

# فصل ۱

## روش پیشنهادی

### ۱.۱ مقدمه

پس از آشنایی با روش‌های پیشین که برای حل مسئله مشابه مورد استفاده قرار گرفته‌اند، حال می‌توانیم به معرفی و تشریح روش‌های پیشنهادی خود برای حل مسئله پیش رو بپردازیم. در این فصل ابتدا داده‌های ورودی مسئله را همراه با فرضیات در نظر گرفته شده بیان می‌کنیم و پس از آن دو روش پیشنهادی متفاوت را بیان خواهیم نمود. در روش اول که به رویکردهای پیشین نزدیک‌تر است با تغییری از جنس روش‌های نوین در مراحل میانی به یک روش جدید می‌رسیم که به علت افزایش سرعت همگرایی می‌توان فرض و داده‌های جدیدی را از طریق حذف و تغییر تعداد کپی<sup>۱</sup> به آن افزود و پاسخ گرفت. اما روش دوم کاملاً متفاوت بوده و با رویکردی جدید در حوزه یادگیری ماشین همراه است که به کمک یادگیری تقویتی به حل مسئله مورد نظر می‌پردازد.

### ۲.۱ معرفی دادگان ورودی

قبل از وارد شدن به بخش روش‌های پیشنهادی نیاز است تا دادگان ورودی را مشخص و معرفی نماییم. دادگان ورودی در این پایان‌نامه همگی به صورت فایل‌های خام اسکی<sup>۲</sup> هستند که حاوی اطلاعات جهش‌های ماتریس

---

<sup>۱</sup>Copy number variation (CNV)

<sup>۲</sup>Ascii

ژن-سلول (SNV) و اطلاعات مربوط به حذف و تغییر تعداد کپی هستند.

در ادامه جدول ۱.۱ را برای معرفی اندیس‌های بکار گرفته شده در روابط مربوط به روش پیشنهادی اول معرفی می‌نماییم.

جدول ۱.۱: اندیس‌های به کار رفته در روابط روش پیشنهادی اول

$D$	ماتریس داده نویزی در دسترس که مقادیر ۰ و ۱ در آن قرار دارد
$E$	ماتریس داده حقیقی بدون نویز که به دنبال آن هستیم
$T$	درخت فیلوژنی جهش‌ها
$\sigma$	بردار انتصابات
$\wp$	بردار پذیرش فقدان
$X_T$	ماتریس متناظر درخت $T$
$N$	تعداد سلول‌های نمونه
$M$	تعداد جهش‌ها
$\mathcal{N}$	مجموعه سلول‌های متمایز از هم
$\mathcal{M}$	مجموعه جهش‌های متمایز از هم
$\alpha$	نرخ خطای مثبت کاذب
$\beta$	نرخ خطای منفی کاذب

### ۳.۱ روش پیشنهادی برای مدیریت داده‌های از دست رفته

در ادامه این بخش به معرفی روش‌های پیشنهادی پرداخته خواهد شد اما در ابتدا به دلیل وجود داده‌های از دست رفته در پایگاه‌داده‌های مورد استفاده لازم است تا به بررسی و ارائه رویکردی برای حل این مشکل پرداخته شود و در ادامه پس از معرفی روش پیشنهادی برای مدیریت این داده‌های از دست رفته، هر کدام از روش‌های پیشنهادی به تفصیل شرح داده شود.

همان‌گونه که در داده‌های حقیقی مشاهده شد در پایگاه داده‌های حقیقی ما با اطلاعات از دست رفته مواجه هستیم و به همین دلیل نیز سعی کردیم تا در پایگاه داده مجازی تولید شده نیز به مشابه داده‌های حقیقی، شامل اطلاعات از دست رفته باشد. در این بخش به رویکرد روش محاسبه استاتیک برای مدیریت این داده‌های از دست رفته می‌پردازیم و در بخش بعد به معرفی روشی برای بدست آوردن درخت فیلوژنی پرداخته خواهد شد. همان‌گونه

که در ادامه بررسی خواهد شد، این اطلاعات از دست رفته در پایگاه داده‌های مختلف نرخ‌های متفاوتی دارد که تاثیر این تغییرات نیز در روشی پیشنهادی بررسی خواهد شد.

### ۱.۳.۱ روش محاسبه استاتیک

در این روش قصد داریم تا به یک‌باره بتوانیم مقادیر مناسب برای داده‌هایی که از دست رفته‌اند را تخمین بزنیم. در این روش باید توجه شود که ما لزوماً به دنبال جایگذاری مقدار از دست رفته با مقدار درست واقعی نیستیم. اگرچه چنین بیانی در نگاه اول ممکن است تعجب‌آور باشد اما با دقت بیشتر متوجه خواهیم شد که ما در آینده برای خطاهای موجود در پایگاه داده مدل‌سازی‌های محدودی داریم. مدل‌هایی که بهترین آن‌ها نیز ممکن است با واقعیت نویز افزوده شده به دادگان متفاوت باشد. در نتیجه اگر مطمئن بودیم که تمام داده‌هایی که موجود می‌باشند بدون خطا هستند در آن صورت ما نیز به دنبال یافتن جایگذاری با مقدار واقعی بودیم اما در حال حاضر که درصدی از داده‌های در دسترس خود همراه با خطا می‌باشند، ما به دنبال جایگذاری ای هستیم که بتواند در مجموع با مدل‌سازی خطایی که در نظر می‌گیریم بیشترین سازگاری را داشته باشد کما اینکه ممکن است در حقیقت جایگزاری اشتباهی انجام داده باشیم. حال با توجه به توضیحی که بیان شد به تشریح این روش می‌پردازیم.

با توجه به فرض مدل مکان‌های بی‌نهایت می‌دانیم که جهش‌های اتفاق افتاده در والد در تمامی نسل‌های آینده باقی خواهد ماند. بنابراین اگر تمامی جهش‌های نمونه (سلول)  $a$  در نمونه‌ای دیگر مانند  $b$  قرار داشته باشد، بنابراین می‌توان نتیجه گرفت که  $a$  یکی از اجداد  $b$  خواهد بود. همین فرضیه هسته اصلی روش پیشنهادی در نظر گرفته شده را تشکیل می‌دهد. بنابراین اگر جهش  $i$  در سلول  $a$  از دست رفته است، با توجه به اینکه آن جهش در سلول  $b$  چه وضعیتی دارد می‌توان تصمیم‌گیری کرد. اگر  $b(i) = 0$  باشد، در این صورت  $a(i)$  حتماً باید ۰ باشد وگرنه فرض اولیه مدل مکان‌های بی‌نهایت نقض خواهد شد. اما اگر  $b(i) = 1$  باشد، آنگاه نتیجه خاصی نمی‌توان گرفت و باید به دنبال نمونه والد  $a$  یعنی نمونه  $d$  باشیم. حال اگر  $d(i) = 1$  باشد، آنگاه  $a(i)$  حتماً باید ۱ باشد. اما اگر  $d(i) = 0$  بود آنگاه انتخاب هر مقداری برای  $a(i)$  تقریباً آزاد خواهد بود زیرا با فرض اولیه تناقضی ندارد و اینکه ساختار فیلوژنی را تغییر نمی‌دهد. اما از آنجایی که خود داده‌های در دسترس شامل خطا می‌باشند و هر نمونه‌ای که حاوی اطلاعات از دست رفته است لزوماً یک نواده یا یک والد ندارد، مجموعه‌ای از سلول‌های فرزند یا والد خواهند بود که متناسب با پارمترهای خطایی که در نظر می‌گیریم و فاصله ژنی‌ای که دارند می‌توانند در

تصمیم‌گیری تاثیرگذار باشند. صورت دقیق‌تر توضیحات داده شده را می‌توان به صورت فرمولی که در ادامه آمده است به نمایش درآورد.

در ابتدا تابعی به نام  $F_s(D_{ij})$  تعریف می‌کنیم که به نوعی با توجه به ارزشی که به سلول‌های نواده شده از سلول  $j$  می‌دهد سعی دارد تا اطمینان به بودن داده از دست رفته  $D_{ij}$  را بیان کند. برای محاسبه این تابع می‌دانیم که ابتدا سلول‌های مختلف با توجه به احتمال نواده بودنشان باید رتبه‌بندی شوند و وزن بگیرند. پس از آن هر سلول متناسب با ارزش تاثیرگذاری خود می‌تواند در مورد جایگاه جهش  $i$  برای سلول  $j$  نظر دهد.

$$F_s(D_{ij}) = \sum_{n \in \mathcal{N}} (1 - D_{mj}) \prod_{m=1}^M W(D_{mn}, D_{mj}) \quad (1.1)$$

در فرمول ۱.۱ مجموعه  $\mathcal{N}$  برابر با مجموعه سلول‌های متمایز از هم است. زیرا که در بسیاری از پایگاه‌داده‌ها از یک نمونه سلول ممکن است چندین نمونه وجود داشته باشد که وجود آن‌ها باعث بایس در محاسبات ما خواهد شد. همچنین تابع  $W_s(c, p)$  به ارزش‌دهی جهش  $c$  در برابر  $p$  به عنوان نواده بودن می‌پردازد که در فرمول ۲.۱ تعریف شده است.

$$W(c, p) = \begin{cases} 1 & \text{if } c = 1, p = 1 \\ 1 - \xi & \text{if } c = 1, p = 0 \\ 0 & \text{if } c = 0, p = 1 \\ 1 & \text{if } c = 0, p = 0 \end{cases} \quad (2.1)$$

مقدار  $\xi$  عددی بین  $(0, 1)$  است که پارامتری در جهت میزان ارزش‌دهی به نوادگان با فواصل مختلف می‌باشد. هرچه این عدد بزرگتر باشد به معنی کم ارزش‌تر شدن نوادگان با فواصل بیشتر است و بالعکس. به همین صورت برای اولاد سلول  $j$  نیز می‌توان مشابه حالت قبل عمل کرد که روابط آن به صورت فرمول ۳.۱ خواهد شد.

$$F_a(D_{ij}) = \sum_{n \in \mathcal{N}} D_{mj} \prod_{m=1}^M W(D_{mj}, D_{mn}) \quad (3.1)$$

حال دو نکته در استفاده از روابط بالا باقی خواهد ماند.

نکته اول وجود داده‌های دیگر از دست رفته در محاسبه توابع است که به دو صورت می‌توان با آن‌ها برخورد نمود. رویکرد اول این است که در آن جایگاه ژنی از محاسبه آن خود داری شود و رویکرد دوم استفاده از مقدار ۰/۵ یا فراوانی نسبی آن جهش در محاسبات است که ما رویکرد اول را در این گزارش استفاده خواهیم کرد.

نکته دوم وجود خطا در داده‌هاست. برای مدیریت این مشکل می‌توان با مدل‌سازی خطا که به صورت فرمول ۴.۱ بیان می‌شود، برخورد کرد.

$$\begin{aligned} P(D_{ij} = 1 | E_{ij} = 0) &= \alpha, & P(D_{ij} = 0 | E_{ij} = 0) &= 1 - \alpha \\ P(D_{ij} = 0 | E_{ij} = 1) &= \beta, & P(D_{ij} = 1 | E_{ij} = 1) &= 1 - \beta \end{aligned} \quad (4.1)$$

پس از تعریف مدل‌سازی خطا می‌توان روابط قبلی را مجدداً به صورتی که در ادامه آمده است بازنویسی کرد.

$$W_e(c, p) = \sum_{i, j \in \{0, 1\}} P(c | E_c = i) P(p | E_p = j) W(i, j) \quad (5.1)$$

که در این صورت توابع  $F_a$  و  $F_p$  نیز به صورت زیر همراه با مدل‌سازی خطا بازتعریف خواهند شد.

$$\begin{aligned} \hat{F}_s(D_{ij}) &= \sum_{n \in \mathcal{N}} [1 - D_{mj}(1 - \alpha)] \prod_{m=1}^M W_e(D_{mn}, D_{mj}) \\ \hat{F}_a(D_{ij}) &= \sum_{n \in \mathcal{N}} D_{mj}(1 - \beta) \prod_{m=1}^M W_e(D_{mj}, D_{mn}) \end{aligned} \quad (6.1)$$

حال پس از محاسبه مقادیر  $\hat{F}_a$  و  $\hat{F}_s$  می‌توان در مورد داده نامعلوم  $D_{ij}$  به صورت فرمول ۷.۱ تصمیم گرفت.

$$D_{ij} = \begin{cases} 0 & \text{if } \hat{F}_s \geq \hat{F}_a \\ 1 & \text{if } \hat{F}_s < \hat{F}_a \end{cases} \quad (7.1)$$

همچنین با کمی دقت در فرمول‌بندی انجام شده اگر برای تمام  $i, j$ ‌های ماتریس  $D$  این مقادیر توابع  $\hat{F}$  محاسبه شوند، خود می‌توانند معیاری برای ارزیابی پایگاه‌داده در دسترس و احتمال درستی فرض مدل مکان‌های بی‌نهایت

باشند.

### ۱.۱.۳.۱ تصادفی

پرو کردن کاملاً تصادفی میس‌ها. در این روش به صورت تصادفی مقادیر از دست رفته را مقدار دهی می‌کنیم. تنها نکته‌ای که در این روش وجود دارد این است که نباید این پرو کردن تصادفی داده‌های از دست رفته باعث شود تا پارامترهای مدل‌سازی‌ای که از قبل در نظر گرفته بودیم با این روش نادقیق شوند.

## ۴.۱ روش پیشنهادی اول (درخت‌بازی)

در این روش ما بر حسب بهتر کردن یک پاسخی که از پیش داشتیم به دنبال رسیدن به بهترین پاسخ ممکن در طی تکرار<sup>۳</sup> پشت سر هم هستیم. برای مشخص شدن نحوه کارکرد روش پیشنهادی در مراحل‌ی که در ادامه بیان خواهد شد به عنوان مثال یک ماتریس

$$D = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 \end{bmatrix} \quad (۸.۱)$$

را به عنوان ورودی مساله به همراه پارامترهای  $\alpha$  و  $\beta$  در نظر بگیرید. (برای راحتی کار فرض کرده‌ایم که داده از دست رفته در  $D$  نداریم.)

