

# تخمین هاپلوتایپ با استفاده از فاکتورسازی ماتریس رتبه پایین در حضور داده‌های پرت

تینا تربالی<sup>۱</sup> سینا مجیدیان<sup>۲</sup> محمد حسین کهای<sup>۳</sup>

۱- دانشجوی کارشناسی ارشد- دانشکده مهندسی برق- دانشگاه علم و صنعت ایران- تهران- ایران

[tarbali.tina@gmail.com](mailto:tarbali.tina@gmail.com)

۲- دانشجوی دکتری- دانشکده مهندسی برق- دانشگاه علم و صنعت ایران- تهران- ایران

[s\\_majidian@yahoo.com](mailto:s_majidian@yahoo.com)

۳- دانشیار- دانشکده مهندسی برق- دانشگاه علم و صنعت ایران- تهران- ایران

[kahaei@iust.ac.ir](mailto:kahaei@iust.ac.ir)

چکیده: تخمین هاپلوتایپ بر اساس اطلاعات DNA برای کشف بیماری‌های ژنتیکی انسان استفاده می‌شود. این مسئله در پردازش ژنومی سیگنال‌ها به صورت یک ماتریس رتبه پایین قابل مدل سازی است که به علت محدودیت‌های موجود در خوانش هاپلوتایپ، فقط تعداد کمی از درایه‌ها مشاهده می‌شوند. در نتیجه یک روش موثر برای بازیابی هاپلوتایپ از مشاهدات ناقص، استفاده از روش‌های تکمیل ماتریس است. در این مقاله به کمک روش‌های تکمیل ماتریس، سعی در تخمین هاپلوتایپ از طریق فاکتورسازی ماتریسی شده است. در مراجع از روش گرادیان کاهشی برای حل مسئله استفاده شده است. اما در روش‌های قبلی داده‌های پرت نیز در محاسبات وارد می‌شود که باعث خطا در تخمین هاپلوتایپ شده است. به عبارتی در این روش‌ها به شروط موجود برای ماتریس‌های هاپلوتایپ توجه نشده است و این موضوع باعث تخمین داده‌های پرت برای هاپلوتایپ شده است. در این مقاله با روش تکمیل ماتریس و با در نظر گرفتن این شروط در ماتریس هاپلوتایپ، یک تابع هزینه جدید به صورت عبارت جریمه برای تخمین هاپلوتایپ معرفی می‌کنیم. عبارت جدید اضافه شده به تابع هزینه باعث می‌شود که اثر داده‌های پرت کاهش یافته و در نتیجه دقت تخمین هاپلوتایپ افزایش می‌یابد. نتایج شبیه سازی نیاز کاهش خطای بازیابی هاپلوتایپ را تایید می‌کند.

واژه‌های کلیدی: تخمین هاپلوتایپ، ماتریس رتبه پایین، تکمیل ماتریس

نوع مقاله: پژوهشی

DOI: 10.29252/jiaeee.18.3.1007

تاریخ ارسال مقاله: ۱۳۹۸/۰۷/۲۸

تاریخ پذیرش مشروط مقاله: ۱۳۹۹/۰۷/۱۳

تاریخ پذیرش مقاله: ۱۳۹۹/۰۹/۱۷

نام نویسنده‌ی مسئول: دکتر محمد حسین کهای

نشانی نویسنده‌ی مسئول: ایران - تهران - بزرگراه رسالت - خیابان هنگام - دانشگاه علم و صنعت ایران - دانشکده‌ی مهندسی برق

## ۱- مقدمه

به علت ارتباط بسیار زیاد بین انواع بیماری‌ها با مسائل ژنتیکی، استخراج اطلاعات از داده‌های ژنتیکی از جمله داده‌های پزشکی و رشته DNA<sup>۱</sup>، نقش مهمی دارد. همچنین به منظور ایجاد توانایی در تشخیص الگوی بیماری‌ها، باید از داده‌های ژنتیکی در دسترس استفاده کرد. به همین علت رشته DNA به یک ابزار همه جانبه و یک عامل کلیدی برای علم پزشکی تبدیل شده است. رشته DNA معمولاً مجموعه داده‌های بسیار زیادی را تولید می‌کند که پردازش آن‌ها چالش‌های محاسباتی پیچیده‌ای را ارائه می‌دهد. در نتیجه یافتن روش مناسب برای استفاده از این داده‌ها، بسیار مورد توجه است.

مسئله تخمین هاپلوتاایپ با استفاده از مدل‌سازی ماتریسی به عنوان یکی از زمینه‌های تحقیقاتی بین رشته‌ای مطرح شده است. تخمین هاپلوتاایپ به کشف ارتباط داده ژنتیکی و بیماری‌ها کمک کرده که منجر به طراحی دارو و روش‌های درمانی مبتنی بر فرد می‌گردد [۴].

برای تخمین هاپلوتاایپ از روش‌هایی مانند الگوریتم انتقال پیام [۵] و SDP<sup>۲</sup> [۶] استفاده شده است. در این روش‌ها از گراف برای مدل‌سازی مسئله استفاده شده است. در صورت وجود نویز، امکان تخمین هاپلوتاایپ به کمک الگوریتم‌های بیان شده میسر نیست. همچنین میتوان از روش‌های پردازش سیگنال<sup>۳</sup> دیجیتال برای تخمین هاپلوتاایپ استفاده کرد. این روش‌ها در سال‌های اخیر برای تجزیه و تحلیل داده‌های ژنتیکی در حال توسعه هستند. از جمله مهم‌ترین مفاهیم مورد استفاده در این روش برای داده‌های ژنتیکی میتوان به تبدیل فوری به گسسته<sup>۴</sup>، فیلتر دیجیتال، تبدیل موجک دیجیتال<sup>۵</sup> اشاره کرد [۷].

یکی از روش‌های تخمین هاپلوتاایپ روش تکمیل ماتریس است [۴]. روش‌های مختلفی برای تکمیل ماتریس از جمله کمینه‌سازی متناوب، ADMiR<sup>۶</sup>، SET<sup>۷</sup> و AP<sup>۸</sup> مطرح می‌شود [۸-۱۰]. یکی از روش‌های تکمیل ماتریس استفاده از فاکتورسازی ماتریسی است که ماتریس رتبه پایین به صورت حاصل ضرب دو ماتریس با ابعاد کمتر تجزیه و در هر مرحله با ثابت فرض شدن یکی از ماتریس‌ها مسئله نسبت به ماتریس دیگر کمینه می‌شود و محدب بودن هر دو ماتریس تجزیه شده به حل مسئله کمک شایانی می‌کند [۱۱، ۱۲].

یکی از روش‌های تخمین هاپلوتاایپ، حل مسئله به کمک فاکتورسازی ماتریسی است. استفاده از این روش در پردازش سیگنال و تخمین پارامترها بسیار متداول است [۱۳]. اما در الگوریتم‌های پیاد شده به کمک این روش امکان حضور داده‌های پرت در جواب مسئله وجود دارد. این امر به این علت است که ویژگی‌های موجود در داده‌های ماتریس تجزیه شده هاپلوتاایپ در مراحل حل مسئله بی اثر بوده است.

