

نشست دو روزه حل مساله در بیوانفورماتیک

نوکلئوتیدی وجود ندارد و تعداد آنها بسیار محدود است. به دلایل تجربی مشخص شده که هاپلوتیپ‌ها را می‌توان به بلوک‌هایی با تنوع محدودی از هاپلوتیپ‌ها تقسیم‌بندی کرد. در هر بلوک تعداد کمی از SNP‌ها (SNP‌های نماینده) کافی خواهد بود که به طور منحصر به فرد هاپلوتیپ‌های مشترک را تشخیص دهند.

مساله اصلی این است: چگونه می‌توان این توالی SNP‌ها را به تعدادی بلوک تقسیم‌بندی کرد به طوری که تعداد SNP‌های نماینده مینیمم باشد.

برای حل مساله چند شرط را باید در نظر گرفت. هاپلوتیپ‌های مشترک، هاپلوتیپ‌هایی هستند که بیشتر از یک بار در یک بلوک دیده شوند. بنابراین ما نیاز داریم که در نتیجهٔ نهایی تقسیم‌بندی، بخش قابل توجهی از هاپلوتیپ‌ها در هر بلوک، هاپلوتیپ‌های مشترک باشند (مشلاً ۸۰٪). فرض کنید که k هاپلوتیپ (از k فرد مختلف) با n SNP متولی داریم. $r_i, i = 1, 2, \dots, n$ یک بردار k بعدی است که $r_i(k) = 0, 1, 2$ دوآل کامین هاپلوتیپ در موقعیت SNP_i است که $r_i(k) = 0, 1, 2$ دوآل مختلف است و صفر به معنای آن است که وضعیت آلل در آن موقعیت مشخص نشده است (به دلایل تجربی ممکن است نتوان نوکلئوتید قرار گرفته در یک موقعیت را تشخیص داد).

یک بلوک با ترتیب r_j, r_i, \dots, r_l تعریف می‌شود. دو هاپلوتیپ k و k' باهم سازگار گفته می‌شوند اگر آلل‌ها در دو هاپلوتیپ در موقعیت‌هایی که صفر نباشند، یکسان باشند. به عبارت دیگر: به ازای هر $j, l, j \leq l, i \leq n$ به طوری که $r_j(k) \neq r_j(k')$ یک هاپلوتیپ در یک بلوک «نمایم» است اگر با دو هاپلوتیپی که باهم ناسازگار هستند سازگار باشد. با در نظر گرفتن شرط سازگار بودن، ما می‌توانیم هاپلوتیپ «نمایم» را به گروه‌های مجزا یا بلوک‌ها تقسیم کنیم. دو هاپلوتیپ در یک گروه خواهند بود اگر باهم سازگار باشند. با تعریف ذکر شده در بارهٔ هاپلوتیپ‌های نمایم می‌توان تابع $f(r_i, \dots, r_l) = \text{block}(r_i, \dots, r_l)$ را تعریف کرد، اگر حداقل α درصد (مشلاً ۸۰ درصد) از هاپلوتیپ‌های نمایم در بلوک بیشتر از یک بار ظاهر شده باشند. بنابراین بلوک‌های به دست آمده در تقسیم‌بندی نهایی باید این شرایط ذکر شده را داشته باشند. همان‌گونه که قبلاً ذکر شد، می‌خواهیم تقسیم‌بندی اپتیمیم بلوک‌ها را پیدا کنیم (مشلاً B_r, \dots, B_l) به طوری که تعداد کل SNP‌های نماینده در طول دو روز برگزاری این نشست، الگوریتم‌های متعددی که توسط دانشجویان مطرح، پیاده سازی و بررسی شد که متأسفانه هیچ کدام از آنها به حل درستی از مساله منجر نشد. پس از آن تعدادی از دانشجویان کار