<sup>3</sup>Iteration

## ۱.۴.۱ پیش پردازش

قبل از شروع باید بر روی داده‌ها یک پیش پردازش اعمال کنیم که وابسته به سیاست در نظر گرفته شده می‌تواند باعث تغییر در پاسخ نهایی نیز شود. به این منظور داده‌هایی که miss شده‌اند با روش‌های زیر می‌تواند برای ورود به مرحله بعد تخمین زده شود. به صورت تصادفی

## ۲.۴.۱ اولین پاسخ (درخت تصادفی)

همان گونه که از قبل می‌دانستیم خروجی نهایی ما برابر با درختی خواهد بود که نودهای آن برابر با جهش‌های ماتریس ورودی ما و برگ‌های آن برابر با نمونه‌های مشاهده شده خواهند بود. در روش پیشنهادی اول ما به دنبال بهتر کردن این درخت به عنوان پاسخ هستیم. از این رو پایه این روش پیشنهادی اول بر مبنای بهتر کردن پاسخ فعلی بنا نهاده شده است. در نتیجه ما همواره پاسخی به عنوان جواب نهایی داریم که تلاش خواهیم نمود تا با استفاده از ابزارهایی بتوانیم ابا ایجاد تغییری در این پاسخ به پاسخی جدید برسیم که قابل مقایسه با پاسخ فعلی برای انجام مراحل بعدی باشد.

با توجه به توضیحاتی که داده شد ما برای شروع الگوریتم پیشنهادی اول خود نیاز به یک پاسخ داریم. این پاسخ که درخت فیلوژنی هست با توجه پارامترهای ورودی و انتخاب یک نود ( $root$  به عنوان ریشه این درخت به صورت زیر حاصل می‌شود.

$$\mathcal{M} = \{1 \dots M\} \quad (9.1)$$

$$\hat{B}_{T_1} = [R_1(\mathcal{M} - |1|), R_2(\mathcal{M} - |2|), \dots, R_{root}(\{\}), \dots, R_M(\mathcal{M} - |M|)]$$

که در این رابطه  $\mathcal{M}$  برابر با مجموعه تمامی جهش‌های متمایز از شماره ۱ تا  $M$  است و  $\hat{B}$  مشخص کننده نود پدر در درخت برای جهش  $i$ ام در این لیست خود است که توسط تابع  $R_i(X)$  به صورت کاملاً یکنواخت<sup>۴</sup> از اعضای مجموعه  $X$  انتخاب می‌شود.

با توجه به مثالی که در رابطه ۸.۱ زده شد فرض کنید مقدار بردار  $\hat{B}$  با ریشه ۲  $root = 2$  به صورت رابطه ۱۰.۱

<sup>۴</sup>Uniform

شود.

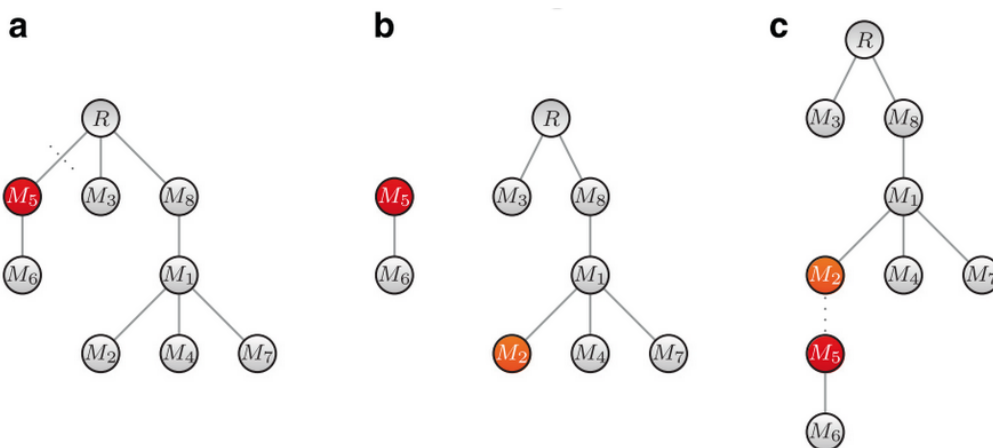
$$\hat{B} = [2, 1, -, 0] \quad (10.1)$$

که درخت شکل ۱.۱ را نتیجه می‌دهد.

شکل ۱.۱: درخت تصادفی اول

### ۳.۴.۱ پاسخ‌ی جدید

تا به اینجا ما یک درخت فیلوژنی به عنوان پاسخ داریم که در این بخش می‌خواهیم با انجام تغییراتی بر روی آن به یک پاسخ جدید برسیم تا در گام‌های بعدی بتوانیم با مقایسه آن‌ها تصمیمات لازم را برای ادامه الگوریتم بگیریم. به همین منظور تقریباً مشابه با روش [۱] به صورت هرس و اتصال دوباره<sup>۵</sup> قصد داریم تا درخت پاسخ فعلی را برای رسیدن به یک پاسخ دیگر تغییر دهیم. در شکل ۲.۱ مثالی از این روش آورده شده است. در این



شکل ۲.۱: نحوه انجام کار روش هرس و اتصال دوباره [۱].

شکل مطابق با قسمت a یکی از نودهای درخت (به جز ریشه) انتخاب می‌شود (در اینجا نود  $M_5$ ) و اتصال از پدرش قطع می‌شود. در این حالت به دو درخت مشابه با شکل b می‌رسیم. حال در درخت باقی مانده یک نود

<sup>5</sup>Prune and reattach



دیگر ( $M_2$ ) به عنوان پدری جدید انتخاب می‌شود تا با این تغییر به درخت جدید شکل C برسیم. در نتیجه با این تکنیک می‌توان به پاسخ‌های جدید رسید اما سوالی که باقی ماند این است که این دو نود چگونه باید انتخاب شوند؟ این دو انتخاب به صورت هوشمند توسط دو شبکه عمیق گرفته می‌شود. این شبکه‌ها که در اصل شبکه‌های یادگیری تقویتی<sup>۶</sup> عمیق هستند، از پیش برای این منظور آموزش داده شده‌اند. شبکه اول برای انتخاب محل برش (هرس) مورد استفاده قرار می‌گیرد و شبکه دوم با دریافت خروجی xxxxxx دو درخت راهبر-پیرو<sup>۷</sup> مکان‌های مناسب برای بازاتصال را ارزش‌گذاری می‌کند. این دو شبکه در قسمت‌های ۶.۴.۱ و ۷.۴.۱ به تفصیل شرح داده شده‌اند.

## ۴.۴.۱ مقایسه و ارزیابی پاسخ‌ها

پس از اینکه از پاسخ فعلی به یک پاسخ جدید رسیدیم حال می‌توان کیفیت این دو پاسخ را باهم مقایسه کرد و پس از آن با توجه به امتیاز دو پاسخ در مورد پذیرش یا عدم پذیرش پاسخ جدید در برابر پاسخ فعلی تصمیم گرفت. این فرآیند شامل دو بخش اصلی است که در دو زیربخشی که در ادامه آمده است بیان شده‌اند.

### ۱.۴.۴.۱ تبدیل درخت پاسخ به ماتریس

همان‌گونه که می‌دانیم

### ۲.۴.۴.۱ مقایسه پاسخ فعلی با پاسخ آرمانی

پس از استخراج ماتریس مناسب از درخت می‌توان به ارزش‌گذاری و محاسبه درست‌نمایی<sup>۸</sup> پرداخت. این عمل به صورت رابطه ۱۱.۱ محاسبه می‌شود.

$$L : P(D|T, \sigma, \varphi, \theta) = \prod_{n=1}^{n=N} \prod_{m=1}^{m=M} P(D_{nm}|X_{nm}) \quad (11.1)$$

<sup>۶</sup>Reinforcement learning

<sup>۷</sup>Master-slave

<sup>۸</sup>Likelihood

که  $X$  برابر ماتریس بدست آمده از درخت  $T$  با توجه به بردارهای  $\sigma$  و  $\wp$  است. این رابطه بیانگر احتمال مشاهده ماتریس داده ورودی  $D$  در صورتی است که درخت فیلوژنی صحیح  $T$  و پارامترهای حقیقی  $\theta$  باشد که توسط بردارهای  $\sigma$  و  $\wp$  تثبیت شده است. هرچه این احتمال بالاتر باشد نمایانگر این است که درخت، پارامترها و بردارهای کنترلی ما بگونه‌ای انتخاب شده‌اند که محتمل‌ترین حالت برای مشاهده داده‌های ورودی ما هست و در این صورت بهترین پاسخ برای ما همان پاسخی خواهد بود که محتمل‌ترین باشد. از این رو با دانستن  $\theta$  ما به دنبال  $T$ ‌ای به همراه بردارهای مربوطه آن هستیم که پاسخ رابطه باشد.

$$(T, \sigma, \wp)_{\text{ML}} = \arg \max_{(T, \sigma, \wp)} P(D|T, \sigma, \wp, \theta) \quad (12.1)$$

اما همانگونه که می‌دانیم ما به دنبال بهترین درخت  $T$  هستیم که  $\sigma$  و  $\wp$  در آن درخت برای ما اهمیت دارند. در واقع هر درخت  $T$  دارای امتیاز  $S(T)$  است که به صورت رابطه تعریف می‌شود.

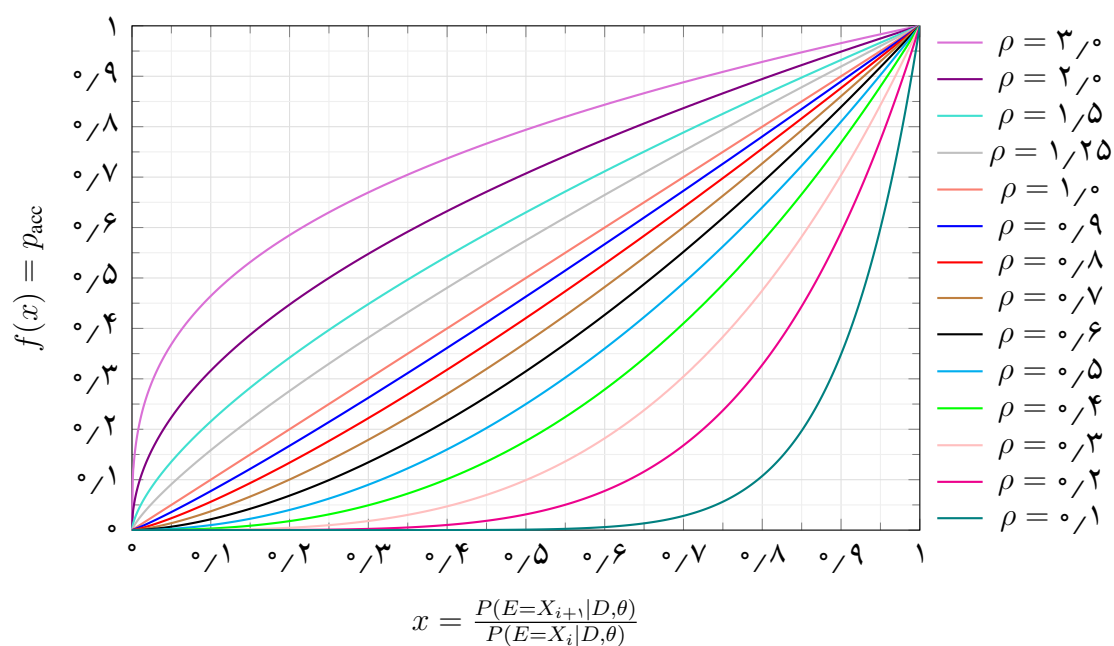
$$S(T) = P(D|T, \sigma^*, \wp^*), \quad \sigma^* = \arg \max_{\sigma} P(D|T, \sigma, \wp) \quad (13.1)$$

#### ۵.۴.۱ پذیرش پاسخ‌های جدید و یافتن بهترین پاسخ

در این مرحله ما دو پاسخ با امتیازهایشان در اختیار داریم که می‌توانیم برحسب آن‌ها برای ورود به تکرار بعد تصمیم‌گیری نماییم. این فرآیند توسط رابطه‌ای که در ادامه آمده است انجام می‌شود.

$$p_{\text{acc}} = \min \left[ 1, \left( \frac{P(E = X_{i+1}|D, \theta)}{P(E = X_i|D, \theta)} \right)^{\rho^{-1}} \right] \quad (14.1)$$

در رابطه ۱۳.۱، اگر پاسخ جدید بهتر از پاسخ فعلی باشد بیان می‌کند که افزایش بهینگی در پاسخ جدید باعث می‌شود تا صورت کسر مقداری بیش از منخرج بگیرد که در این صورت  $p_{\text{acc}}$  که برابر با احتمال پذیرش پاسخ جدید است، برابر ۱ خواهد شد که یعنی حتما پاسخ جدید به عنوان پاسخ پابرجا برای ورود به تکرار بعد در نظر گرفته می‌شود. اما اگر پاسخ جدید (درخت جدید) بهتر از پاسخ فعلی ارزیابی نشود ما آن را مستقیماً رد نمی‌کنیم و به احتمالی کمتر از ۱ ممکن است آن را بپذیریم. دلیل این پذیرش جلوگیری از به دام افتادن الگوریتم در پاسخ‌مان



شکل ۳.۱: نمودار تغییر احتمال پذیرش پاسخ جدید نامطلوب‌تر با توجه به مقدار پارامتر  $\rho$ .