در این مقاله به کمک روش‌های تکمیل ماتریس، سعی در تخمین هاپلوتاایپ از طریق فاکتورسازی ماتریسی شده است. با توجه به اینکه

داده‌های ژنتیکی علاوه بر رتبه پایین بودن ماتریس مشاهدات، دارای ویژگی‌هایی در ساختار ماتریس‌ها است، تابع هزینه جدید به صورت عبارت جریمه برای مسئله بیان می‌شود. تخمین هاپلوتاایپ به کمک الگوریتم گرادیان کاهشی قابل حل است. عبارت جدید اضافه شده به مسئله، به تخمین هاپلوتاایپ با دقت بیشتر کمک کرده و بازیابی را بهبود بخشیده است. زیرا عبارت تنظیم اضافه شده به تابع هزینه باعث کنترل درایه‌های ماتریس و عدم انتخاب داده پرت می‌شود. نتایج حاصل از شبیه‌سازی موید این مطلب است.

## ۲- بیان مسئله

## ۲-۱- معرفی مفاهیم زیستی

نوکلئیک اسید یکی از اجزای تشکیل‌دهنده هسته سلول انسان است که بر دو نوع ریبونوکلئیک اسید<sup>۹</sup> و دئوکسی‌ریبونوکلئیک اسید<sup>۱۰</sup> یافت می‌شود. DNA رشته‌ای از نوکلئیک اسیدها است که شامل بسیاری از اجزای کوچکتر مرتبط با نوکلئوتید می‌باشد. هر نوکلئوتید شامل یکی از چهار اسیدهای آمینه ممکن یعنی آدنین<sup>۱۱</sup>، تیمین<sup>۱۲</sup>، سیتوزین<sup>۱۳</sup> و گوانین<sup>۱۴</sup> است. DNA دارای ساختار دو رشته‌ای است که در جهت مخالف هم اجرا می‌شوند. به طور مثال در انسان تعداد کروموزوم‌ها<sup>۱۵</sup> ۴۶ عدد هستند که شامل دو سری یکسان ۲۳ تایی هستند که هر کدام از آن‌ها از یکی از والدین به ارث رسیده است. کروموزوم جفت را کروموزوم‌های همسان<sup>۱۶</sup> می‌گویند. مجموع DNA در هر سری از این ۲۳ کروموزوم شامل حدود ۳ میلیارد جفت باز می‌باشد. جانداري که دو نسخه‌ی یک DNA را داشته باشد اصطلاحاً دیپلوئید<sup>۱۷</sup> نامیده می‌شود. DNA وظیفه‌ی انتقال و ذخیره اطلاعات ژنتیکی را بر عهده دارد. به عبارتی اطلاعات بنیادی درباره‌ی نحوه‌ی زندگی یک موجود زنده در DNA قرار دارد. شایع‌ترین تغییرات بین دو کروموزوم در یک جفت همولوگ عبارتند از چند ریختی تک نوکلئوتیدی<sup>۱۸</sup> که اسنپ خوانده می‌شود. یک اسنپ جایگاهی از ژن است که در بین افراد جامعه، به کمک نوکلئوتیدهای مختلفی مشاهده می‌شود. در واقع به جایگاه ژن هنگام مقایسه دو رشته، اسنپ گفته می‌شود و برای راحتی می‌توان فقط محل اسنپ را نمایش داد. از طرفی رشته‌ی اسنپ بر روی هر یک از جفت کروموزوم‌ها هاپلوتاایپ نامیده می‌شود. بنابراین، تغییرات بین دو کروموزوم در یک جفت توسط هاپلوتاایپ نشان داده می‌شود [۱۴، ۱۵].

برای به دست آوردن دنباله‌ی هاپلوتاایپ، می‌توان از روش‌های توالی‌یابی استفاده کرد. اولین روش توالی‌یابی توسط ماکسام<sup>۱۹</sup> و گیلبرت<sup>۲۰</sup> انجام شد [۱۶]. در ادامه فردریک سنگر<sup>۲۱</sup> روشی مبتنی بر خاصیت فسفر سانس ارائه کرد [۱۷]. هزینه‌ی بالا و دقت کم این روش‌ها باعث شد که از روش‌های دیگری از جمله روش توالی شاتگان<sup>۲۲</sup> استفاده شود. در این روش ژنوم به طور تصادفی در نقاط مختلف بریده می‌شود که به توالی هر کدام از این تکه‌ها، خوانش<sup>۲۳</sup> می‌گویند. به بیان دیگر در روش شاتگان از هم‌پوشانی<sup>۲۴</sup> قطعات مختلف خوانش، رشته DNA

دنباله موردنظر، روش خوانش جفت‌انتهای است. به کمک شکل (۱) می‌توان ماتریسی  $\mathbb{R}^{m \times n}$  به نام ماتریس خوانش ایجاد کرد که  $n$  نشان دهنده طول هاپلوتایپ و  $m$  مربوط به تعداد قطعات تقسیم شده برای خوانش است. ردیف  $i$ ام  $\mathbf{R}$ ، اساساً اطلاعات مربوط به هاپلوتایپ ارائه شده توسط  $i$ امین خوانش را جمع آوری می‌کند. به علت اینکه در انسان‌ها و دیگر ترکیبات دوگانه، مکان‌های اسنپ، دو قطبی هستند، به این معنی که در هر موقعیت هاپلوتایپ تنها دو از چهار نوکلئوتید A، C، T و G ممکن است رخ دهد. بنابراین می‌توان نوکلئوتیدها را در موقعیت اسنپ با استفاده از نمادهای باینری نمادگذاری کرد [۴]. این نمادگذاری با توجه به شماره‌ی جایگاه اسنپ و همچنین فراوانی یا نادر بودن آن نوکلئوتید در اسنپ تعریف می‌شود. جهت سهولت بررسی، درایه‌هایی از  $r_i$  که فاقد اطلاعات اسنپ هستند، صفر در نظر گرفته می‌شوند. به کمک این نمادگذاری، ماتریس  $\mathbf{R}$  شامل ورودی‌های مبنای سه تایی است. به طور خاص، درایه  $i$ ،  $(i, j)$  (ماتریس  $\mathbf{R}$ ، اطلاعات مربوط به موقعیت  $j$ امین اسنپ توسط  $i$ امین خوانش را ارائه می‌دهد. اگر  $i$ امین ورودی موقعیت  $j$ امین اسنپ را پوشش ندهد، درایه  $(i, j)$  از ماتریس  $\mathbf{R}$  را  $R_{ij} = 0$  نتیجه می‌دهد. در (۱) ماتریسی از خوانش قطعه‌ای اسنپ ارائه شده است [۴].