در روزهای ۲۲ و ۲۳ دی ماه ۸۴ یک دوره دو روزه حل مساله به پیشنهاد رئیس پژوهشکده علوم کامپیوتر، دکتر حمید سربازی آزاد و به سرپرستی دکتر روزبه توسرکانی در پژوهشگاه دانش‌های بنیادی برگزار شد که هدف از آن، کسب یک تجربه جدید در استفاده از توان دانشجویان بر جسته علوم ریاضی و کامپیوتر برای حل یک مساله باز و مهم روز بود. با توجه به اهمیت زیاد بیوانفورماتیک به عنوان یک رشتہ جدید و در حال رشد و همچنین ماهیت بین‌رشته‌ای آن که به علوم ریاضی و کامپیوتر وابسته است، بررسی یک مساله بیوانفورماتیک و تلاش برای حل آن به عنوان موضوع فعالیت انتخاب شد. بدین منظور از دکتر مهدی صادقی از پژوهشگاه ملی مهندسی ژنتیک و دکتر چنگیز اصلاحچی از دانشکده ریاضی دانشگاه شهید بهشتی که در این زمینه فعالیت می‌کنند، دعوت شد تا با طرح یک مساله، با این دانشجویان در جهت پیدا کردن راه حل آن همکاری کنند. مساله‌ای با عنوان بخش‌بندی بلوک‌های هاپلوتیپی (Haplotype Block Partitioning) برای دانشجویان مطرح شد که می‌توان آن را به طور خلاصه به شرح زیر بیان کرد.

ماده ژنتیکی یا DNA را می‌توان به صورت یک رشتہ پیوسته از چهار مولکول متفاوت (نوکلئوتید) که با نمادهای A, C, G, T مشخص می‌شوند تعریف کرد. طول این رشتہ در انسان حدود ۳ میلیارد است که در ۲۳ چفت کروموزوم با اندازه‌های مختلف تقسیم شده است. در واقع هر کروموزوم یک زنجیره بسیار طویل با اندازهٔ متوسط حدود ۱۵۰ میلیون نوکلئوتید است.

نشان داده شده است که در افراد مختلف حدود سه میلیون موقعیت در رشتہ DNA وجود دارد که در این موقعیت‌ها، نوع نوکلئوتید قرار گرفته در انسان‌ها متفاوت است و در بقیه موقعیت‌ها (۹۹٪) نوکلئوتیدها شبیه به هم هستند. همچنین نشان داده شده است که در هر کدام از این موقعیت‌ها فقط دو نوکلئوتید از چهار نوکلئوتید می‌توانند قرار گیرند. این موقعیت‌ها چند ریختی تک نوکلئوتیدی (Single Nucleotide Polymorphism) یا به اختصار SNP نامیده می‌شوند و دو نوکلئوتید قرار گرفته در این موقعیت‌ها آلل (Allele) نام دارند. بنابراین، چنانچه ترتیب سه میلیارد نوکلئوتید در یک فرد شناخته شده باشد در افراد دیگر فقط با دانستن نوکلئوتید قرار گرفته در موقعیت‌های ترتیب کل توالی نوکلئوتیدها شناخته خواهد شد. از آنجایی که در هر فرد در موقعیت‌های SNP یکی از دو نوکلئوتید می‌تواند قرار گیرد، تعداد کل توالی‌های متفاوت ممکن در انسان 2^n است که n تعداد این موقعیت‌های SNP است. ترتیب نوکلئوتیدها در این موقعیت‌ها در یک فرد، هاپلوتیپ آن فرد نامیده می‌شود. از راه‌های تجربی و نظری نشان داده شده است که در بین انسان‌ها همه‌این 2^n تنوع

بر روی حل مسئله را ادامه دادند که نهایتاً یک فرمول‌بندی بر اساس نظریه اطلاعات به نتایج قابل توجهی منجر شد که کار بر روی آن و استخراج اطلاعات از یک پایگاه داده واقعی مربوط به ژئوم انسان و نهایتاً تهیه یک

کامپیوتر، دانشگاه صنعتی شریف در حال انجام است که انتظار می‌رود به زودی تکمیل شود.