در بیشینه‌های محلی<sup>۹</sup> است. در شکل تاثیر تغییر پارامتر  $\rho$  در احتمال پذیرش پاسخ‌های جدیدی که مطلوب‌تر از پاسخ فعلی نیستند نمایش داده شده است.

<sup>۹</sup>Local maxima

#### ۶.۴.۱ شبکه هرس‌کننده

#### ۷.۴.۱ شبکه بازاتصال‌کننده

#### ۸.۴.۱ جمع‌بندی

روش پیشنهادی ما الگو گرفته شده از روش ارائه شده در مقاله سایت بود اما با این تفاوت که در آنجا فرض مکان‌های بی‌نهایت بود ولی ما فرض اسکارلت را جایگزین کردیم که در این بین چون فضای جست و جو بزرگتر شد در نتیجه مجبور شدیم نحوه برداشتن گام‌های خود را تغییر بدیم و هوشمندانه تر جلو برویم که در این فضای بزرگتر بتوانیم به جواب مناسب برسیم. در سایت برای رسیدن به درخت بهتر به دنبال تنظیم کردن پارامترهای خطا بود در حالی که ممکن بود با جواب واقعی فاصله داشته باشد اما چون دنبال جواب با امتیاز بالا بود در نتیجه این پارامتری بودن و جست و جو برای مقادیر بهینه خطا در روش آن وجود داشت اما ما برای بهتر کردن امتیاز به جای تغییر پارامترهای خطا به ازای یک درخت جست و جوی ضمیمه کردن‌های مختلف سلول‌ها و لاس شدن جهش‌ها را جست و جو میکنیم.

## ۵.۱ روش پیشنهادی دوم

۱.۵.۱ مقدمه

۲.۵.۱ دادگان ورودی

۳.۵.۱ تبدیل داده‌ها به بردار ویژگی

۴.۵.۱ مدل سازی مساله

۱.۴.۵.۱ معماری شبکه

۵.۵.۱ تابع هزینه

۶.۵.۱ اصلاح خطا و یافتن درخت جواب

## ۶.۱ نتیجه گیری

جدول ۲.۱: پارامترهای مدل ریاضی

زمان خدمت دهی به بیمار در مرحله $k$ ام	$t_{ik}$
زمان فاری خدمت دهی به بیمار در محله $k$ ام	$\tilde{t}_{ik}$
مقدار بدبینانه (حداکثر) برای زمان خدمت دهی به بیمار در مرحله $k$ ام	$t_{ik}^p$
محتمل ترین مقدار برای زمان خدمت دهی به بیمار در مرحله $k$ ام	$t_{ik}^m$
مقدار خوشبینانه (حداقل) برای زمان خدمت دهی به بیمار در مرحله $k$ ام	$t_{ik}^o$

جدول ۳.۱: متغیرهای مدل ریاضی

متغیر صفر-یک تخصیص بیمار به تخت/اتاق عمل	$X_{ild_k}$
زمان شروع خدمت‌دهی به بیمار	$S_{ild_k}$
متغیر صفر-یک توالی بیماران	$Y_{ijkl_k}$
متغیر صفر-یک تخصیص جراح به بیمار	$V_{ni}$

## فصل ۲

# نتایج تجربی

### ۱.۲ پایگاه داده‌های ورودی

قبل از اینکه وارد روش پیشنهادی شویم به تشریح وردی‌های مسئله و داده‌هایی که مورد استفاده قرار خواهیم داد می‌پردازیم. داده‌های ورودی برابر ماتریس  $D_{m \times n}$  می‌باشد که بعد اول  $M$  برابر با ژن‌ها و بعد دوم  $N$  برابر سلول‌های نمونه‌برداری شده می‌باشد. در هر خانه  $d_{i,j}$  یک بردار داده قرار دارد که حاوی اطلاعات ژن  $j$  در سلول  $i$  می‌باشد.

#### ۱.۱.۲ پایگاه داده مصنوعی<sup>۱</sup>

با توجه به این نکته که از درخت فیلوژنی حقیقی<sup>۲</sup> داده‌های حقیقی موجود اطلاعی نداریم، به سراغ ساخت پایگاه داده مصنوعی می‌رویم. با استفاده از این پایگاه داده مصنوعی می‌توانیم در مورد روش‌هایی که در ادامه بیان خواهیم کرد یک معیار ارزیابی نسبتاً مناسبی داشته باشیم و تا حدودی از مشکلات روش‌های پیشنهادی آگاه شویم و به تصحیح آن بپردازیم. برای ساخت پایگاه داده مصنوعی که همان ماتریس ورودی  $D_{m \times n}$  می‌باشد، از دو روش مختلف با دو فرض مختلف استفاده خواهیم کرد که در ادامه به تشریح هر کدام خواهیم پرداخت. برای ایجاد پایگاه داده در این حالت ابتدا درختی تصادفی با پارامترهای  $n$ ،  $\zeta$  ایجاد می‌کنیم که  $n$  تعداد ژن‌ها

<sup>1</sup>Synthetic Dataset

<sup>2</sup>Ground-truth Phylogeny Tree

(جهش‌ها) بوده و  $\infty$  عددی در بازه  $(0, \infty)$  است که یک پارامتر کنترلی است که وظیفه‌اش کنترل کلی تعداد نسل‌های مختلف را از یک جمعیت در درخت فیلوژنی می‌باشد. حال برای تولید پایگاه داده مصنوعی به ترتیب سه گام زیر باید انجام شود.

- ایجاد یک درخت فیلوژنی تصادفی
  - تبدیل درخت فیلوژنی به ماتریس اطلاعات سلول-ژن ( $E$ )
  - اضافه کردن نویز به ماتریس  $E$  و تبدیل آن به ماتریس نویزی  $D$
- در ادامه هر بخش به صورت جداگانه به تفصیل شرح داده خواهد شد.

### ۱.۱.۱.۲ ساخت درخت تصادفی

برای ساخت درخت تصادفی از دو روش مختلف استفاده شده است که هرکدام جداگانه توضیح داده شده است.

#### روش اول: با استفاده از درخت تصادفی دودویی ژنولوژی<sup>۳</sup>

در این روش همان‌گونه که از نام آن مشخص است با استفاده از درخت تصادفی دودویی ژنولوژی به ساخت ماتریس داده ورودی مسئله می‌پردازیم که برای ساخت این دادگان از فرض‌های که در ادامه آمده است استفاده خواهیم کرد.

- در مرحله اول که ساخت درخت است به این صورت عمل می‌کنیم که به تعداد  $n$  گونه (سلول) در نظر می‌گیریم. سپس به ترتیب مراحل زیر را انجام می‌دهیم تا به درخت تصادفی مورد نظر برسیم.
- به هر کدام از  $n$  گونه متمایز در ابتدا وزن  $w_i = 1$  را اختصاص می‌دهیم که متناسب با احتمال انتخاب هر گونه در مراحل بعدی خواهد بود.

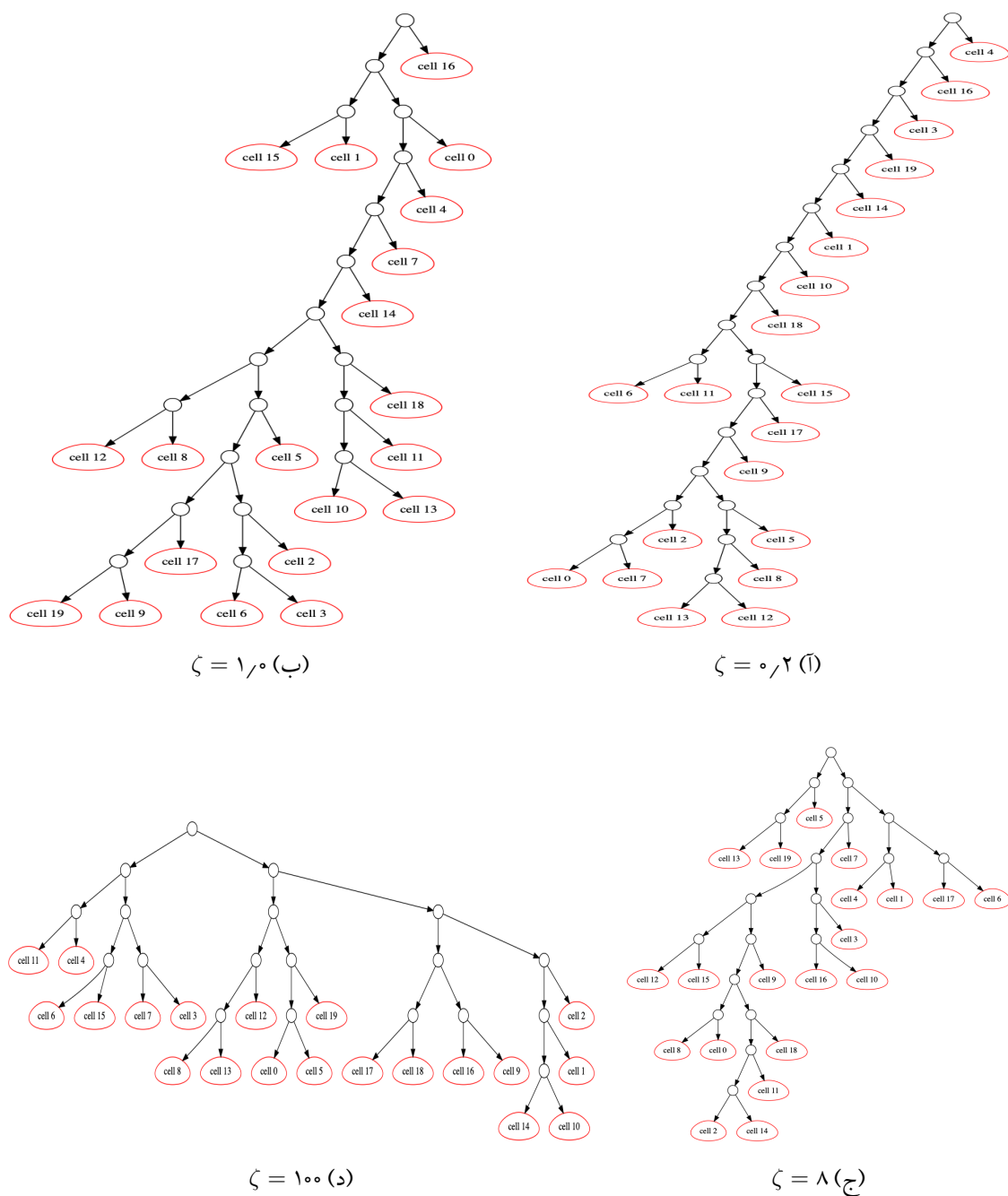
• برای هر گونه  $i$  تابع جرم احتمال را در ادامه به صورت  $F_i = \frac{w_i}{\sum_{i=1}^n w_i}$  در نظر می‌گیریم

- با استفاده از  $F$  دو گونه متمایز  $u, v$  را انتخاب می‌کنیم و به هم متصل می‌کنیم

<sup>3</sup>Random Binary Genealogical Tree



- به جای دو گونه  $u, v$  یک گونه جدید  $uv$  با وزن  $w_{uv} = \frac{w_u + w_v}{\sqrt[4]{\zeta}}$  را قرار می‌دهیم.
- تعداد گونه‌ها یک واحد کم شده است. بررسی می‌کنیم اگر تعداد گونه‌های باقی مانده از ۲ کمتر باشد درخت تصادفی ساخته شده است و پایان کار است. در غیر این صورت به مرحله اول بازمی‌گردیم.
- پارامتر  $\zeta$  به گونه‌ای کنترل‌کننده میزان ناپایداری در طی نسل‌ها می‌باشد. بطوریکه نمونه‌ای از نتایج مقادیر مختلف آن برای  $n = 20$  در شکل ۱.۲ آورده شده است. پس از ساخت درخت تصادفی به سراغ مرحله بعد یعنی تبدیل درخت به ماتریس ژن-سلول  $E$  می‌رویم.
- در ادامه با توجه به اینکه تعداد دلخواه جهش‌ها چه عددی بوده است یکی از گام‌های زیر را برمی‌داریم.
- اگر تعداد جهش‌ها  $M > N$  بوده باشد در آن صورت به صورت تصادفی به تعداد دفعات اختلاف یکی از انشعاب‌ها در درخت را به صورت تصادفی انتخاب کرده و آن جهش اضافه شده را تا تمامی نوادگان پیش خواهیم برد.
- اگر تعداد جهش‌ها  $M < N$  بوده باشد آنگاه مجدداً به اندازه تعداد اختلاف انشعاب‌هایی را انتخاب کرده و این بار جهش در آن انشعاب را تا تمامی نوادگان حذف می‌کنیم.
- به این ترتیب تمامی سلول‌ها را با تعداد جهش‌های انتخابی خواهیم داشت. در نهایت برای آخرین تغییر در جهش‌ها می‌توان یک گام دیگر برداشت که آن تولید یک عدد تصادفی کوچکتر از  $\frac{M}{4}$  است که به آن تعداد می‌توان جهش‌های موجود را از انشعابی برداشت و بر روی انشعابی دیگر قرار داد. با این کار ممکن است تعداد جهش‌ها در انشعاب‌های مختلف تغییر کند و چه بسا به مدل‌های واقعی نزدیکتر شود که البته در این پایان‌نامه از گام آخر صرف‌نظر کرده‌ایم.
- حال کار ما با پخش تصادفی جهش‌ها در پایگاه داده مجازی پایان یافته است. تا به اینجا ما در فرض خود از هر نمونه جمعیت مختلف یک سلول داشته‌ایم. اما در بعضی مواقع در پایگاه داده‌های واقعی ممکن است از یک جمعیت بیش از یک نمونه وجود داشته باشد که البته این امر لزوماً درست نیست به این دلیل که بعد از افزوده شدن نوین به داده‌ها ممکن است برخی سلول‌ها جهش‌هایشان مشابه هم شود. اما به هر حال اگر چنین چیزی را بخواهیم که داشته باشیم با انتخاب تصادفی برخی سلول‌ها (برگ‌ها) در درخت و کپی کردن آن‌ها می‌توان به چنین مقصودی رسید.