$$\mathbf{R} = \begin{bmatrix} 0 & 1 & 0 & 0 & -1 & 0 & 0 \\ -1 & 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & -1 & -1 & 0 \\ 1 & 0 & 0 & -1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & -1 & 1 & 0 & -1 & 0 & 0 \end{bmatrix} \quad (1)$$

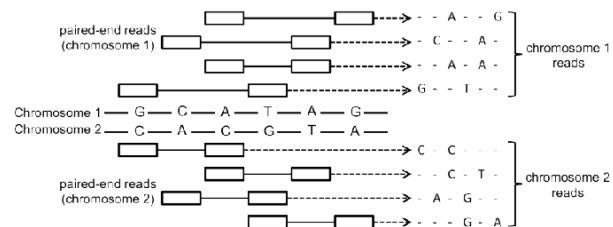
همانطور که بیان شد در ماتریس حاصل از خوانش جفت‌انتهای تمام اطلاعات در دسترس نیست. به همین علت به کمک روش‌های بازایی تکمیل ماتریس رتبه‌پایین سعی در تخمین هاپلوتایپ شده است. در مسئله ماتریس  $\mathbf{R}$ ، ماتریس مشاهدات نویزی و ماتریسی رتبه‌پایین است و  $\mathbf{M}$  ماتریس هدف است که ماتریس بازایی شده از  $\mathbf{R}$  است. هر سطر از ماتریس  $\mathbf{R}$  مربوط به یک نمونه از هاپلوتایپ است. تجزیه‌ی ماتریسی به صورت  $\mathbf{M} = \mathbf{UV}^T$  صورت می‌گیرد. تخمین هاپلوتایپ از تخمین ماتریس  $\mathbf{M}$  حاصل می‌شود

$$\mathbf{M} = \mathbf{UV}^T \quad (2)$$

$\mathbf{U}$  و  $\mathbf{V}$  به ترتیب ماتریس‌های  $n \times k$  و  $m \times k$  و  $k$  نشان‌دهنده‌ی  $k$ -پلوئیدی بودن جاندار است. ستون‌نام از  $\mathbf{V}$ ،  $v_j$ ، دنباله‌ای از  $i$ امین هاپلوتایپ است.  $i$ امین ردیف از  $\mathbf{U}$ ،  $u_i$ ، نشان می‌دهد که این خوانش مربوط به کدام یک از جفت کروموزوم پدری و مادری است. توجه شود که هر سطر از ماتریس  $\mathbf{M}$ ،  $m_i$ ، یک رشته کامل هاپلوتایپ است. به عنوان مثال،  $m_i \in \mathbf{H}$ ، از آنجایی که  $k < m, k < n$  است طبق رابطه‌ی  $rank(\mathbf{M}) \leq \min\{rank(\mathbf{U}), rank(\mathbf{V})\}$ ، رتبه‌ی ماتریس  $\mathbf{M}$  حداکثر برابر با  $k$  خواهد بود. به همین علت این مدل، رتبه‌پایین است [۴]. شایان ذکر است که مدل فاکتورسازی ماتریسی پیش از این در مسائل

حاصل شده است. به علت انتخاب تصادفی این قطعات، بهتر است که هر نقطه از ژنوم توسط تعداد معینی از قطعات هم‌پوشانی شود تا احتمال خطا در این نتیجه‌گیری کم شود [۱۸].

روش توالی‌یابی شانگان قادر به خوانش رشته DNA از طریق جفت‌انتهای است. این روش مجموعه‌ای از خوانش‌ها را تولید می‌کند که در آن هر خوانش به مجموعه‌ای از اطلاعات جزئی در مورد کروموزوم که از آن تولید شده است، مربوط می‌شود. در مواردی که تخمین هاپلوتایپ مورد نظر است این هم‌پوشانی و مقایسه به تفاوت جایگاه بین کروموزوم‌ها می‌پردازد. در این روش بدون دانستن محتوای نوکلئوتید و فقط با شناختن طول، نمونه‌برداری صورت می‌گیرد. به عبارت دیگر یک خوانش با طول مشخص در دسترس است که ابتدا و انتهای آن مشخص است اما میان آن مشخص نیست [۴، ۱۹]. شکل (۱) نشان‌دهنده‌ی خوانش جفت‌انتهای برای بازایی هاپلوتایپ است.



شکل (۱): چگونگی خوانش جفت کروموزوم‌ها [۴]

در شکل (۱) برای هر کروموزوم چهار خوانش انجام شده است و در هر خوانش دو جایگاه مشاهده شده و نوکلئوتید مربوط به آن جایگاه ارائه شده است و بقیه جایگاه‌ها با '-' نمایش داده شده‌اند. منظور از '-' آن است که خوانش، اطلاعاتی درباره‌ی آن جایگاه نوکلئوتید ارائه نکرده است. برای توضیح بیشتر شکل می‌توان توالی کروموزوم ۱ را در نظر گرفت. کروموزوم ۱ شامل نوکلئوتیدهای GCATAG است که در خوانش اول نوکلئوتیدهای دو جایگاه ارائه شده است که به ترتیب A و G را نمایش می‌دهد و دیگر جایگاه‌های آن با '-' نمایش داده شده است که نشان از عدم خوانش کروموزوم است. به همین ترتیب خوانش دوم نیز دو جایگاه کروموزوم ۱ را نشان داده شده است که به ترتیب نوکلئوتید C و A ارائه شده است. همین روند برای هر دو کروموزوم و تمام خوانش‌ها انجام شده است و شکل (۱) نشان از اطلاعات این خوانش‌ها دارد.

## ۲-۲- مدل ماتریس و ساختار مسئله

برای استفاده از روش‌های محاسبات مهندسی و پردازشی نیاز است که رشته DNA به عدد تبدیل شود [۱، ۱۷]. یعنی به هریک از نوکلئوتید اسیدهای آدنین، سیتوزین، گوانین و سیتوزین یک عدد نسبت داده شود. انتخاب مناسب این اعداد باعث بهبود عملکرد روش‌های حل مسئله می‌شود. روش‌هایی از جمله واس<sup>۲۶</sup> یا دودویی [۷] و EIIP<sup>۲۷</sup> برای این نمادگذاری مناسب است. یکی دیگر از روش‌های عددی برای بیان

درایه‌های  $\pm 1$  است که می‌توان بدون یافتن اطلاعات دقیق از هر درایه-ی ماتریس، مقدار  $\|U\|_F^2$  و  $\|V\|_F^2$  را به کمک نرم فروبنیوس به صورت رابطه (۵) و (۶) به دست آورد.

$$\|U\|_F^2 = \sum_{i,j} U_{ij}^2 = m \quad (5)$$

$$\|V\|_F^2 = \sum_{i,j} V_{ij}^2 = n \quad (6)$$

همچنین در صورت انتخاب  $k = 2$ ، ماتریس  $U$  که دارای یک درایه‌ی یک در هر سطر است و دیگر درایه‌های آن صفر است رابطه  $\|U\|_F^2 = m$  برقرار است. اما ماتریس  $V$  که دارای درایه‌های  $\pm 1$  است، رابطه (۷) جایگزین رابطه (۶) می‌شود.