چه کشوری سریعترین رشد را در تولیدات علمی دارد؟

ممکن است بگویید چین؛ ولی اشتباه می‌کنید. طبق تحقیقی که اخیراً مؤسسه بازرگانی و صنعت بریتانیا (DTI) درباره رتبه‌بندی کشورها از لحاظ علمی انجام داده، کشور دیگری حائز این رتبه است. در این بررسی، ا نوع شاخص‌های ورودی و خروجی برای بدست آوردن تصویر نسبتاً کاملی از توزیع فعالیت و تأثیر علمی در سراسر جهان مورد توجه قرار گرفته است. علی‌رغم پاسخ تعبیج‌آوری که عنوان این نوشته دارد (و در پایان مطلب آن را افشا خواهم کرد)، این گزارش رشد سریع چین را به مثابه یک قدرت علمی که برتری سابقاً بی‌چون و چرای آمریکا را به چالش طلبیده، تأیید می‌کند و نیز (در چشم‌انداز محدود ملی) مؤید عملکرد نسبتاً قوی بریتانیاست که نسبت به رقباًش پول کمتری صرف تحقیق می‌کند و پژوهشگران کمتری دارد ولی در مقایسه، دانش‌بیشتری با تأثیر زیاد تولید می‌کند.

در زمینه بودجه تحقیقات علمی، صعود چین بیش از هر کشوری چشمگیر است. بودجه پژوهشی چین (با در نظر گرفتن ملاحظات مربوط به قدرت خرید) طی دهه گذشته چهار برابر شده و اکنون بیشتر از هر کشوری در جهان بجز آمریکا و ژاپن است و به نصف کل بودجه تحقیقات در اتحادیه اروپا رسیده است. از لحاظ میزان انتشارات علمی، سهم چین ۵٪ کل جهان است و در دهه گذشته سه برابر شده و امروز بیشتر از فرانسه است. اگر هر کشوری جداگانه در نظر گرفته شود، ایالات متحده آمریکا هنوز هم در صدر قرار دارد و تقریباً یک سوم کل محصول علمی جهان را تولید می‌کند، ولی اتحادیه اروپا در مجموع از آمریکا پیشی‌گرفته و سهم آن ۳۷/۹٪ از کل انتشارات است. بریتانیا با حدود ۹٪ مقام دوم را در بین تک تک کشورها دارد و اخیراً ژاپن را پشت سر نهاده است. گروه کشورهای آسیلپاسیفیک (چین، کره، چین، کره، ژاپن، و سنگاپور) روی هم رفته ۱۰٪ از تولید علمی جهان را در دست دارند. و اما درباره کیفیت و تأثیر، در این مورد آمریکا هنوز به وضوح پیشگام است؛ اگر سهم هر کشور از بر ارجاع‌ترین مقاله‌ها (با در نظر گرفتن ۶۱٪ بالا در هر رشته) به عنوان شاخصی از تأثیر جهانی قلمداد شود، آمریکا با ۶۱٪ در رأس قرار دارد. بریتانیا سهم خود را از مقاله‌های پر ارجاع افزایش داده و به ۱۳٪ رسانده است. چن از این لحاظ هنوز نسبت به صدرنشینان عقب‌تر است ولی فاصله‌اش از آنها در حال کم شدن است به خصوص که شمار ارجاعات، یک شاخص تأخیری است یعنی مدتی طول می‌کشد تا شمرة صرف بودجه در تحقیقات به صورت انتشار مقاله‌ها ظاهر شود و پس از آن، پژوهشگران دیگر به آن مقاله‌ها استناد کنند.

کشوری که رشد انتشارات علمی آن طی دهه گذشته بیش از همه بوده، ایران است که میزان تولید آن ده برابر شده است (هرچند نسبت به میزانی اندک در آغاز دهه). باید دید که آیا با توجه به تحولات سیاسی اخیر، ایران می‌تواند همچنان به عملکرد خوب خود در این زمینه ادامه دهد یا خیر.

منبع:

<http://www.softmachines.org.wordpress/>