شکل ۱.۲: درخت فیلوژنی تصادفی تولید شده برای  $n = 20$  و  $\zeta$ های مختلف

## روش دوم: با استفاده از درخت تصادفی جهش‌های ژنی<sup>۴</sup>

این روش نیز تا حدود زیادی مشابه روش قبل است با این تفاوت که در اینجا به جای اینکه درخت تصادفی را

<sup>۴</sup>Random Mutation History Tree

با توجه سلول‌ها از پایین به بالا بسازیم، ابتدا یک درخت تصادفی بدون در نظر گرفتن سلول‌ها ایجاد می‌کنیم و سپس به تخصیص جهش‌ها به آن می‌پردازیم و در نهایت برای آخرین مرحله به تعداد دلخواه سلول را به درخت اضافه کرده و درخت را تکمیل می‌کنیم. در گام اول به تعداد  $M + 1$  نود در نظر می‌گیریم. مشابه حالت قبل با طی مراحل بکه در ادامه آمده است به ساختار یک درخت تصادفی می‌رسیم.

- به هر کدام از  $m$  نود متمایز در ابتدا وزن  $w_i = 1$  را اختصاص می‌دهیم که متناسب با روند حرکتی تومور به سمت آن جهش‌ها در مراحل بعدی خواهد بود.

- برای هر نود  $i$  تابع جرم احتمال را در ادامه به صورت  $F_i = \frac{w_i}{\sum_{i=1}^n w_i}$  بیان می‌شود در نظر می‌گیریم.

- با استفاده از  $F$  دو نود متمایز  $u, v$  را انتخاب می‌کنیم و به هم متصل می‌کنیم

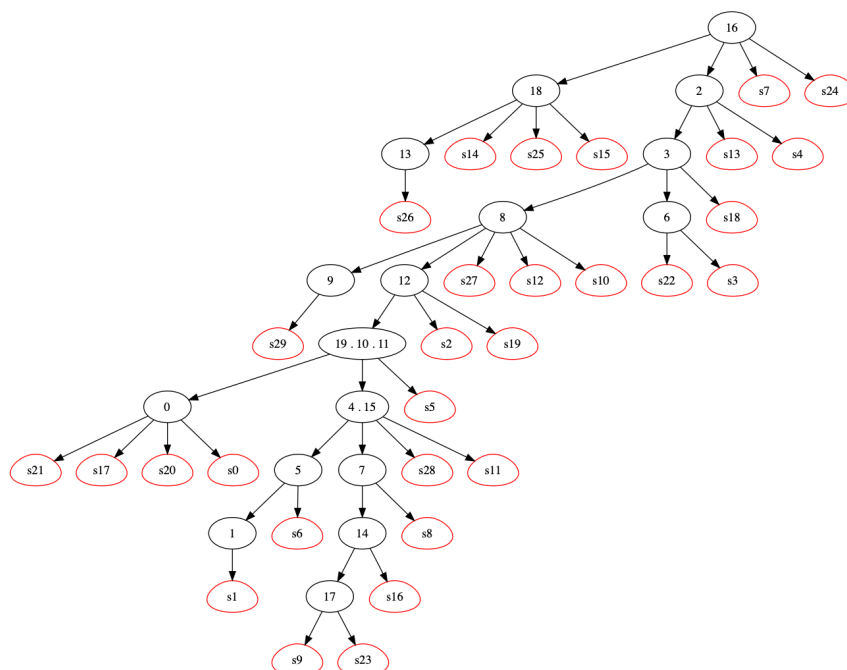
- به جای دو گونه  $u, v$  یک نود جدید  $uv$  با وزن  $w_{uv} = \frac{w_u + w_v}{\sqrt{2}}$  را قرار می‌دهیم.

- تعداد نودها یک واحد کم شده است. بررسی می‌کنیم اگر تعداد نودهای باقی مانده از ۲ کمتر باشد به مرحله بعد می‌رویم و در غیر این صورت به مرحله اول بازمی‌گردیم.

- در این مرحله تمامی برگ‌های درخت ساخته شده را حذف می‌کنیم و تنه باقی مانده را به عنوان درخت تصادفی جهش‌ها در نظر می‌گیریم.

پس از به پایان رسیدن مراحل که بیان شد درخت تصادفی آماده است و حال نوبت به تخصیص دادن خود ژن‌ها به هرکدام از این نودهای درخت است. برای این منظور به هرکدام از  $M$  نود یک ژن را به صورت تصادفی تخصیص می‌دهیم. پس از آن برای نهایی سازی درخت جهش‌ها از پارامتر دلخواه  $A = \lfloor \gamma * (M - 1) \rfloor$  استفاده می‌کنیم که  $\gamma$  عددی بین  $(0, 1)$  است و  $A$  تعداد یال‌هایی است که در درخت باید برداشته شود و دو نود آن با یکدیگر ادغام شود. این کار باعث می‌شود تا در درخت جهش‌ها در برخی نودها به جای یک جهش چند جهش داشته باشیم که بتواند به مدل داده‌های واقعی نزدیکتر باشد.

پس از تکمیل درخت جهش‌ها نوبت قرار دادن نمونه‌هایی بر روی آن است. به همین منظور با فرض اینکه  $N \geq M$  است. به تعداد  $M$  تا از سلول‌ها را به هر کدام از نودهای درخت جهش به عنوان برگ‌های جدید اضافه می‌کنیم و برای  $N - m$  سلول باقی مانده همین کار را این بار به صورت تصادفی انجام می‌دهیم. در نهایت درخت تصادفی جهش‌ها ساخته شده است که نمونه‌ای از آن را در شکل ۲.۲ قابل مشاهده است.



شکل ۲.۲: درخت جهش تصادفی با پارامترهای  $N = 30, M = 20, \zeta = 1, \gamma = 0.15$

### ۲.۱.۱.۲ تبدیل درخت به ماتریس ژن-سلول

با داشتن درخت (تولید شده با هرکدام از روش‌ها تفاوتی ندارد) در ادامه از فرض‌های مختلف در تولید ماتریس  $E$  می‌توان استفاده کرد.

#### فرض مدل مکان‌های بی‌نهایت<sup>۵</sup>

در این حالت فرض می‌کنیم که هر جهش اتفاق افتاده در درخت فیلوژنی در تمامی نسل‌های پس از آن باقی می‌ماند و هیچ‌گاه از بین نمی‌رود. در چنین حالتی درخت حاصل از این روش درختی یکتا بوده که به نام درخت فیلوژنی کامل<sup>۶</sup> شناخته می‌شود.

در این قسمت باید با استفاده از درخت تصادفی تولید بتوانیم ماتریس جهش‌ها را برای سلول‌های مختلف با فرض مکان‌های بی‌نهایت بدست آوریم. در ابتدا ماتریس  $E$  را به ابعاد  $M \times N$  ایجاد می‌کنیم و برای هر درایه  $i, j$  در

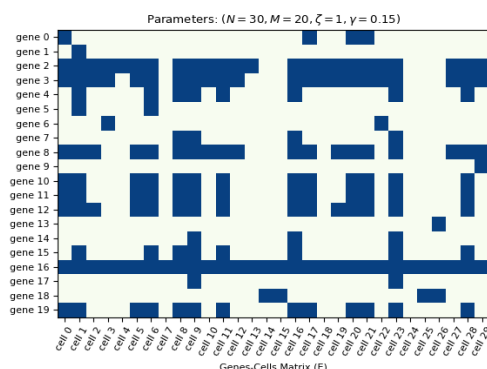
<sup>۵</sup>Infinite Site Models

<sup>۶</sup>Perfect Phylogeny Tree

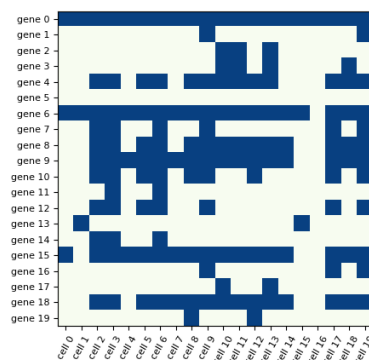
آن که  $i$  شماره جهش و  $j$  شماره سلول است به صورت فرمولی که در ادامه آمده است مقداردهی می‌کنیم.

$$E_{i,j} = \begin{cases} 1 & \text{if mutation } i \text{ is an ancestor of cell } j \\ 0 & \text{o.w} \end{cases} \quad (۱.۲)$$

به این ترتیب با فرض مدل مکان‌های بی‌نهایت ماتریس بدون خطا  $E$  را داریم که برای تصاویر دو روش درخت مرحله قبل در شکل ۳.۲ بدست آمده‌اند.



(ب) ماتریس درخت شکل ۲.۲



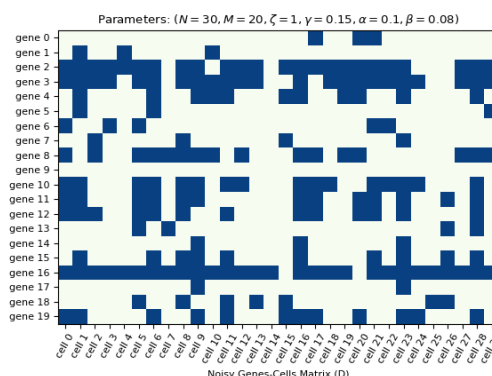
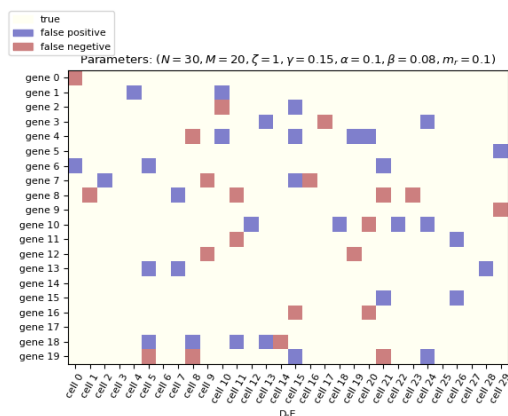
(آ) ماتریس درخت شکل ۱.۲ ب

شکل ۳.۲: ماتریس‌های ژن-سلول ( $E$ ) بدست آمده از درخت‌های تصادفی ساخته شده

### ۳.۱.۱.۲ اضافه کردن نویز به ماتریس ژن-جهش

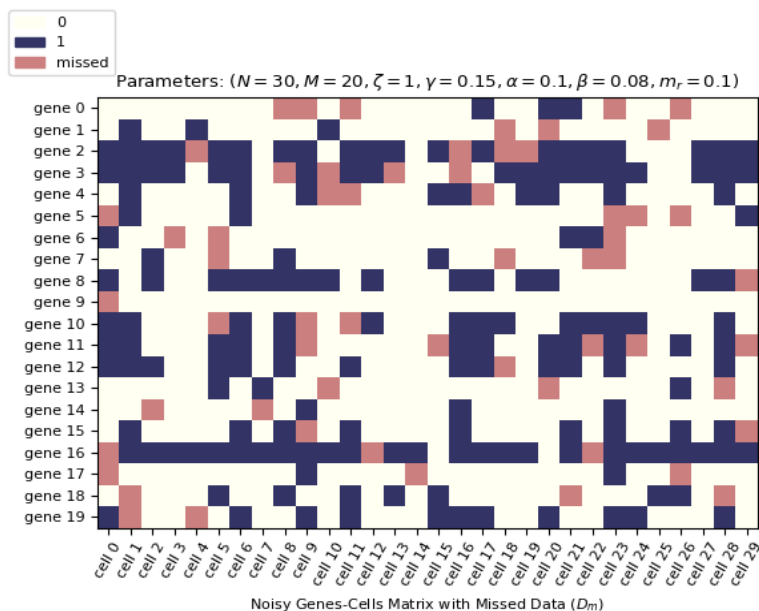
برای قسمت نهایی آماده سازی پایگاه داده مجازی نیاز است تا به ماتریس  $E$  با پارامتر  $\Theta = (\alpha, \beta, m_r)$  نویز اضافه کنیم و آن را به ماتریس  $D$  تبدیل کنیم که  $\alpha = P(D_{ij} = 1 | E_{ij} = 0)$  و  $\beta = P(D_{ij} = 0 | E_{ij} = 1)$  است و همچنین  $m_r \in (0, 1)$  که نرخ داده‌های از دست رفته را مشخص می‌کند. برای این منظور به ازای تمامی درایه‌های  $E$  ماتریس  $E$  هربار یک عدد تصادفی با توزیع یکنواخت بین  $(0, 1)$  بوجود می‌آوریم و اگر عدد تولید شده کوچکتر از  $\alpha$  بود آنگاه آن درایه در ماتریس  $D$  را برابر با ۱ قرار می‌دهیم. به همین ترتیب مجدداً این بار برای درایه‌های ۱ ماتریس  $E$  این کار را تکرار می‌کنیم و اگر عدد تصادفی تولید شده کوچکتر از  $\beta$  شد، درایه متناظر را در ماتریس  $D$  برابر با ۰ قرار می‌دهیم.

پس از اتمام کار نوبت به اضافه کردن داده‌های از دست رفته است. برای این منظور با نرخ  $m_r$  بعضی از درایه‌های ماتریس  $D$  را برابر با ۲ قرار می‌دهیم که به منزله در دسترس نبودن اطلاعات است. نام ماتریس نهایی را که شامل داده‌های از دست رفته است  $D_m$  می‌گذاریم. در ادامه تصاویر اضافه شدن نویز به ماتریس شکل ۳.۲.۲ ب در شکل ۴.۲ آمده است.