$$\|V\|_F^2 = \sum_{i,j} V_{ij}^2 = 2n \quad (7)$$

در ادامه از روابط (۵) و (۶) استفاده شده است. از طرفی به کمک فاکتورسازی ماتریسی در مسئله، یکی از شروط مسئله برای بازیابی ماتریس‌ها  $(U, V)_{ij} = (UV^T)_{ij} = \pm 1$  است. هدف مسئله تخمین  $M$  است. برای محاسبه‌ی آن با تجزیه‌ی ماتریسی به دنبال  $U$  و  $V$  هستیم که شرط  $(UV^T)_{ij} = \pm 1$  را ا قناع کند. البته نه هر نوع ماتریسی که در این رابطه صدق کند. به طور مثال، اگر درایه‌های ماتریس  $U_{ij} = 1/1000$  و  $V_{ij} = 1000$  باشد، به کمک ضرب ماتریسی،  $(M)_{ij} = \pm 1$  حاصل می‌شود. اما با وجود اینکه جواب حاصل دارای شرط لازم برای درایه‌های ماتریس  $M$  است، مناسب داده‌های ژنتیکی مسئله نیست. زیرا درایه‌های ماتریس  $U$  و  $V$  نباید مقادیر بزرگی داشته باشند و این مقادیر با اطلاعاتی که در فرض مسئله وجود دارد همخوانی ندارد. می‌توان نتیجه گرفت که مسئله مورد نظر علاوه بر شرط لازم، مناسب بودن میزان بزرگی مقادیر درایه‌های ماتریس  $U$  و  $V$  به عنوان شرط کافی، نیز نیاز دارد.

در نتیجه بر اساس رابطه‌های (۵) و (۶)، مسئله‌ی بهینه‌سازی جدیدی را پیشنهاد می‌کنیم.

$$f(U, V) = \|P_\Omega(R - UV^T)\|_F^2 + \lambda_1 (\|U\|_F^2 - m) + \lambda_2 (\|V\|_F^2 - n) \quad (8)$$

که  $U \in \mathbb{R}^{m \times k}$  و  $V \in \mathbb{R}^{n \times k}$  و ماتریس رتبه پایین مشاهده شده  $R \in \mathbb{R}^{m \times n}$ ، طول هاپلوتاایپ است. در مسئله پیشنهادی پارامترهای اضافه شده، باعث لحاظ شدن شرط کافی برای مقادیر ماتریس‌های  $U$  و  $V$  هنگام کمینه‌سازی تابع هزینه، شده است.

الگوریتم گرادیان کاهشی به کمک تابع پیشنهادی بیان می‌شود.

در رابطه‌ی (۸) هدف یافتن  $U$  و  $V$  است که تابع  $f(U, V)$  را کمینه کند. محاسبه  $U$  و  $V$  به ازای مقادیر  $t=0, 1, 2, \dots$  به صورت رابطه زیر بیان می‌شود:

$$V_{t+1} = V_t - \alpha \nabla f(V_t) \quad (9)$$

مختلفی از جمله تخمین جهت منابع استفاده شده است [۲، ۳، ۷، ۲۰].

ماتریس  $M$  به کمک اپراتور معرفی شده زیر قابل توصیف است.

$$P_\Omega(M) = \begin{cases} M_{ij}, & \text{if } (i, j) \in \Omega \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

در صورتی که مسئله دارای خطا نباشد رابطه‌ی  $R = P_\Omega(M)$  برقرار است و درایه‌های  $R$  به کمک شرط  $(i, j) \in \Omega$  قابل بیان است. مجموعه‌ی  $\Omega$  خوانش‌های مربوط به کروموزوم‌های مورد بررسی را نشان می‌دهد. در حالت کلی خوانش رشته DNA دارای خطاست در نتیجه  $P_\Omega(R) \neq P_\Omega(M)$ . به عنوان مثال ورودی  $(i, j) \in \Omega$  در  $R_{ij}$ ، به صورت زیر تعیین می‌شود:

$$R = \begin{cases} M_{ij} & \text{With Probability : } 1-p \\ -M_{ij} & \text{With Probability } p \end{cases} \quad (4)$$

در رابطه (۴)،  $p$  میزان خطا دنباله است. وجود خطا درایه‌ها را منفی می‌کند. زیرا اعمال نویز با کاهش یا افزایش مقدار در درایه‌ها بی معنی است.

هر درایه از سطر  $R$  نمونه‌ای از هاپلوتاایپ است. به عنوان مثال،  $R_{ij} = -1$  به این معنی است که  $i$ امین ورودی موقعیت  $j$ امین اسنیپ را پوشش می‌دهد و اطلاعاتی را که توسط  $-1$  کدگذاری شده است را فراهم می‌کند. با این حال، اینکه کدام هاپلوتاایپ  $k$  توسط  $i$ امین ورودی نمونه‌برداری شده است، مجهول است. به طور کلی، ماتریس  $R$  را می‌توان به صورت نمونه‌گیری از ماتریس  $M$  با مقداری خطا استخراج کرد [۲۱].

### ۳- روش پیشنهادی

هدف در این مقاله بهبود خطای تخمین هاپلوتاایپ به کمک ویژگی‌های موجود در ساختار مسئله است. در ابتدا ویژگی‌های ماتریس‌های ژنتیکی مطرح می‌شود، سپس تابع هزینه جدید ارائه شده است.

ماتریس بیان شده  $U$  در بخش قبل دارای ابعاد  $m \times k$  است که با توجه به انتخاب  $K$  دارای ویژگی‌های منحصر به فردی خواهد بود. در صورتیکه  $k=1$  انتخاب شود، ماتریس  $U$  برداری با درایه‌های  $\pm 1$  خواهد بود. درایه‌ی  $h_1$ ، نشان‌دهنده خوانش مربوط به دنباله‌ی هاپلوتاایپ  $h_1$  است. همچنین درایه‌ی  $h_2$  به معنی خوانش از هاپلوتاایپ  $h_2$  است. در صورتیکه  $k=2$  انتخاب شود، ماتریس  $U$  دارای درایه‌های  $\{0, 1\}$  است که با در نظر گرفتن قواعد خوانش ژنتیکی، در هر سطر یک درایه با مقدار ۱ و دیگر درایه‌های آن صفر است. ماتریس  $V$  دارای ابعاد  $n \times k$  است و مقادیر درایه‌های آن  $\pm 1$  است.

