1) را بر مبنای مدل MCIH (نام دیگر MEC/GI) را طراحی کرده است. این شبکه عصبی پیشرو از انتشار خطای گره های خروجی تنها در لایه ی اول استفاده می نماید، پس دارای m گره ورودی (هر کدام برای یک قطعه)، دو گره ی میانی (هاپلوتیپها و جواب نهایی مسئله) و یک گره ی خروجی به عنوان ژنوتیپ می باشد. روش ارائه شده در [5]، دارای مشکلاتی از قبیل عدم استقلال شبکه از تعداد قطعه ها و زمان اجرای زیاد می باشد. از طرف دیگر محدود بودن به اطلاعات ژنوتیپی باعث عدم کارایی آن در مدل MEC شده است. برای رفع این مشکلات، دو روش خوشه بندی به ازای دو مدل MEC (دارای دو قسمت شبکه ی عصبی و پس پردازش) و MEC/GI توسط شبکه های عصبی در این مقاله ارائه شده است.

**شبکه ی عصبی:** مسئله ی بازسازی هاپلوتیپها جزء مسائل غیر نظارتی طبقه بندی می شود. شبکه طراحی شده ما دو لایه و غیر نظارتی می باشد (شکل 1-ب). قطعات SNP به عنوان ورودی در نظر گرفته می شوند. لایه اول این شبکه، دارای n گره ورودی، که هر گره به یک SNP مربوط می شود. در لایه ی دوم دو گره، که هر گره مربوط به یک هاپلوتیپ می شود. هاپلوتیپها در طی آموزش به تدریج بر روی وزن های شبکه، به صورت اعداد اعشاری شکل می یابند و در نهایت از این وزن ها استخراج می شوند (الگوریتم 1).

**پس پردازش:** ابتدا با استفاده از فرمول 4 هاپلوتیپهای پیشنهادی اولیه که آنها را شبه هاپلوتیپ (*semi-haplotype*) می نامیم، را بدست می آوریم. سپس با فرض اینکه P قطعه های SNP را به دو خوشه  $C_1$  و  $C_2$  تقسیم می کند (فرمول 4)، هاپلوتیپهای نهایی را از طریق رأی گیری (فرمول 5) محاسبه می کنیم.



روش‌های دوم (مبتنی بر MEC/GI): اضافه کردن اطلاعات ژنوتیپ به داده‌های ورودی، بازسازی هاپلوتیپ‌های صحیح‌تری را نتیجه می‌دهد. در این قسمت، از روش خوشه‌بندی در مدل MEC، به همراه استفاده از اطلاعات ژنوتیپ استفاده شده است. اگر برابر با 'A' ('B') باشد، نشان‌دهنده این است که همه عناصر موجود بر روی ستون  $i$  در  $M$  باید 'A' ('B') بوده و مقدار 'B' ('A') در این ستون‌ها نشان‌دهنده وجود خطاست. با این روش تقریباً نیمی از ستون‌ها حذف خواهند شد. پس از این مرحله، شبکه‌ی عصبی برای خوشه‌بندی قطعه‌ها، اعمال می‌شوند. برای استنتاج هاپلوتیپ‌ها از دو دسته حاصل، تابع رأی‌گیری جدیدی تعریف می‌کنیم (فرمول 6).

**Algorithm: Pseudo code for new approach**  
**Input:** SNP fragments **Output:** two haplotypes  
**Step1:** Initialize weights ( $w_{ij}=0$ ).  
 Set topological neighborhood parameters ( $R = 0$ ).  
 Set learning rate parameters ( $\alpha = 1/n$ ).  
**Step2:** While stopping condition is false do, step 3-7.  
**Step3:** For each input vector  $x$ , do step 4-6.  
**Step4:** For each  $l = \{1, 2\}$ , compute:  $D(l) = \sum_j (w_{lj} \cdot x_j)$   
**Step5:** Find index  $L$  such that  $D(L)$  is minimum.  
**Step6:** For unit  $L$ , and for all  $i$ :  $w_{iL}(new) = w_{iL}(old) + \alpha \cdot x_i$   
**Step7** Test stopping condition.  
 Post-Processing.

الگوریتم 1

$$semi - haploype_{ik} = \begin{cases} B & (w_{ik} > 0) \\ A & (w_{ik} < 0) \\ Random & otherwise \end{cases}$$

$$k = \{1, 2\}$$

$$C_1 = \{f_i : HD(h_1, m_i) < HD(h_2, m_i)\}$$

$$C_2 = \{f_i : HD(h_2, m_i) \leq HD(h_1, m_i)\}$$

$$i = 1, 2, 3, \dots, m$$

$$v_{ik} = \begin{cases} A & N_A^i(C_k) > N_B^i(C_k) \\ B & N_A^i(C_k) < N_B^i(C_k) \\ Random & otherwise \end{cases}$$

$$k = \{1, 2\}$$

(فرمول 5)

$$V_{ij} = \begin{cases} A & N_A^j(C_i)/N_B^j(C_i) > N_A^j(C_{i'})/N_B^j(C_{i'}) \\ B & N_A^j(C_i)/N_B^j(C_i) < N_A^j(C_{i'})/N_B^j(C_{i'}) \\ K_{ij} & otherwise \end{cases}$$

$$K_{ij} = \begin{cases} A & N_A^j(C_i) - N_A^j(C_{i'}) \geq N_B^j(C_i) - N_B^j(C_{i'}) \\ B & otherwise \end{cases}$$

$$i, i' \in \{1, 2\}, i \neq i'$$

$$0 \leq j < n$$

(فرمول 6)