(ب) نویزی اضافه شده با پارامترهای  $\alpha = 0.1, \beta = 0.08$

(آ) ماتریس نویزی با  $\alpha = 0.1, \beta = 0.08$

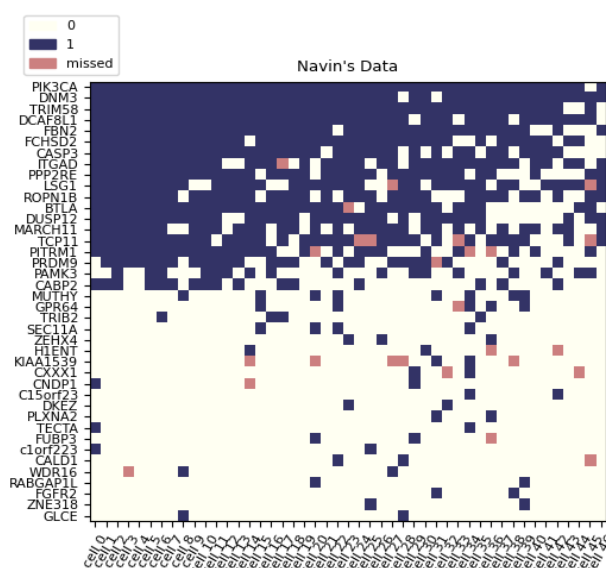


(ج) ماتریس نویزی به همراه داده‌های از دست رفته با پارامترهای  $\alpha = 0.1, \beta = 0.08, m_r = 0.1$

شکل ۴.۲: ماتریس‌های ژن-سلول همراه با نویز و داده‌های از دست رفته شکل ۳.۲.۲ ب که برای ورودی مسئله آماده شده است.

۲.۱.۲ پایگاه داده حقیقی<sup>۷</sup>

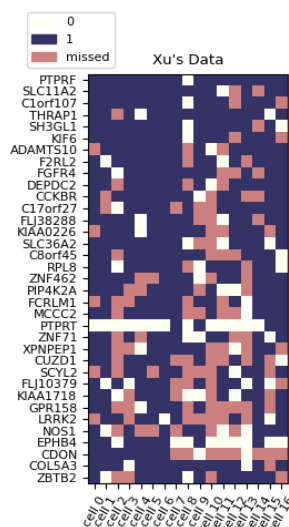
به عنوان پایگاه داده حقیقی از پایگاه داده استفاده شده در مقاله SCITE به عنوان پایگاه داده حقیقی اصلی استفاده خواهیم کرد که ماتریس داده ورودی آن به صورت شکل ۵.۲ می باشد. همچنین پایگاه داده حقیقی Xu



شکل ۵.۲: داده‌های حقیقی Navin در مقاله SCITE

نیز که در مقاله SCITE مورد استفاده قرار گرفته است در شکل ۶.۲ آمده است.

<sup>۷</sup>Real Dataset



شکل ۶.۲: داده‌های حقیقی Xu در مقاله SCITE

## ۲.۲ روش پیشنهادی بدست آوردن درخت فیلوژنی

پس از تخمین داده‌های از دست رفته، در این بخش به معرفی روش پیشنهادی برای یافتن درخت فیلوژنی می‌پردازیم. در روش‌های گذشته که رویکرد آن در ادامه بیان شده است به موفقیت نرسیدیم و این بار در نظر داریم تا با استفاده از یک درخت در ساختار شبکه ژن‌ها بتوانیم به یک درخت فیلوژنی مناسب دست یابیم.

- استفاده از شبکه ژن‌ها

۷ استفاده از یک گراف ابتدایی و سپس تغییر و هرس کردنش تا رسیدن به درخت جهش‌ها

۳ استفاده از یک درخت نمونه نابینه و تغییر اتصالات تا رسیدن به درخت بهینه جهش‌ها

- استفاده از شبکه سلول‌ها

۷ استفاده از یک گراف سلول‌ها و بهینه کردن ارتباطات بین آن‌ها و سپس تبدیل آن به درخت فیلوژنی

در رویکرد اول ما از شبکه‌های ژنی استفاده خواهیم نمود. این شبکه‌ها نودهایی معادل با یک ژن متمایز را در نظر می‌گیرند. در گذشته با استفاده از شبکه‌ای کامل با وزن‌های متفاوت که بر حسب اطلاعات ورودی به الگوریتم تعیین می‌شد، متاسفانه به موفقیت خاصی نرسیدیم. همچنین مشابه همین رویکرد را در ساختار



شبکه‌های سلولی دنبال کردیم که مجددا پیشرفت قابل ملاحظه‌ای حاصل نشد. به همین جهت این بار در این گزارش با تغییری اساسی به دنبال یافتن روشی مناسب برای استنتاج درخت فیلوژنی می‌باشیم.

## ۱.۲.۲ استفاده از شبکه ژن‌ها برای یافتن درخت فیلوژنی

در این رویکرد با استفاده از شبکه‌ای که نودهایی معادل ژن‌ها داشته باشد سعی داریم تا به درخت فیلوژنی بهینه برسیم.

### ۱.۱.۲.۲ استفاده از یک درخت نمونه نابینه و تغییر اتصالات تا رسیدن به درخت فیلوژنی بهینه

در این روش قصد داریم تا با شروع از یک درخت نمونه که در ابتدا به صورت تصادفی از اتصال ژن‌ها بوجود آمده است، به بهینه‌ترین درخت ممکن برسیم. این روش به صورت تکرارواره با تغییر اتصالات درخت سعی در بدست آوردن درختی مطلوب‌تر دارد که شرایط و روابط تاثیرگذار در آن به تفصیل شرح داده خواهد شد. در واقع این روش پیشنهادی یک جست‌وجوی حریصانه می‌باشد که طی شرایطی می‌توان انتظار داشت که به پاسخ بهینه دست یافته شود. این روش به نام روش زنجیره مارکو مونت-کارلو<sup>۸</sup> شناخته می‌شود که در بسیاری از مقالات مرتبط نیز مورد استفاده قرار گرفته شده است.

برای شروع یک درخت تصادفی  $T$  را با نودهایی معادل ژن‌های پایگاه داده ورودی در نظر می‌گیریم که در گام اول به صورت تصادفی ساخته شده است. در گام‌های بعدی یک نود  $n_1$  را از درخت  $T$  به صورت تصادفی انتخاب می‌کنیم. سپس زیردرخت با ریشه این نود را از درخت کم می‌کنیم. حال در درخت باقی‌مانده یک نود دیگر  $n_2$  را به صورت تصادفی انتخاب می‌کنیم و آن زیر درخت قبلی با ریشه  $n_1$  را به  $n_2$  متصل می‌کنیم و درخت جدید را  $T_n$  نام‌گذاری می‌کنیم. پس از آن با احتمال،

$$P = \min \left( 1, \frac{Eng(T)}{Eng(T_n)} \right) \quad (2.2)$$

درخت جدید بدست آمده  $T_n$  را به عنوان نتیجه این گام می‌پذیریم و در غیر این صورت درخت این گام نیز همان درخت سابق  $T$  باقی خواهد ماند. در رابطه ۲.۲، تابع  $Eng$  برای یک درخت در واقع انرژی آن درخت را محاسبه

<sup>8</sup>Markov Chain Monte Carlo

می‌کند و ما به دنبال پایدارترین درخت هستیم که کمترین انرژی را داشته باشد. تعریف این تابع برای یک درخت به این صورت است که با توجه به نمونه‌هایی که در دادگان ورودی  $D$  قرار دارد و اینکه کدام ژن بالاتر یا پایین‌تر از دیگر ژن‌ها قرار دارد به درخت یک نمره انرژی منصوب می‌کند که به صورت فرمول ۳.۲ بیان می‌شود.

$$Eng(T) = ||E - \hat{E}|| \quad (۳.۲)$$

که در اینجا  $\hat{E}$  ماتریس تخمین زده شده روش پیشنهادی با توجه به درخت نهایی بدست آمده خواهد بود. در واقع ماتریس  $E$  همان ماتریس صحیح بدون خطا است که جهش‌های مختلف را به ازای سلول‌های مختلف مشخص می‌کند. هنگامی که ما درخت ساخته شده فرضی  $T$  را داشته باشیم می‌توانیم در دو گام به  $\hat{E}$  برسیم. توجه به این نکته ضروری است که اگر در واقعیت فرض ما که همان مکان بی‌نهایت بود کامل برقرار باشد و  $E$  را داشته باشیم، حتماً باید بتوانیم به درختی با  $Eng(T) = 0$  دست یابیم. اما از آنجایی که ما  $D$  را به عنوانی از تخمین  $E$  داریم بنابراین محاسبه خطای ما نیز تخمینی از خطای واقعی خواهد بود نه خود آن که در اصل به صورت فرمول ۴.۲ می‌شود.

$$Eng(T) \approx Err(T) = ||D - \hat{E}|| \quad (۴.۲)$$

می‌دانیم که هر نود از این درخت  $T$  یک مکان برای اتصال سلولی می‌تواند باشد که در این صورت معنی آن اینگونه خواهد بود که سلول ضمیمه شده به آن نود تمام جهش‌های والد خود را داشته است. بنابراین در گام اول نیاز است تا هر مکان از درخت مشخص شود که چه نمونه‌هایی می‌تواند تولید نماید. این اطلاع توسط ماتریس  $A$  مشخص می‌شود که به صورت زیر از روی درخت ساخته خواهد شد.

$$A_{i,j} = \begin{cases} 1 & \text{اگر } j = i \text{ یا جهش } i \text{ والد جهش } j \text{ باشد} \\ 0 & \text{در غیر این صورت} \end{cases} \quad (۵.۲)$$

حال در گام دوم کافی می‌توانیم با توجه به یک معیار بهترین انتخاب را برای ضمیمه کردن سلول‌ها (نمونه‌های) موجود به درخت داشته باشیم. به همین جهت در ماتریس  $D$  که هر ستون آن برابر با نمایش یک سلول است، می‌تواند با هر ستون از ماتریس  $A$  مقایسه شود و بهترین ستونی که از  $A$  انتخاب شود برابر با جایگاه مناسب

ضمیمه شدن نمونه با مقداری خطا به درخت  $T$  است. حال با توجه به اینکه فرض مکان‌های بی‌نهایت را داشتیم ماتریس  $E$  را به صورت زیر می‌سازیم.

$$\hat{E}_{i,j} = A_{i,\sigma_j} \quad (۶.۲)$$

که  $\sigma_i$  برابر با بهترین نود (ژن) برای اتصال نمونه بردار  $d_j$  است که بهترین جایگاه به صورت فرمول زیر انتخاب می‌شود.

$$\sigma_j = \arg \max_{x \in [1 \rightarrow M]} \sum_{i=1}^M \left[ \begin{aligned} &A_{i,x} D_{i,j} (1 - \beta) + (1 - A_{i,x}) (1 - D_{i,j}) (1 - \alpha) + \\ &A_{i,x} (1 - D_{i,j}) \beta + (1 - A_{i,x}) D_{i,j} \alpha \end{aligned} \right] \quad (۷.۲)$$

حال با داشتن ماتریس  $\hat{E}$  می‌توان خطای درخت را محاسبه نمود و با هدایت mcmc طبق فرمول ۸.۲ به درخت بهینه  $T_{op}$  رسید.

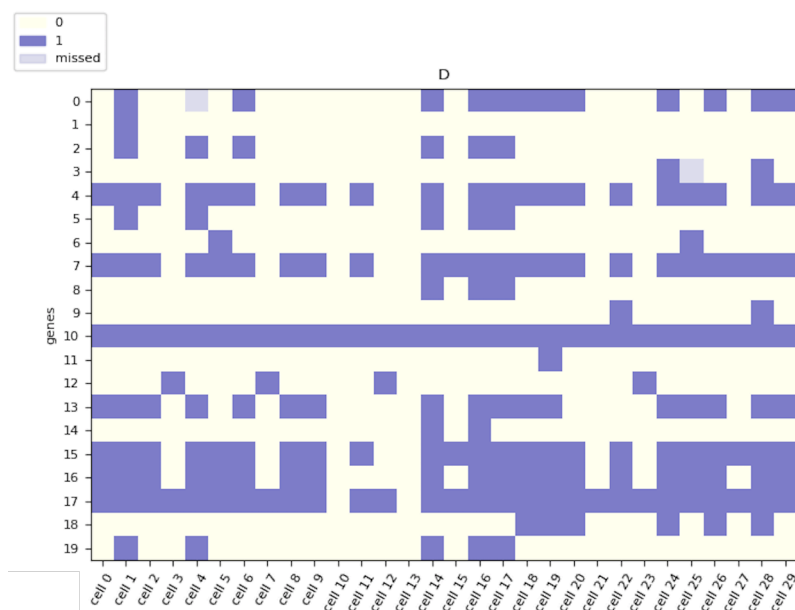
$$T_{op} = \min_{T \in \text{All possible } T} \left( ||D - \hat{E}_T|| \right) \quad (۸.۲)$$

## ۳.۲ نتایج تجربی

در این بخش به نتایج بدست آمده برای روش پیشنهادی می‌پردازیم و برای هر دو داده مصنوعی و حقیقی نتایج بدست آمده را تحلیل خواهیم نمود.