### ۳-۱- مسئله بهینه‌سازی پیشنهادی

به کمک اطلاعاتی که از درایه‌های ماتریس‌های  $U$  و  $V$  عنوان شد، در صورتیکه  $k=1$  انتخاب شود، ماتریس‌های  $U \in \mathbb{R}^{m \times k}$  و  $V \in \mathbb{R}^{n \times k}$  دارای

$$U_{t+1} = U_t - \beta \nabla f(U_t) \quad (10)$$

به کمک تابع جدید بیان شده، مشتقات جزئی تابع هزینه  $f(U, V)$ ،  $\nabla f(U_t, V_t)$  و  $\nabla f(U_t, V_{t+1})$  به صورت زیر تعریف می‌شود:

$$\nabla f(V_t) = -2 \left( P_{\Omega} \left( R - U_t V_t^T \right) \right)^T U_t + \lambda_2 V_t \quad (11)$$

$$\nabla f(U_t) = -2 P_{\Omega} \left( R - U_t V_{t+1}^T \right) V_{t+1} + \lambda_1 U_t \quad (12)$$

حال به کمک رابطه‌های بیان شده (۱۱) و (۱۲) و تابع جدید مطرح شده (۷)، الگوریتم گرادین کاهشی [۴] به کمک شرایط موجود در مسئله پیشنهادی، جدول (۱) بازنویسی می‌شود.

جدول (۱): الگوریتم گرادین کاهشی با تابع پیشنهادی جدید

گام اول	مقدار دهی اولیه $U_0$ و $V_0$ به کمک روش تجزیه مقدار تکین
گام دوم	$\nabla f(V_t) = -2 \left( P_{\Omega} \left( R - U_t V_t^T \right) \right)^T U_t + \lambda_2 V_t$ $V_{t+1} = V_t - \alpha \nabla f(V_t)$ $\nabla f(U_t) = -2 P_{\Omega} \left( R - U_t V_{t+1}^T \right) V_{t+1} + \lambda_1 U_t$ $U_{t+1} = U_t - \beta \nabla f(U_t)$
گام سوم	گرد کردن $V_{t_{last}}$ به $\pm 1$

### ۳-۲- همگرایی مسئله

هدف ما در این بخش محاسبه‌ی ضریب آلفا موجود در رابطه (۹) است تا به کمک آن‌ها الگوریتم گرادین کاهشی جدول (۱) حل شود. تابع هزینه مسئله و گرادین ماتریس  $V$  که قبلاً در رابطه‌ی (۸) و (۱۱) بیان شد

در رابطه‌ی (۱۱) ماتریس  $U$  ثابت فرض شده و از تابع هدف نسبت به ماتریس  $V$  مشتق گرفته می‌شود.

رابطه (۸) برای ماتریس‌های  $(U_t, V_{t+1})$  و  $(U_t, V_t)$  به صورت زیر بازنویسی می‌شود:

$$f(U_t, V_{t+1}) = \left\| P_{\Omega} \left( R - U_t V_{t+1}^T \right) \right\|_F^2 + \lambda_2 \|V_{t+1}\|_F^2 + \lambda_1 \|U_t\|_F^2 \quad (13)$$

$$f(U_t, V_t) = \left\| P_{\Omega} \left( R - U_t V_t^T \right) \right\|_F^2 + \lambda_2 \|V_t\|_F^2 + \lambda_1 \|U_t\|_F^2 \quad (14)$$

به کمک رابطه (۹) برای ماتریس تکرار شونده  $V$  به ساده‌سازی رابطه‌های (۱۳) و (۱۴) می‌پردازیم.

هدف محاسبه  $f(U_t, V_{t+1}) - f(U_t, V_t)$  است که به کمک فرض همگرایی موجود برای مسئله‌ی بهینه‌سازی با تابع هزینه (۸) بتوان مفدار آلفا را به دست آورد. در نتیجه ابتدا رابطه  $\left\| P_{\Omega} \left( R - U_t V_{t+1}^T \right) \right\|_F^2 - \left\| P_{\Omega} \left( R - U_t V_t^T \right) \right\|_F^2$  را بسط می‌دهیم. برای جلوگیری از پیچیدگی محاسبات، تفاضل  $\left\| P_{\Omega} \left( R - U_t V_{t+1}^T \right) \right\|_F^2 - \left\| P_{\Omega} \left( R - U_t V_t^T \right) \right\|_F^2$  را نام گذاری می‌کنیم.

$$Q = \left\| P_{\Omega} \left( R - U_t (V_t - \alpha \nabla f(V_t))^T \right) \right\|_F^2 - \left\| P_{\Omega} \left( R - U_t V_t^T \right) \right\|_F^2 \quad (15)$$

رابطه (۱۵) به کمک رابطه (۹) حاصل شده که به صورت ساده‌تر نیز قابل بیان است:

$$Q = \left\| P_{\Omega} \left( R - U_t V_t^T + \alpha U_t \nabla f(V_t)^T \right) \right\|_F^2 - \left\| P_{\Omega} \left( R - U_t V_t^T \right) \right\|_F^2 \quad (16)$$

به کمک تعریف نرم فروبنیوس، نرم به صورت سیگما بیان می‌شود و محدوده‌ی سیگما به کمک تعریف اپراتور  $P_{\Omega}(\cdot)$  که در رابطه (۳) تعریف شد، به دست می‌آید. در نتیجه رابطه (۱۶) به صورت زیر حاصل می‌شود:

$$Q = 2\alpha \sum_{(i,j) \in \Omega} \left( R - U_t V_t^T \right)_{ij} \left( U_t \nabla f(V_t)^T \right)_{ij} + \alpha^2 \sum_{(i,j) \in \Omega} \left( U_t \nabla f(V_t)^T \right)_{ij}^2 \quad (17)$$

بخش دوم تفاضل (۱۳) و (۱۴) که شامل جملات با ضریب  $\lambda$  است، در زیر بیان شده است:

$$\lambda_2 \|V_{t+1}\|_F^2 - \lambda_2 \|V_t\|_F^2 = \lambda_2 \|V_t - \alpha \nabla f(V_t)\|_F^2 - \lambda_2 \|V_t\|_F^2 \quad (18)$$

رابطه (۱۸) به کمک تعریف نرم فروبنیوس به صورت زیر حاصل می‌شود:

$$\lambda_2 \|V_{t+1}\|_F^2 - \lambda_2 \|V_t\|_F^2 = \lambda_2 \sum_{i,j} \left[ \left( V_t - \alpha \nabla f(V_t) \right)_{ij}^2 \right] - \lambda_2 \|V_t\|_F^2 \quad (19)$$

با طرف راست رابطه (۱۹) می‌توان نوشت:

$$\lambda_2 \|V_{t+1}\|_F^2 - \lambda_2 \|V_t\|_F^2 = \lambda_2 \alpha^2 \sum_{i,j} \left( \nabla f(V_t) \right)_{ij}^2 - 2\alpha \lambda_2 \sum_{i,j} \left( V_t \right)_{ij} \left( \nabla f(V_t) \right)_{ij} \quad (20)$$