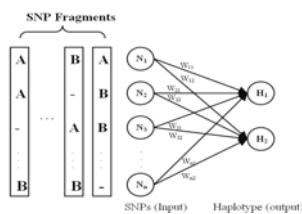
(فرمول 4)

## نتایج و بحث

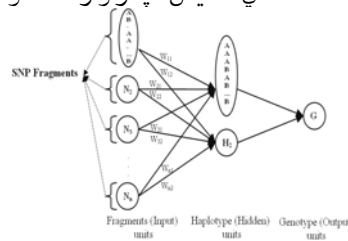
نتایج حاصل که بر روی پایگاه داده ACE آزمایش شده بود در نمودار 3 نشان داده شده است. این نمودار حاکیست که روش ما در خطای بالا (0.3) و دارای RR بیشتر میباشد و به مراتب بهتر از روش ارائه شده در [5] عمل می‌کند. از طرفی وارد کردن قطعات SNP به صورت سری، باعث عدم وابستگی این روش به تعداد قطعات ورودی می‌شود. نتایج حاصله در مدل MEC نیز، با توجه به عدم وجود اطلاعات ژنوتیپی قابل قبول می‌باشد (نمودار 1).

## تشکر و قدردانی

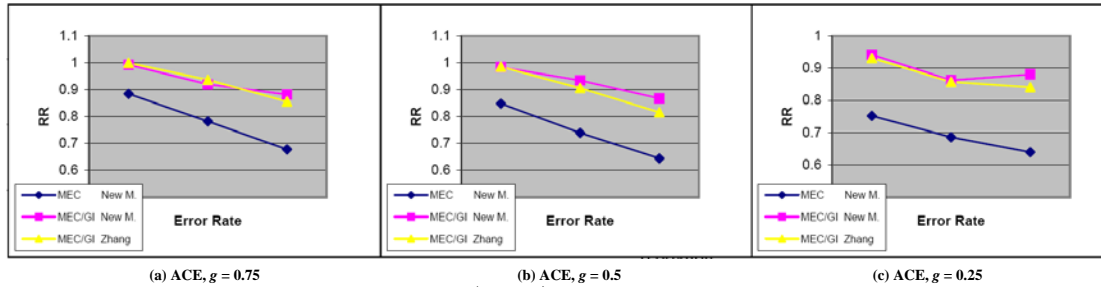
بخشی از هزینه‌های این پروژه توسط دانشگاه تهران تامین گردیده است.



شکل 1- ب



شکل 1- الف



نمودار 1: مقایسه‌ی RR در روش‌ها و مدل‌های مختلف

## مراجع

- [1] Venter, J.C. and Adams, M.D. et al. 2001. The sequence of the human genome. *Science*, 291(5507):1304–1351.
- [2] Terwilliger, J. & Weiss, K. (1998). Linkage disequilibrium mapping of complex disease: Fantasy and reality? *Curr. Opin. Biotechnology*, 579–594.
- [3] Chakravarti, A. 1998. It's raining, hallelujah? *Nature Genetics*, 19:216–217.
- [4] Wang, R., Wu, L., Li, Z. and Zhang, X. 2005. Haplotype reconstruction from SNP fragments by Minimum Error Correction. *Bioinformatics*, 21(10):2456–2462.
- [5] X.Zhang, R.Wang, L.Wu, W.Zhang. 2006. Minimum conflict individual Haplotyping from SNP fragments and related Genotype, *Bioinformatics* 271-280.
- [6] X.Zhang, R.Wang, L.Wu, W.Zhang. 2007. A clustering algorithm based on two distance functions for MEC model, *Computational biology and chemistry*:148,150.
- [7] Rieder, M., Taylor, S., Clark, A. and Nickerson, D. 1999. Sequence variation in the human angiotensin converting enzyme. *Nature genetics*, 22:59–62.

### Abstract

Most positions of the human genome are typically invariant (99%) and only some positions (1%) are commonly invariant which are associated with complex genetic diseases. Haplotype reconstruction is to divide aligned SNP fragments, which is the most frequent form of difference to address genetic diseases, into two classes, and thus inferring a pair of haplotypes from them. Minimum error correction (MEC) is an important model for this problem but only effective when the error rate of the fragments is low. MEC/GI as an extension to MEC employs the related genotype information besides the SNP fragments and so results in a more accurate inference. The haplotyping problem, due to its NP-hardness, may have no efficient algorithm for exact solution. In this paper, we design an unsupervised neural network based on MEC and MEC/GI model. As numerical results on real biological data, this neural network approach works well and an increase in the rate of similarity between the real haplotypes and the reconstructed ones is gained.