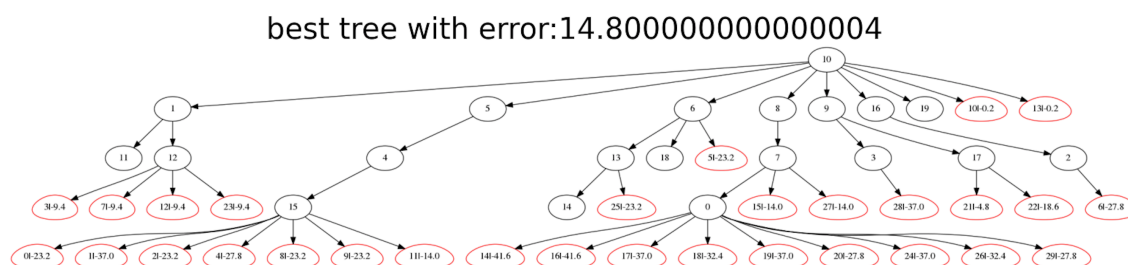
### ۱.۳.۲ نتایج بر روی پایگاه داده مصنوعی

همان‌گونه که در بخش دوم توضیح داده شد با توجه به سختی دسترسی به پایگاه داده‌های حقیقی و اینکه در آن‌ها نیز حقیقت داده‌ها ( $E$ ) وجود ندارد تصمیم به ایجاد پایگاه داده‌ای مصنوعی گرفته شد که با کمک آن بتوان ارزیابی مناسبی از روش پیشنهادی و میزان کارایی و مقاومت روش را نسبت به تغییر پارامترها سنجید. فرض کنید ماتریس ورودی شکل ۷.۲ را در اختیار داریم و می‌خواهیم بهترین درخت فیلوژنی را برای آن بیابیم.



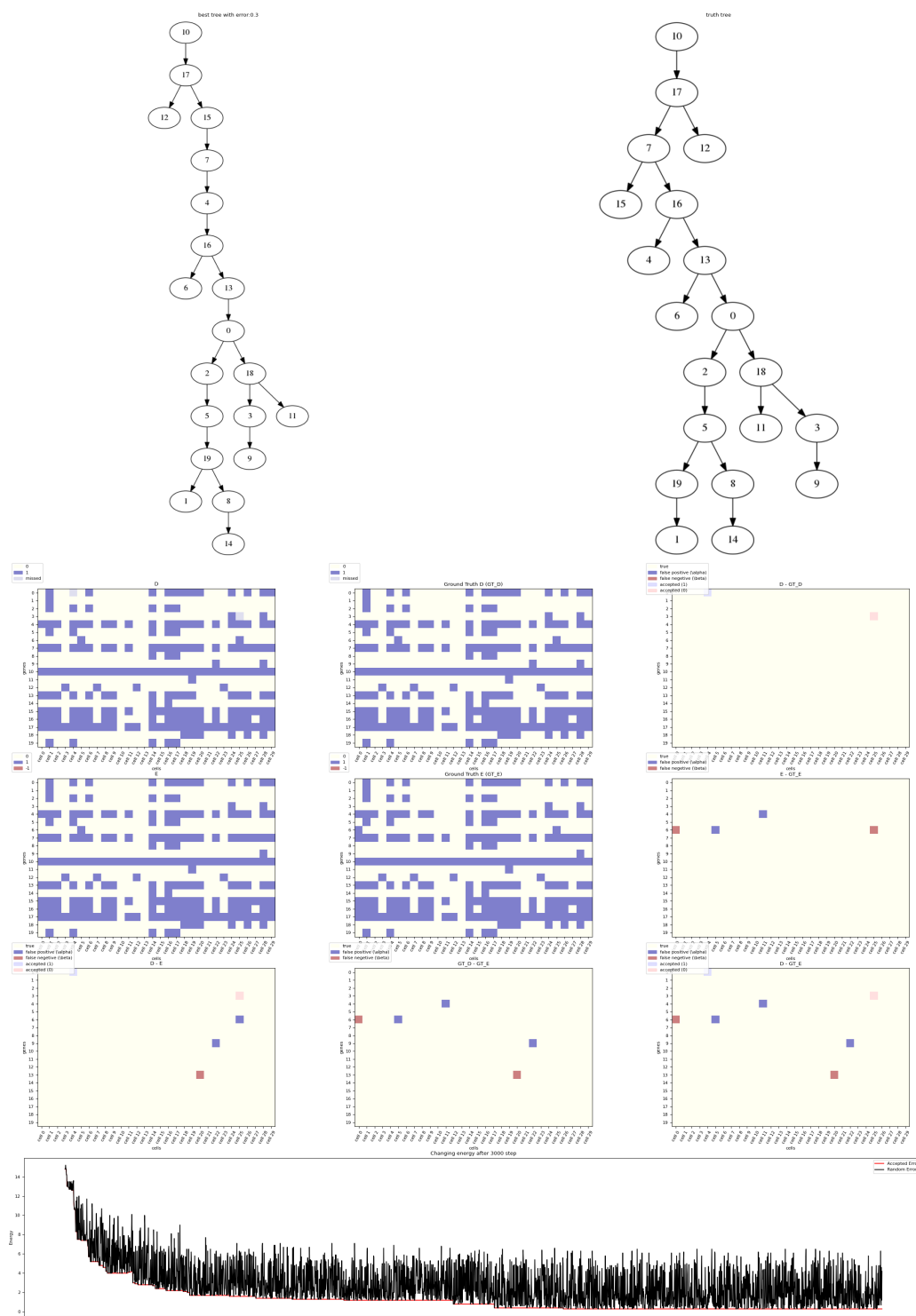
شکل ۷.۲: نمونه‌ای تصادفی از ماتریس ورودی  $D$

حال یک درخت تصادفی به صورت شکل ۸.۲ می‌سازیم. در درخت شکل ۸.۲ نمونه‌ها (سلول‌ها) با رنگ قرمز به درخت متصل شده‌اند که البته این ضمیمه بهترین ضمیمه ممکن است و میزان انرژی (خطای) هر ضمیمه نیز در کادر قرمز رنگ سلول‌ها به صورتی عددی منفی نوشته شده است.



پس از این مرحله اگر ۳۰۰۰ گام MCMC را اجرا نماییم می‌توانیم نتیجه حاصله را در شکل ۹.۲ مشاهده کنیم. در این شکل دو درخت وجود دارد که درخت سمت راستی درخت حقیقی است که به دنبال آن بودیم و درخت سمت چپ بهترین درخت یافته شده است. همچنین در پایین شکل، ۹ ماتریس مشاهده می‌شود که ماتریس‌ها سمت راست و پایین به نوعی بیان‌کننده میزان خطا بین ۴ ماتریس سمت چپ بالا می‌باشند. در بالای هر ماتریس نام آن نوشته شده است و در نهایت در انتهای تصویر نیز روند کاهش خطا و تلاش‌های MCMC در گام‌های مختلف قابل مشاهده است. فقط نکته‌ای که وجود دارد این است که خطای نوشته شده در تصاویر برابر  $0.1$  مقیاس نوشته شده است.

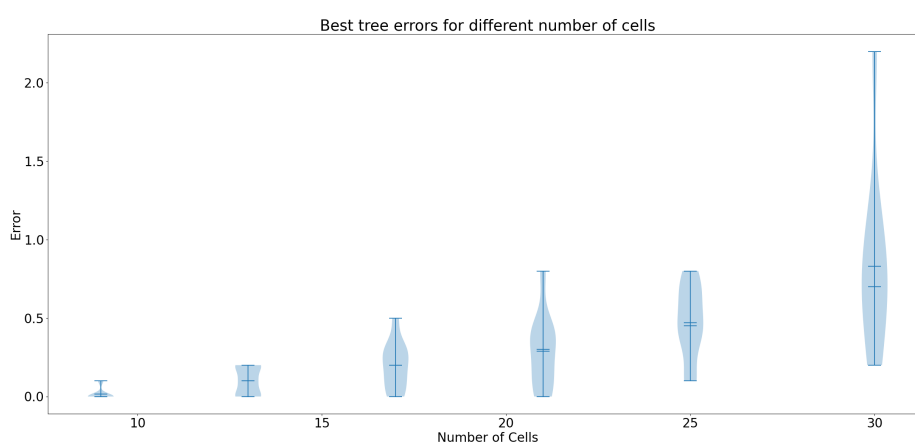
۲۹



شکل ۹.۲: نتیجه اجرای روش پیشنهادی برای ماتریس شکل ۷.۲

اما برای بررسی مناسب‌تر تعدادی تست را به ازای  $M$  و  $N$ ‌های مختلف اجرا نمودیم که به صورت خلاصه نتایج حاصل از آن در ادامه قابل مشاهده است.

در شکل ۱۰.۲ مقدار خطای درخت بهینه یافته شده قابل مشاهده است که نشان می‌دهد هر چه تعداد نمونه‌ها افزایش پیدا میکند و اندازه ماتریس ورودی بزرگ‌تر می‌شود، مقدار خطا نیز افزایش می‌یابد. در این اجرا تعداد جهش‌ها نیز عددی بین تعداد نمونه‌ها و نصف تعداد نمونه‌ها بوده است. حال برای اینکه متوجه شویم آیا این

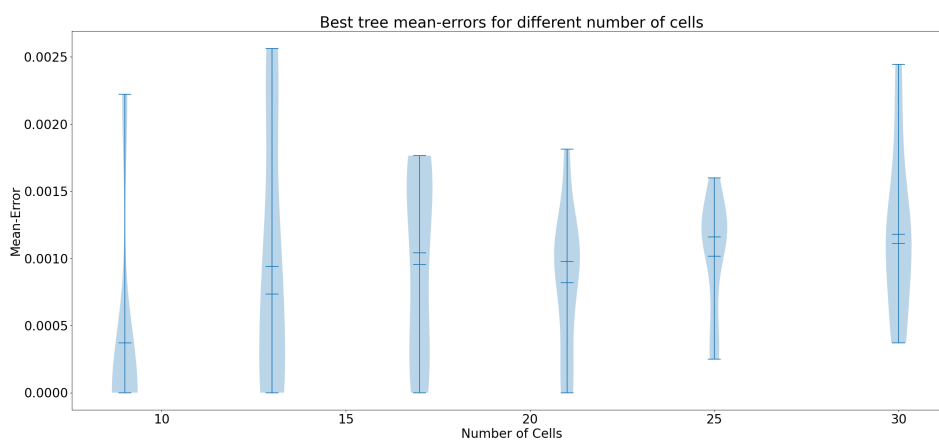


شکل ۱۰.۲: نتیجه اجرای روش پیشنهادی برای تعداد نمونه‌های مختلف

افزایش خطا بخاطر ضعف روش پیشنهادی است یا ماهیت داده‌های ورودی میزان خطا را در هر اجرا بر تعداد خانه‌های ماتریس  $D$  تقسیم می‌کنیم که در آن صورت به نمودار شکل ۱۱.۲ می‌رسیم. در این نمودار جدید مشخص می‌شود که با افزایش اندازه ماتریس ورودی روش پیشنهادی سعی می‌کند تا خطا را به ازای هر داده کنترل کند که نشان از کارآمدی روش پیشنهادی می‌باشد.

## ۲.۳.۲ نتایج بر روی داده‌های حقیقی

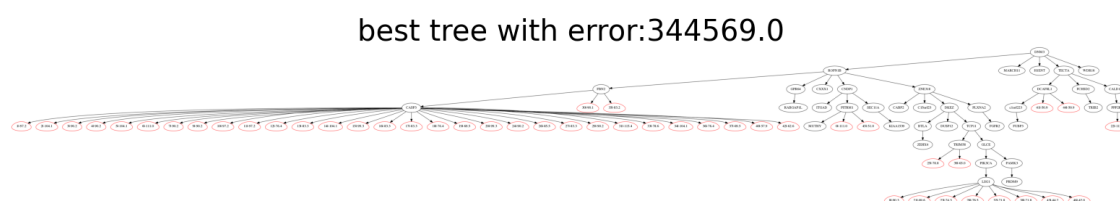
در این قسمت به ارائه گزارش و نتایج حاصل از روش‌های پیشنهادی با استفاده از داده‌های حقیقی برای بدست آوردن درخت فیلوژنی خواهیم پرداخت.



شکل ۱۱.۲: نتیجه اجرای روش پیشنهادی برای تعداد نمونه‌های مختلف

### ۱.۲.۳.۲ نتایج بهینه‌سازی درخت ژنی

همان‌طور که در فصل قبل بیان شد، یکی از روش‌های بدست آوردن درخت فیلوژنی استفاده از یک درخت تصادفی نابینه ژنی بود که طی تکرار گام‌هایی سعی در تغییر اتصالات و یافتن درخت بهینه داشت که بتواند روند صحیح تغییرات ژنی را در تومور مورد نظر نمایش دهد. نتیجه بدست آمده بر روی پایگاه داده حقیقی Navin به شرح زیر می‌باشد که عکس ۱۲.۲ درخت تصادفی اولیه الگوریتم را نشان می‌دهد که انرژی آن نیز بالای تصویر نوشته شده است.