به کمک روابط ساده شده‌ی (۱۷) و (۲۰)، رابطه  $f(U_t, V_{t+1}) - f(U_t, V_t)$  به صورت زیر بیان می‌شود:

$$f(U_t, V_{t+1}) - f(U_t, V_t) = \alpha^2 \left\| P_{\Omega} \left( U_t \nabla f(V_t)^T \right) \right\|_F^2 + \lambda_2 \alpha^2 \sum_{i,j} \left( \nabla f(V_t) \right)_{ij}^2 + 2\alpha \sum_{(i,l) \in \Omega} \left( R - U_t V_t^T \right)_{il} \left( U_t \nabla f(V_t)^T \right)_{il} - \lambda_2 \alpha^2 \sum_{i,j} \left( \nabla f(V_t) \right)_{ij}^2 \quad (21)$$

با توجه به اینکه همگرایی مسئله مورد نظر است از نامساوی زیر کمک گرفته می‌شود:

$$f(U_t, V_{t+1}) - f(U_t, V_t) \leq 0 \quad (22)$$

در صورتی که رابطه (۲۱) برابر با صفر قرار گیرد، ریشه‌های معادله به صورت  $\left\{ 0, \frac{\alpha_1}{\alpha_2} \right\}$  حاصل می‌شود که  $\frac{\alpha_1}{\alpha_2}$  به صورت زیر تعریف می‌شود:

$$\alpha_1 = 2[\lambda_2 \sum_{i,j} (V_t)_{ij} (U_t \nabla f(V_t))_{ij} - \sum_{(i,j) \in \Omega} \left( R - U_t V_t^T \right)_{ij} (U_t \nabla f(V_t))_{ij}] \quad (23)$$

$$MEC = \min \left( D(\mathbf{h}, \hat{\mathbf{h}}), D(\mathbf{h}, -\hat{\mathbf{h}}) \right) \quad (32)$$

در رابطه (32)،  $\mathbf{h}$  هاپلوتاپ واقعی و  $\hat{\mathbf{h}}$  نسخه‌ی بازایی شده هاپلوتاپ است. تساوی  $D(\mathbf{h}, \hat{\mathbf{h}}) = \sum_{l=1}^n d(\mathbf{h}_l, \hat{\mathbf{h}}_l)$  در رابطه (32) برقرار است که  $n$  طول هاپلوتاپ است.

خطای تخمین به صورت زیر بیان می‌شود:

$$MEC_r = \frac{MEC}{n} \quad (33)$$

برای ارائه نتایج شبیه سازی تابع هزینه پیشنهادی در تخمین هاپلوتاپ بررسی و نتایج آن با روش مرجع [4] مقایسه می‌شود. روش‌های گرادیان کاهشی بین روش‌های موجود برای تخمین هاپلوتاپ از جمله FAST، SPH، HapCompass و Belief Propagation بهترین جواب را نتیجه می‌دهد. به همین علت روش گرادیان کاهشی به عنوان روش مرجع انتخاب شده است [19، 22، 23].

نتایج به دست آمده کاهش خطای تخمین هاپلوتاپ را نشان می‌دهد. در هر مرحله اجرا، داده‌هایی شبیه سازی‌ها برای داده‌های با سه نرخ نويز  $\{0.1, 0.2, 0.3\}$  و سه نرخ ورودی  $\{0.1, 0.2, 0.3\}$  و طول هاپلوتاپ  $\{50, 100\}$  هاپلوتاپ اجرا شد که در تمام موارد تابع پیشنهادی نسبت به مقاله مرجع کاهش معیار  $MEC_r$  به دست آمد. در ادامه بخشی از این نتایج بیان می‌شود:

جدول (2): مقایسه عملکرد تابع پیشنهادی با مرجع در نرخ ورودی

0/2

مرجع [4]	پیشنهادی $\lambda = 2$		پیشنهادی $\lambda = 10$	
	$MEC_r$	درصد بهبود نسبت به مرجع	$MEC_r$	درصد بهبود نسبت به مرجع
نرخ نويز 0.1	0.0024	0.0019	0.0021	8
0.2	0.0303	0.0296	0.030	0.99
0.3	0.2197	0.210	0.2090	4.87

جدول (3): مقایسه عملکرد تابع پیشنهادی با مرجع در نرخ نويز 0/2

مرجع [4]	پیشنهادی $\lambda = 2$		پیشنهادی $\lambda = 10$	
	$MEC_r$	درصد بهبود نسبت به مرجع	$MEC_r$	درصد بهبود نسبت به مرجع
ورودی 0.1	0.1692	0.149	0.160	5.43
0.2	0.213	0.198	0.207	2.81
0.3	0.0095	0.0086	0.0093	2.10

$$\alpha_2 = \lambda_2 \left\| \nabla f(\mathbf{V}_t) \right\|_F^2 + \left\| P_{\Omega} \left( \mathbf{U}_t \nabla f(\mathbf{V}_t)^T \right) \right\|_F^2 \quad (24)$$

به کمک ریشه‌های به دست آمده و نامساوی (22) که مربوط به همگرایی مسئله است، می‌توان نتیجه گرفت که اگر آلفا بین دو ریشه انتخاب شود، در نامساوی (22) صدق می‌کند.

در نتیجه با توجه به بازه‌های انتخابی برای مقدار ثابت  $C \in (0, 1)$ ، ضریب آلفا به صورت زیر به دست می‌آید:

$$\alpha = C \frac{\alpha_{\text{numerator}}}{\alpha_{\text{denominator}}} \quad (25)$$

$$\alpha_{\text{numerator}} = 2 \left[ \lambda_2 \sum_{i,j} (\nabla f(\mathbf{V}_t))_{ij} - \sum_{(i,j) \in \Omega} \left( \mathbf{R} - \mathbf{U}_t \mathbf{V}_t^T \right)_{ij} \left( \mathbf{U}_t \nabla f(\mathbf{V}_t)^T \right)_{ij} \right]$$

$$\alpha_{\text{denominator}} = \lambda_2 \left\| \nabla f(\mathbf{V}_t) \right\|_F^2 + \left\| P_{\Omega} \left( \mathbf{U}_t \nabla f(\mathbf{V}_t)^T \right) \right\|_F^2 \quad (27)$$

آلفای به دست آمده در رابطه (25) برای جایگذاری در رابطه (9) استفاده می‌شود تا به کمک الگوریتم موجود در جدول (1) به تخمین هاپلوتاپ کمک کند. محاسبه ضریب بتا هم مشابه ضریب آلفا انجام خواهد گرفت و رابطه نهایی به صورت زیر بیان می‌شود:

$$\beta = C \frac{\beta_{\text{numerator}}}{\beta_{\text{denominator}}} \quad (28)$$

$$\beta_{\text{numerator}} = 2 \left[ \lambda_1 \sum_{i,j} (\nabla f(\mathbf{U}_t, \mathbf{V}_{t+1}))_{ij} - \sum_{(i,j) \in \Omega} \left( \mathbf{R} - \mathbf{U}_t \mathbf{V}_{t+1}^T \right)_{ij} \left( \nabla f(\mathbf{U}_t, \mathbf{V}_{t+1}) \mathbf{V}_{t+1}^T \right)_{ij} \right] \quad (29)$$

$$\beta_{\text{denominator}} = \lambda_1 \sum_{i,j} \left( \nabla f(\mathbf{U}_t, \mathbf{V}_{t+1}) \right)_{ij}^2 + \sum_{(i,j) \in \Omega} \left( \nabla f(\mathbf{U}_t, \mathbf{V}_{t+1}) \mathbf{V}_{t+1}^T \right)_{ij}^2 \quad (30)$$

مقدار ثابت  $C \in (0, 1)$  برای شبیه سازی‌ها اعمال می‌شود.