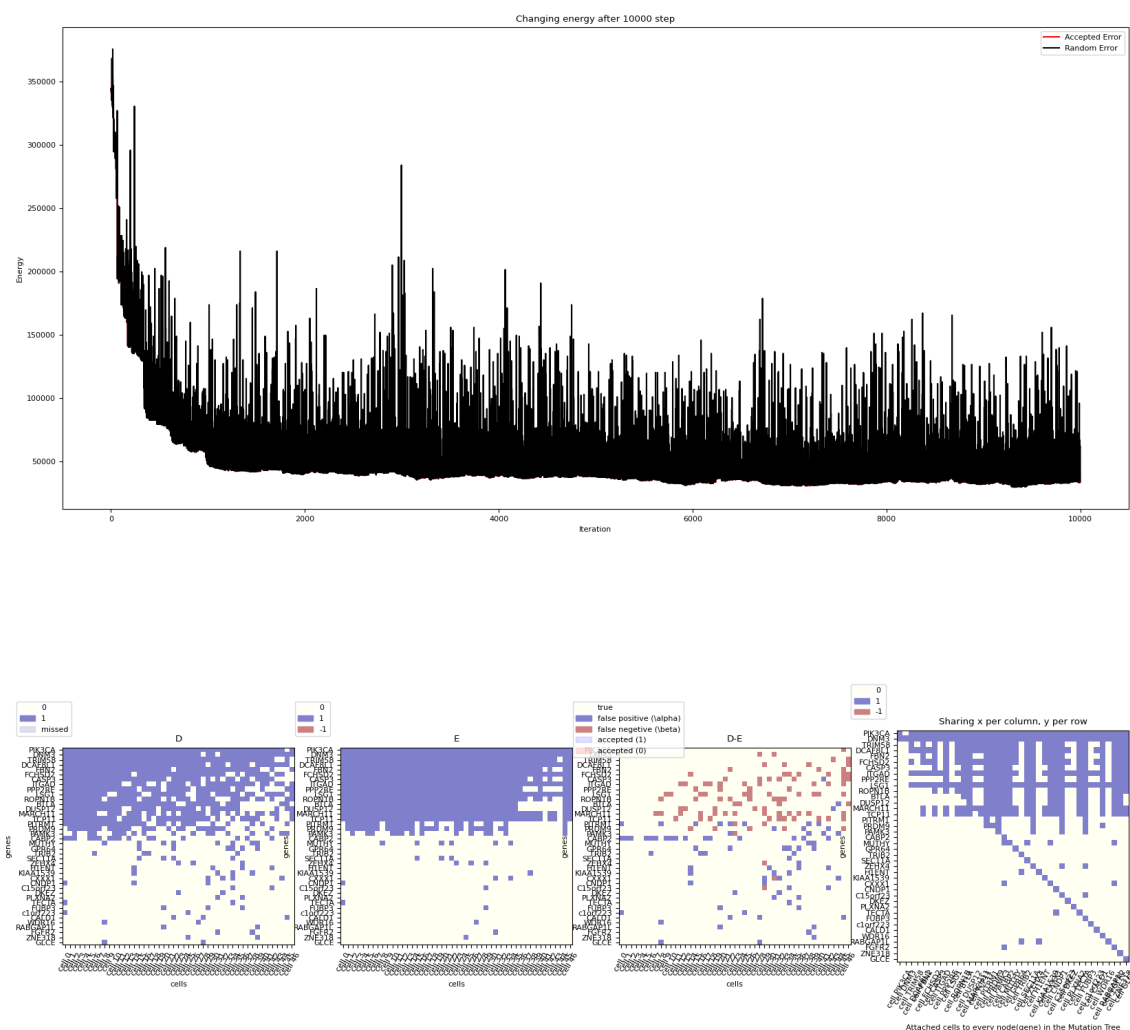


شکل ۱۲.۲: درخت تصادفی اولیه

تصویر ۱۳.۲ نیز نمودار تغییر انرژی را طی گام‌های مختلف نمایش پیشنهادی مشخص می‌کند. در نهایت تصویر بهترین درخت یافته شده به همراه انرژی آن.

در نهایت برای مقایسه نیز تصویر درخت حاصله در مقاله اصلی SCITE را در شکل ۱۵.۲ نمایش داده شده





شکل ۱۳.۲: نمودار تغییر انرژی در طی گام‌های مختلف

است. همان‌طور که مشخص است مقدار انرژی بدست آمده برای خروجی الگوریتم پیشنهادی بهتر (کمتر) از انرژی درخت SCITE می‌باشد که دلیل بر بهینه‌تر بودن درخت روش پیشنهادی ارائه شده در این گزارش است.

best tree with error:29929.0



شکل ۱۴.۲: بهترین درخت یافته‌شده و خروجی الگوریتم برای مقاله SCITE

## ۴.۲ گام‌های آتی

در ادامه برای تکمیل روش پیشنهادی در دو قسمت نیاز به بهبود وجود دارد. قسمت اول مربوط به درخت اولیه است و قسمت دیگر مربوط به سرعت MCMC می‌باشد.

### ۱.۴.۲ بهبود در ساخت درخت اولیه

در حال حاضر ما درخت اولیه را به صورت تصادفی انتخاب می‌کنیم که می‌توان در این مرحله درخت اولیه را با استفاده از مفروضات مدل مکان‌های بی‌نهایت و با توجه به ماتریس ورودی بهبود بخشید. این کار باعث می‌شود تا شروع الگوریتم از نقطه بهتری باشد که در این صورت هم گام‌های لازم برای رسیدن به درخت بهینه می‌تواند کمتر شود و هم اینکه احتمال قرار گرفتن در نقاط اکسترمم نسبی را کاهش می‌دهیم.

### ۲.۴.۲ افزایش سرعت همگرایی MCMC

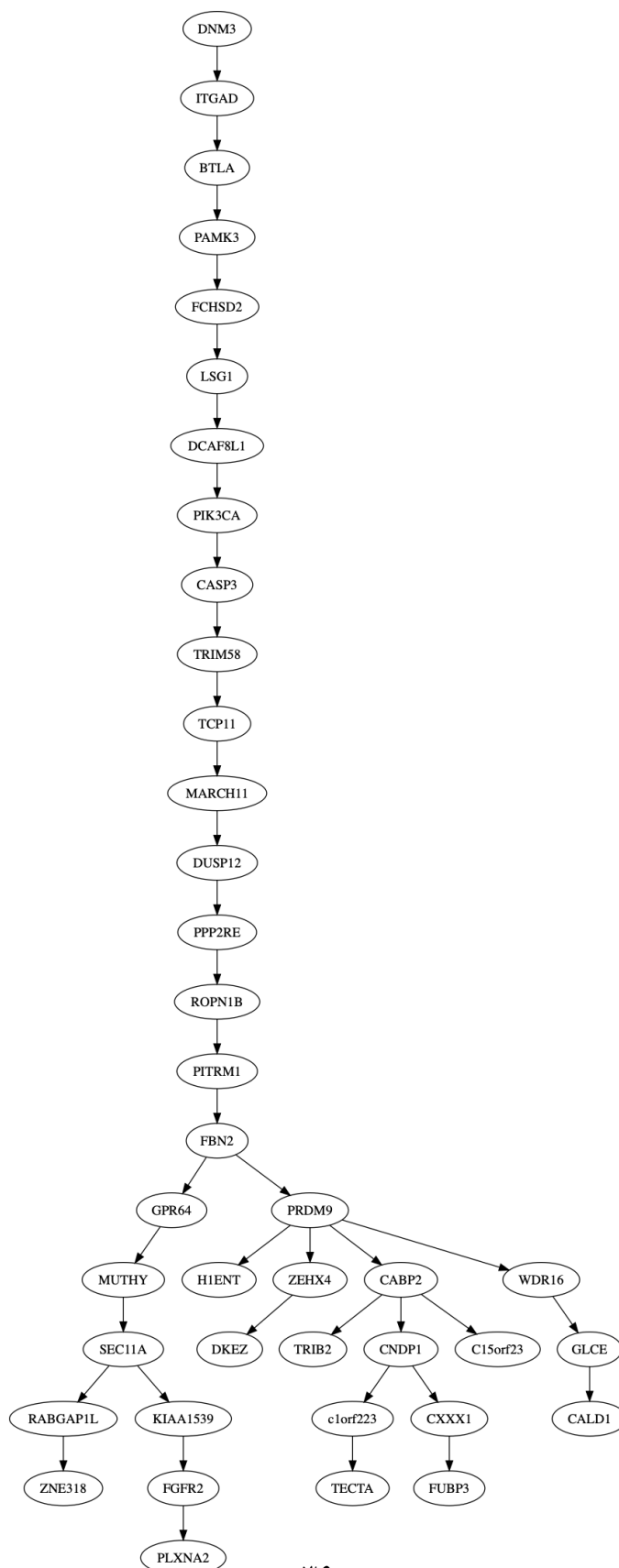
در این بخش نیاز است تا در دو قسمت روش پیشنهادی بهبود یابد.

#### ۱.۲.۴.۲ تنوع در گام‌ها با استراتژی معقول

برای این بخش کاری که باید انجام شود این است که بتوان برای افزایش سرعت همگرایی از روش‌های مختلف در گام‌ها استفاده کرد. برای مثال در حال حاضر می‌توان از سه روش مختلف در هر گام استفاده نمود. روش اول تعویض دو نود در درخت می‌باشد. روش دوم جدایی یک زیر درخت و اتصال آن به محلی دیگر می‌باشد و در نهایت روش سوم تعویض دو زیر درخت با یکدیگر می‌باشد. با انتخاب یک استراتژی مناسب بین هر کدام از این روش‌ها در گام‌های مختلف احتمالاً بتوان سرعت همگرایی را افزایش داد.

#### ۲.۲.۴.۲ قرار دادن احتمال وزن‌دار به ازای هر انتخاب

در حال حاضر ما در هر کدام از روش‌های مختلف که در بخش قبل برای گام‌های MCMC بیان کردیم، انتخاب نودها را به صورت کاملاً یکنواخت انجام می‌دهیم. در صورتی که احتمالاً بتوان با تعریف فرمولی مناسب این احتمال انتخاب بین نودهای مختلف در درخت را از حالت یکنواخت خارج کرد و در نتیجه مجدداً سرعت همگرایی الگوریتم را افزایش داد.



## مراجع

- [1] Davis, Alexander and Navin, Nicholas E. Computing tumor trees from single cells. *Genome biology*, 17(1):1–4, 2016.



# واژه‌نامه فارسی به انگلیسی

## ب

loss-supported	با در نظر گرفتن حذف خطا	Adenine	آدنین
feedback	بازخورد	Detection	آشکارسازی
Bias	بایاس	Binary cross-entropy	آنتروپی متقاطع دودویی
reversion	برگشت	Pathology	آسیب‌شناسی
Dimension	بعد	conflict	ابهام
Back-propagation	بک‌پروپگیشن	Joining-based	اتصال اکتشافی همسایگی
Optimization	بهینه‌سازی	heuristic	
Bayesian	بیزی	Consensus	اجماعی
Over-fitting	بیش‌پردازش	Cross validation	ارزیابی‌های متقابل
Maximum-likelihood	بیشینه درست‌نمایی	Gold standard	استاندارد مرجع
Local maxima	بیشینه‌های محلی	Surrounding stroma	استرومای اطراف
Computer Vision	بینایی کامپیوتر	Inference	استنباط
Biopsy	بیوپستی	Ascii	اسکی
		X-ray	اشعه ایکس

## پ

Consistency	پایداری	Missing data	اطلاعات از دست رفته
Low sequence coverage	پوشش توالی کم	concordant gain of	افزایش همزمان جهش‌ها
Read coverage	پوشش خوانش	mutations	
Tumor progression	پیشرفت تومور	Concatenation	الحاق
Spontaneous	پیشرونده	Standard deviation	انحراف استاندارد
Pixel	پیکسل	Stride	اندازه گام
		Oncogene	انکوژن
		Indel	ایندل

Somatic mutation . . . . . جهش بدنی	ت
Single nucleotide . . . . . جهش تک نوکلئوتیدی	
mutation	تابع فعالیت . . . . . Activation function
Driver mutation . . . . . جهش راننده	تابع واحد اصلاح شده خطی Rectified linear unit (ReLU)
Single Somatic Mutation . . . . . جهش ساده بدنی	تابع هزینه . . . . . Cost function
Passenger mutation . . . . . جهش مسافر	ترک آلل . . . . . Allele dropout
Point mutation . . . . . جهش نقطه‌ای	ترک آلی . . . . . Allelec dropout
	تشخیص به کمک کامپیوتر . . . . . Computer-aided diagnosis
چ	تعداد خوانش‌ها . . . . . Read counts
Multi-state . . . . . چند حالت	تغییر ساختاری . . . . . Single variant
	تغییر همبستگی داخلی . . . . . Internal covariate shift
ح	تغییرات بدنی تک نوکلئوتیدی . . . . . Somatic single nucleotide variation
Long Short Term . . . . . حافظه‌ی کوتاه‌مدت بلند	تغییرات شماره کپی . . . . . Copy number alteration
Memory (LSTM)	تکامل تومور . . . . . Tumor Evolution
State . . . . . حالت	تکثیر . . . . . Proliferative
Maximum parsimony . . . . . حداکثر زوجیت	تکرار . . . . . Iteration
Deletion . . . . . حذف	تمایز . . . . . Differentiation
Dropout . . . . . حذف تصادفی	تنوع . . . . . Variant
Copy number variation . . . . . حذف و تغییر تعداد کپی (CNV)	توالی‌یابی انبوه . . . . . Bulk sequencing
Single nucleotide . . . . . حذف و تغییر تک‌نوکلئوتیدها	توالی‌یابی تک سلولی . . . . . Single cell sequencing
variants (SNV)	توالی‌یابی نسل بعدی . . . . . Next generation sequencing
خ	تومور . . . . . Tumor
Melanocytic nevi . . . . . خال ملانوسیتیک	تومور بالقوه بدخیم . . . . . Potentially malignant tumor
Loss . . . . . خطا	تومور بدخیم . . . . . Malignant tumor
Limited loss . . . . . خطای محدود	تیمین . . . . . Thymine
Read . . . . . خوانش	ج
Overlapping read . . . . . خوانش همپوشانی	جسمی . . . . . somatic
K-Centroids . . . . . خوشه‌بندی k هسته‌ای	



Dollo parsimony . . . . . زوجیت دولو	د
Subtrees . . . . . زیردرخت‌ها	
Divergent subclones . . . . . زیرکلون‌های واگرا	
	Bulk sequencing data . . . . . داده توالی‌یابی انبوه
	Cytotoxic . . . . . داروهای سمیت سلولی
	Insertion . . . . . درج
	Phylogeny tree . . . . . درخت فیلوژنی
ژ	Likelihood . . . . . درست‌نمایی
	Gate . . . . . دروازه
ژنوم جوانه‌زنی	Hierarchical clustering . . . . . دسته‌بندی سلسله‌مراتبی
	Accuracy . . . . . دقت
	Single nucleotide variant . . . . . دگرگونی تک‌هسته‌ای (SNV)
س	Structural variant . . . . . دگرگونی‌های ساختاری
Cancer . . . . . سرطان	Detected variants . . . . . دگرگونی‌های شناسایی‌شده
Thrombocythemia . . . . . سرطان خون	Binary . . . . . دودویی‌ها
Carcinoma In Situ . . . . . سرطان در محل	epoch . . . . . دوره
Breast cancer . . . . . سرطان سینه	Dollo . . . . . دولو
Infiltrating immune cell . . . . . سلول ایمنی نفوذی	DNA . . . . . دی‌ان‌ای
Cytosine . . . . . سیتوزین	
Sigmoid . . . . . سیگماید	ر
	Radiologist . . . . . رادیولوژیست
ش	Master-slave . . . . . راهبر-پیرو
	Resolution . . . . . رزولوشن
شاخه-مرز	Regression . . . . . رگرسیون
شبکه عصبی بازگشتی	Angiogenic . . . . . رگ‌زایی
Network	Regularization . . . . . رگولاریزیشن
شبکه عصبی کاملاً متصل پیش‌خور	Colorectal carcinoma . . . . . روده بزرگ
connected feed forward neural network	Data driven . . . . . رویکرد مبتنی بر داده
Convolutional neural . . . . . شبکه عصبی کانولوشنی	ز
network	
شدت	Markov Chain . . . . . زنجیره مارکوف مونت کارلو
شناسایی	Monte Carlo (MCMC)

ط لایه‌گذاری صفر . . . . . Zero padding

طبقه‌بندی . . . . . Classification

## م

ف ماتریس فیلوژنی کامل . . . . . Perfect phylogenetic matrix

مماوگرام . . . . . Mammogram

متاستاز . . . . . Metastases

مثبت کاذب . . . . . False positive

مجموع . . . . . Total

محلی . . . . . Local

مدالیت تصویربرداری . . . . . Imaging modality

مرگ . . . . . Death

مکان هندسی . . . . . Loci

مکان‌های بی‌نهایت . . . . . Infinite sites

منفی کاذب . . . . . False negative

موازی . . . . . Antiparallel

## ق

قطعه‌بندی . . . . . Segmentation

## ک

کانال . . . . . Channel

ن کانولوشن . . . . . Convolution

کرنل . . . . . Kernel

کلونی . . . . . Clonal

کواریانس . . . . . Covariance

ناحیه ادراک محلی . . . . . Local receptive field

ناهمگن . . . . . Heterogenetic

ناهمگنی تومور . . . . . Tumor heterogeneity

نرمال‌سازی دسته . . . . . Batch normalization

نشانگر زیستی . . . . . Biomarker

نقشه ویژگی . . . . . Feature map

نگاشت . . . . . Mapping

نورون . . . . . Neuron

نوکلئوتید . . . . . Nucleotid

## گ

گسترش متاستاتیک . . . . . Metastatic spread

گلیوبلاستوم . . . . . Glioblastomas

گوانین . . . . . Guanine

## ل

## و

لایه . . . . . Layer

وزن شبکه . . . . . Network weight

## ششم

## ه

Hyperbolic tangent . . . . . هایپربولیک  
Prune and reattach . . . . . هرس و اتصال دوباره  
Alignment . . . . . هم‌ترازی  
Artificial Intelligence (AI) . . . . . هوش مصنوعی

## ی

Reinforcement learning . . . . . یادگیری تقویتی  
Deep learning . . . . . یادگیری عمیق  
Machine learning . . . . . یادگیری ماشین  
Uniform . . . . . یکنواخت

## L

Leave one out cross validation . . LOOCV



# واژه‌نامه انگلیسی به فارسی

**A** داده توالی‌یابی انبوه . . . . Bulk sequencing data

**C**

Cancer . . . . . سرطان  
Carcinoma In Situ . . . . . سرطان در محل  
Channel . . . . . کانال  
Classification . . . . . طبقه‌بندی  
Clonal . . . . . کلونی  
Colorectal carcinoma . . . . . روده بزرگ  
Computer Vision . . . . . بینایی کامپیوتر  
Computer-aided . . . . . تشخیص به کمک کامپیوتر  
diagnosis

Concatenation . . . . . الحاق  
concordant gain of . . . . . افزایش همزمان جهش‌ها  
mutations

conflict . . . . . ابهام  
Consensus . . . . . اجماعی  
Consistency . . . . . پایداری  
Convolution . . . . . کانولوشن  
Convolutional neural . . . . . شبکه عصبی کانولوشنی  
network  
Copy number alteration . . . . . تغییرات شماره کپی  
Copy number variation . . . . . حذف و تغییر تعداد کپی  
(CNV)  
Cost function . . . . . تابع هزینه  
Covariance . . . . . کواریانس

Accuracy . . . . . دقت  
Activation . . . . . فعال‌سازی  
Activation function . . . . . تابع فعالیت  
Adenine . . . . . آدنین  
Alignment . . . . . هم‌ترازی  
Allele dropout . . . . . ترک آلل  
Allelec dropout . . . . . ترک آلی  
Angiogenic . . . . . رگ‌زایی  
Antiparallel . . . . . موازی  
Artificial Intelligence (AI) . . . . . هوش مصنوعی  
Ascii . . . . . اسکی

**B**

Back-propagation . . . . . بک‌پروپگیشن  
Batch normalization . . . . . نرمال‌سازی دسته  
Bayesian . . . . . بیزی  
Bias . . . . . بایاس  
Binary . . . . . دودویی‌ها  
Binary cross-entropy . . . . . آنتروپی متقاطع دودویی  
Biomarker . . . . . نشانگر زیستی  
Biopsy . . . . . بیوپستی  
Branch-bound . . . . . شاخه-مرز  
Breast cancer . . . . . سرطان سینه  
Bulk sequencing . . . . . توالی‌یابی انبوه

<b>G</b>	Cross validation . . . . . ارزیابی‌های متقابل
	Cytosine . . . . . سیتوزین
Gate . . . . . دروازه	Cytotoxic . . . . . داروهای سمیت سلولی
Germline genome . . . . . ژنوم جوانه‌زنی	
Glioblastomas . . . . . گلیوبلاستوم	<b>D</b>
Gold standard . . . . . استاندارد مرجع	Data driven . . . . . رویکرد مبتنی بر داده
Guanine . . . . . گوانین	Death . . . . . مرگ
	Deep learning . . . . . یادگیری عمیق
<b>H</b>	Deletion . . . . . حذف
Heterogenetic . . . . . ناهمگن	Detected variants . . . . . دگرگونی‌های شناسایی شده
Hierarchical clustering . . . . . دسته‌بندی سلسله‌مراتبی	Detection . . . . . آشکارسازی
Hyperbolic tangent . . . . . هایپربولیک	Differentiation . . . . . تمایز
	Dimension . . . . . بعد
<b>I</b>	Divergent subclones . . . . . زیرکلون‌های واگرا
Imaging modality . . . . . مدالیته تصویربرداری	DNA . . . . . دی‌ان‌ای
Indel . . . . . ایندل	Dollo . . . . . دولو
Inference . . . . . استنباط	Dollo parsimony . . . . . زوجیت دولو
Infiltrating immune cell . . . . . سلول ایمنی نفوذی	Driver mutation . . . . . جهش راننده
Infinite sites . . . . . مکان‌های بی‌نهایت	Dropout . . . . . حذف تصادفی
Insertion . . . . . درج	
Intensity . . . . . شدت	<b>E</b>
Internal covariate shift . . . . . تغییر همبستگی داخلی	epoch . . . . . دوره
Iteration . . . . . تکرار	
	<b>F</b>
<b>J</b>	False negative . . . . . منفی کاذب
Joining-based . . . . . اتصال اکتشافی همسایگی	False positive . . . . . مثبت کاذب
heuristic	Feature map . . . . . نقشه ویژگی
	feedback . . . . . بازخورد
	Fully-connected feed forward neural network . . . . . شبکه عصبی کاملاً متصل پیش‌خور

- K**  
 Metastatic spread . . . . . گسترش متاستاتیک  
 Missing data . . . . . اطلاعات از دست رفته  
 Multi-state . . . . . چند حالت  
 K-Centroids . . . . . خوشه‌بندی k هسته‌ای  
 Kernel . . . . . کرنل
- N**  
 Network weight . . . . . وزن شبکه  
 Neuron . . . . . نورون  
 Next generation . . . . . توالی‌یابی نسل بعدی  
 sequencing  
 Nucleotid . . . . . نوکلئوتید
- O**  
 Oncogene . . . . . انکوژن  
 Optimization . . . . . بهینه‌سازی  
 Over-fitting . . . . . بیش‌پردازش  
 Overlapping read . . . . . خوانش همپوشانی
- P**  
 Passenger mutation . . . . . جهش مسافر  
 Pathology . . . . . آسیب‌شناسی  
 Perfect phylogenic . . . . . ماتریس فیلوژنی کامل  
 matrix  
 Perfect phylogeny . . . . . فیلوژنی کامل منحصر بفرد  
 Phylogeny tree . . . . . درخت فیلوژنی  
 Pixel . . . . . پیکسل  
 Point mutation . . . . . جهش نقطه‌ای  
 Potentially malignant . . . . . تومور بالقوه بدخیم  
 tumor  
 Proliferative . . . . . تکثیر  
 Prune and reattach . . . . . هرس و اتصال دوباره
- L**  
 Layer . . . . . لایه  
 Leave one out cross validation . . . . . LOOCV  
 Likelihood . . . . . درست‌نمایی  
 Limited loss . . . . . خطای محدود  
 Local . . . . . محلی  
 Local maxima . . . . . بیشینه‌های محلی  
 Local receptive field . . . . . ناحیه ادراک محلی  
 Loci . . . . . مکان هندسی  
 Long Short Term . . . . . حافظه‌ی کوتاه‌مدت بلند  
 Memory (LSTM)  
 Loss . . . . . خطا  
 loss-supported . . . . . با در نظر گرفتن حذف خطا  
 Low sequence coverage . . . . . پوشش توالی کم
- M**  
 Machine learning . . . . . یادگیری ماشین  
 Malignant tumor . . . . . تومور بدخیم  
 Mammogram . . . . . ماموگرام  
 Mapping . . . . . نگاشت  
 Markov Chain . . . . . زنجیره مارکوف مونت کارلو  
 Monte Carlo (MCMC)  
 Master-slave . . . . . راهبر-پیرو  
 Maximum parsimony . . . . . حداکثر زوجیت  
 Maximum-likelihood . . . . . بیشینه درست‌نمایی  
 Melanocytic nevi . . . . . خال ملانوسیتیک  
 Metastases . . . . . متاستاز

<b>R</b>	تغییرات بدنی تک نوکلئوتیدی . . . Somatic single nucleotide variation
Radiologist . . . . .	رادیولوژیست
Read . . . . .	خوانش
Read counts . . . . .	تعداد خوانش‌ها
Read coverage . . . . .	پوشش خوانش
Recognition . . . . .	شناسایی
Rectified linear unit . . . . .	تابع واحد اصلاح شده خطی (ReLU)
Recurrent Neural Network . . . . .	شبکه عصبی بازگشتی
Regression . . . . .	رگرسیون
Regularization . . . . .	رگولاریزیشن
Reinforcement learning . . . . .	یادگیری تقویتی
Replication . . . . .	تکثیر
Resolution . . . . .	رزولوشن
reversion . . . . .	برگشت
<b>S</b>	
Segmentation . . . . .	قطعه‌بندی
Sigmoid . . . . .	سیگموید
Single cell sequencing . . . . .	توالی‌یابی تک سلولی
Single nucleotide mutation . . . . .	جهش تک نوکلئوتیدی
Single nucleotide variant (SNV) . . . . .	دگرگونی تک هسته‌ای (SNV)
Single nucleotide variants (SNV) . . . . .	حذف و تغییر تک نوکلئوتیدها
Single Somatic Mutation . . . . .	جهش ساده بدنی
Single variant . . . . .	تغییر ساختاری
somatic . . . . .	جسمی
Somatic mutation . . . . .	جهش بدنی
<b>T</b>	
Thrombocythemia . . . . .	سرطان خون
Thymine . . . . .	تیمین
Total . . . . .	مجموع
Tumor . . . . .	تومور
Tumor Evolution . . . . .	تکامل تومور
Tumor heterogeneity . . . . .	ناهمگنی تومور
Tumor progression . . . . .	پیشرفت تومور
<b>U</b>	
Uniform . . . . .	یکنواخت
Uterine fibroid . . . . .	فیروئید رحمی
<b>V</b>	
Variant . . . . .	تنوع
Variant allele frequency . . . . .	فراوانی تغییرات آلل



## X

اشعه ایکس ..... X-ray

## Z

لایه‌گذاری صفر ..... Zero padding



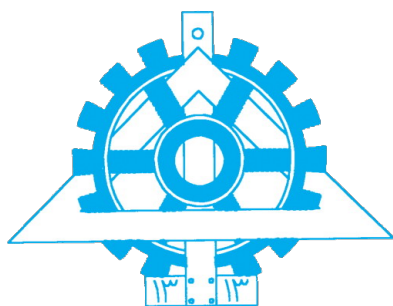


## **Abstract**

This thesis studies on writing projects, theses and dissertations using tehran-thesis class. It ...

**Keywords** SNV, CNV, Phylogenetic, Tree, Q-learning, Deep learning





University of Tehran  
College of Engineering  
Faculty of New Science and  
Technology  
Network



# **Inference of Phylogenetic Tree for Inter Tumor using Single Cell Mutations and CNV**

A Thesis submitted to the Graduate Studies Office  
In partial fulfillment of the requirements for  
The degree of Master of Science  
in Information Technology - Network Science

By:

**Afshin Bozorgpour**

Supervisors:

**Dr. Saman Haratizadeh and Dr. Abolfazl Motahari**

Jul 2021