#### 4- نتایج شبیه سازی

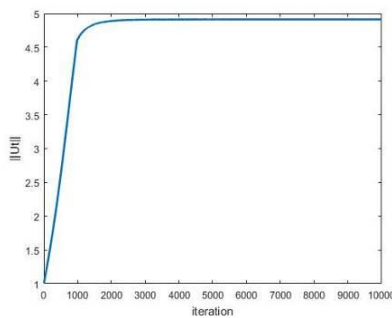
در این بخش در ابتدا معیار ارزیابی نتایج ارائه شده است. روشی عملکرد بهتری دارد که معیارها خطای کمتری را در بازایی داده‌ها نتیجه دهد. در این مقاله از معیار خطای تخمین و کمینه خطای تصحیح استفاده می‌شود.

در ابتدا فاصله همینگ<sup>28</sup> به صورت زیر تعریف می‌گردد:

$$d(a, b) = \begin{cases} 1 & a \neq b \\ 0 & a = b \end{cases} \quad (31)$$

کمینه خطای تصحیح<sup>29</sup> که معیار عملکرد روش پیشنهادی است به صورت زیر تعریف می‌شود:





شکل (۲): همگرایی نرم  $U$  برای نرخ نویز  $0.1$  و نرخ ورودی  $0.3$

## ۵- نتیجه گیری

الگوریتم گردادیان کاهشی یکی از روش‌های تخمین هاپلوتاایپ به کمک تجزیه‌ی ماتریس است. در روش‌های تخمین گزارش شده، از اطلاعات درایه‌های ماتریس استفاده نشده‌است. به همین علت در این مقاله با گذاشتن شرط بر روی میزان بزرگی مقادیر درایه‌های ماتریس  $U$  و  $V$ ، تابع هزینه جدید پیشنهاد شد. تابع پیشنهادی مانع از تخمین ماتریس‌های پرت و نامناسب  $U$  و  $V$  شد که باعث تخمین دقیق‌تر هاپلوتاایپ شد. نتایج حاصل از شبیه‌سازی، نشان‌دهنده‌ی کاهش خطای تخمین در تابع هزینه پیشنهادی نسبت به روش‌های پیشین است. این نتایج در داده‌های متفاوت ماتریس خوانش برقرار است که برتری تابع پیشنهادی را نشان می‌دهد.

## مراجع

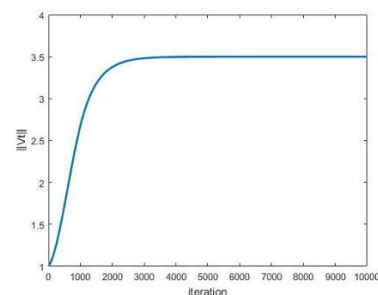
- [۱] فرازکیش راضیه، مدل‌سازی قابلیت اطمینان در نانوربات‌های زیستی. مجله مهندسی برق و الکترونیک ایران. ۱۳۹۹؛ ۱۷ (۳): ۱۶-۱۱.
- [۲] آتشبار محمود، کهائی محمدحسین. جهت یابی چند گونه‌ده با استفاده از روش WCSSDOA. مجله مهندسی برق و الکترونیک ایران. ۱۳۹۵؛ ۱۳ (۲): ۶۱-۷۴.
- [۳] مجیدیان سینا، حدادی فرزانه. تخمین جهت منابع با استفاده از زیرفضای ختری-راو. مجله مهندسی برق و الکترونیک ایران. ۱۳۹۶؛ ۱۴ (۲): ۳۷-۴۷.
- [4] C. Cai, S. Sanghavi, and H. Vikalo, "Structured low-rank matrix factorization for haplotype assembly," IEEE Journal of Selected Topics in Signal Processing, vol. 10, pp. 647-657, 2016.
- [5] Z. Puljiz and H. Vikalo, "A message passing algorithm for haplotype assembly," in 2013 Asilomar Conference on Signals, Systems and Computers, 2013, pp. 1726-1729.
- [6] S. Das and H. Vikalo, "SDhaP: haplotype assembly for diploids and polyploids via semi-definite programming," BMC genomics, vol. 16, p. 260, 2015.
- [7] T. Inbamalar and R. Sivakumar, "Study of DNA sequence analysis using DSP techniques," J. Autom. Control Eng., vol. 1, 2013.

در مراجع روش‌های مختلفی برای انتخاب ضریب  $\lambda$  بیان می‌شود که در آن‌ها دلیل و اثبات ریاضی برای انتخاب این مقدار وجود ندارد. انتخاب  $\lambda$  می‌تواند به کمک نتایج حاصل از شبیه‌سازی انجام گیرد و یا به کمک اثبات‌های دیگر موجود در مسئله، شرایطی برای مقدار دهی در نظر گرفت. ما نیز در این مقاله به کمک نتایج حاصل از شبیه‌سازی مقادیر  $0.2$  و  $0.1$  برای ضریب  $\lambda$  انتخاب کردیم [۲۰، ۲۴-۲۷].

در این شبیه‌سازی  $m=n=50$  انتخاب شده است و تابع برای دو مقدار  $\lambda$  پیشنهادی اجرا و با مقاله مرجع مقایسه می‌شود. نویز اعمال شده بر اساس تعریف نویز در مسئله‌ی هاپلوتاایپ به صورت رابطه (۳) می‌باشد و نرخ نویز نسبت درایه‌های نویزی در  $R$  به کل درایه‌های دارای مقدار  $R$  است. نرخ ورودی از نسبت تعداد درایه‌های دارای مقدار در  $R$  نسبت به کل درایه‌های آن تعریف می‌شود. همانطور که در جدول‌های (۲) و (۳) مشخص است به دلیل وجود ترم‌های مرتبط با ضرایب  $\lambda_1$  و  $\lambda_2$  کاهش خطای تخمین هاپلوتاایپ و در نتیجه تخمین هاپلوتاایپ با دقت بیشتری انتظار می‌رفت. بنابراین این نتایج صحت تابع پیشنهادی (۸) را تایید می‌کند. زیرا معیار عملکرد رابطه (۳۳) در جدول‌های (۲) و (۳) کاهش میزان  $MEC_r$  را نشان می‌دهد. یعنی وجود ضرایب  $\lambda_1$  و  $\lambda_2$  باعث شده است که برای تخمین هاپلوتاایپ از انتخاب داده‌های پرت و نامناسب برای ماتریس‌های  $U$  و  $V$  جلوگیری شود.

با توجه به برتری روش بیان شده نسبت به گردادیان کاهشی و یکسان بودن نوع داده‌ها در شبیه‌سازی‌ها می‌توان نتیجه گرفت که روش پیشنهادی از دیگر الگوریتم‌های مورد استفاده برای تخمین هاپلوتاایپ عنوان شده در ابتدا این بخش هم برتری دارد.

در روش پیشنهادی همگرایی مسئله هم مورد ارزیابی قرار گرفته است که نتایج حاصل در شکل‌های (۱) و (۲) قابل مشاهده است. همگرایی مسئله برای طول هاپلوتاایپ  $50$ ، نرخ ورودی  $0.3$ ، نرخ نویز  $0.1$  و ضریب  $\lambda$   $0.1$  اجرا شده است. محور افقی تعداد تکرارها و محور عمودی مربوط به نرم  $U$  و  $V$  است. همگرایی مسئله در مثال بیان شده قابل مشاهده است.



شکل (۱): همگرایی نرم  $V$  برای نرخ نویز  $0.1$  و نرخ ورودی  $0.3$

- [26] F. Shang, J. Cheng, Y. Liu, Z.-Q. Luo, and Z. Lin, "Bilinear factor matrix norm minimization for robust PCA: Algorithms and applications," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, pp. 2066-2080, 2017.
- [27] R. Cabral, F. De la Torre, J. P. Costeira, and A. Bernardino, "Unifying nuclear norm and bilinear factorization approaches for low-rank matrix decomposition," in *Proceedings of the IEEE International Conference on Computer Vision*, 2013, pp. 2488-2495.
- زیر نویس ها**
- 
- <sup>1</sup>Deoxyribonucleic Acid  
<sup>2</sup>Semi-Definite Programming  
<sup>3</sup>Digital Signal Processing (DSP)  
<sup>4</sup>Discrete Fourier Transform (DFT)  
<sup>5</sup>Discrete Wavelet Transform (DWT)  
<sup>6</sup>Atomic Decomposition for Minimum Rank Approximation  
<sup>7</sup>Subspace Evolution and Transfer  
<sup>8</sup>Alternating Projection  
<sup>9</sup>Ribonucleic Acid (RNA)  
<sup>10</sup>Deoxyribonucleic Acid(DNA)  
<sup>11</sup>Adenin(A)  
<sup>12</sup>Thyamine(T)  
<sup>13</sup>Cytosine(C)  
<sup>14</sup>Guanine(G)  
<sup>15</sup>Chromosomes  
<sup>16</sup>Homologous  
<sup>17</sup>Diploid  
<sup>18</sup>Single Nucleotide Polymorphisms (SNP)  
<sup>19</sup>AllanMaxam  
<sup>20</sup>WalterGilbert  
<sup>21</sup>Frederick Sanger  
<sup>22</sup>Shotgun Sequencing  
<sup>23</sup>Read  
<sup>24</sup>Overlap  
<sup>25</sup>Paired-end  
<sup>26</sup>Voss  
<sup>27</sup>Electron-Ion Interaction Potential  
<sup>28</sup>Hamming distance  
<sup>29</sup>Minimum Error Correction (MEC)
- [8] X. Jiang, Z. Zhong, X. Liu, and H. C. So, "Robust matrix completion via alternating projection," *IEEE Signal Processing Letters*, vol. 24, pp. 579-583, 2017.
- [9] H. C. So and W.-J. Zeng, "Outlier-Robust Matrix Completion via lp-Minimization," 2018.
- [10] X. P. Li, L. Huang, H. C. So, and B. Zhao, "A survey on matrix completion: perspective of signal processing," *arXiv preprint arXiv:1901.10885*, 2019.
- [11] P. Jain, P. Netrapalli, and S. Sanghavi, "Low-rank matrix completion using alternating minimization," in *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, 2013, pp. 665-674.
- [12] J. Yu, G. Zhou, C. Li, Q. Zhao, and S. Xie, "Low Tensor-Ring Rank Completion by Parallel Matrix Factorization," *IEEE transactions on neural networks and learning systems*, 2020.
- [13] Q. Wang, X. He, X. Jiang, and X. Li, "Robust Bi-stochastic Graph Regularized Matrix Factorization for Data Clustering," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [14] S. Barik and H. Vikalo, "Matrix Completion and Performance Guarantees for Single Individual Haplotyping," *IEEE Transactions on Signal Processing*, vol. 67, pp. 4782-4794, 2019.
- [15] R. Rizzi, V. Bafna, S. Istrail, and G. Lancia, "Practical algorithms and fixed-parameter tractability for the single individual SNP haplotyping problem," in *International Workshop on Algorithms in Bioinformatics*, 2002, pp. 29-43.
- [16] A. M. Maxam and W. Gilbert, "A new method for sequencing DNA," *Proceedings of the National Academy of Sciences*, vol. 74, pp. 560-564, 1977.
- [17] N. M. Haan and S. J. Godsill, "Bayesian models for DNA sequencing," in *2002 IEEE International Conference on Acoustics, Speech, and Signal Processing*, 2002, pp. IV-4020-IV-4023.
- [18] A. S. Motahari, G. Bresler, and N. David, "Information theory of DNA shotgun sequencing," *IEEE Transactions on Information Theory*, vol. 59, pp. 6273-6289, 2013.
- [19] D. Aguiar and S. Istrail, "HapCompass: a fast cycle basis algorithm for accurate haplotype assembly of sequence data," *Journal of Computational Biology*, vol. 19, pp. 577-590, 2012.
- [20] Z. Zhou, X. Li, J. Wright, E. Candes, and Y. Ma, "Stable principal component pursuit," in *2010 IEEE international symposium on information theory*, 2010, pp. 1518-1522.
- [21] H. Si, H. Vikalo, and S. Vishwanath, "Information-theoretic analysis of haplotype assembly," *IEEE Transactions on Information Theory*, vol. 63, pp. 3468-3479, 2017.
- [22] F. Geraci, "A comparison of several algorithms for the single individual SNP haplotyping reconstruction problem," *Bioinformatics*, vol. 26, pp. 2217-2225, 2010.
- [23] Z. Puljiz and H. Vikalo, "Decoding genetic variations: Communications-inspired haplotype assembly," *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, vol. 13, pp. 518-530, 2010.
- [24] E. J. Candès, X. Li, Y. Ma, and J. Wright, "Robust principal component analysis?," *Journal of the ACM (JACM)*, vol. 58, pp. 1-37, 2011.
- [25] M. Mardani, G. Mateos, and G. B. Giannakis, "Decentralized sparsity-regularized rank minimization: Algorithms and applications," *IEEE Transactions on Signal Processing*, vol. 61, pp. 5374-5388, 2013.