

**Article:** EMu: probabilistic inference of mutational processes and their localization in the cancer genome

**Journal:** BMC genome biology

**Authors:** Fischer, A. et al.

**Year:** 2013

**Citations:** 42

**link:** <https://genomebiology.biomedcentral.com/articles/10.1186/gb-2013-14-4-r39>

هدف:

- استخراج فرایندهای جهش‌زا ابتدایی و امضای آن‌ها از داده‌های جهش سرطانی کل ژنوم

- محدود کردن فعالیت فرایند به یک ناحیه در ژنوم‌های سرطان

انواع مختلفی از فرایندهای فیزیکی، شیمیایی و بیولوژیکی و همچنین غیرفعال شدن مکانیسم تعمیر mismatch منجر به جهش در سرطان می‌شود. ما در اینجا وظیفه شناسایی فرایندهای جهش‌زا فعال را در طول توسعه سرطان در نظر می‌گیریم. ما فرض می‌کنیم که این فرایندها مجزا هستند، به این ترتیب، هر فرایند علامت مشخصه متفاوتی را بر روی ژنوم سرطان باقی می‌گذارد. با توجه به داده‌های توالی از تومورهای متعدد، که در آن فرایندهای جهش‌زا در ابعاد مختلف فعال هستند، هدف ما استنباط تعداد فرایندهای جهش‌زا اولیه، امضای آن‌ها و سهم هر فرایند تا طیف جهش‌های مشاهده شده در هر تومور است. شناسایی فرایندهای جهش‌زا گامی در جهت شناخت مکانیسم‌های ایجادکننده سرطان است.

محدودیت روش‌های قبلی:

1. فرصت وقوع جهش‌ها در یک توالی ژنتیکی مشخص صریحاً مورد توجه قرار نگرفته است. این مهم است زیرا نتیجه مشاهده شده از یک فرایند جهش‌زا به محتویات توالی‌ای که در آن عمل می‌شود بستگی دارد. به عنوان مثال، فرایندی که انتقال‌های  $C > T$  را در محل‌های CpG تولید می‌کند، ممکن است با فعالیت یکنواخت در سراسر ژنوم انجام دهد، اما بیشتر در مناطقی با چگالی بیشتر از این دونوکلوئوتیدها مشاهده می‌شود.

۲. در حالی که با توجه به فرآیندهای جهش‌زا در کل ژنوم، مطالعات قبلی عملکرد فرآیندهای مختلف را در مناطق کوچک‌تر ژنوم را بررسی نکرده اند. **غیریکنواختی در میزان جهش در یک ژنوم سرطان مشاهده شده است، که با تفاوت و bias در عملکرد مکانیسم‌های ترمیم DNA و [5] با توزیع علائم اصلاح histone همراه است [14]**، نشان می‌دهد که فرآیندهای جهش‌زا ممکن است در مناطق مختلف ژنوم **متفاوت عمل کنند**. تغییر در فرصت جهش‌زایی، که در بالا ذکر شد، به ویژه در مناطق کوچک ژنوم از اهمیت ویژه‌ای برخوردار است. Copy number variation، که می‌تواند در سرطان شایع باشد [15]، به طور قابل ملاحظه‌ای فرصت جهش‌زایی در مناطقی از تغییر پلوئیدی (ploidy) (تعداد دسته‌های کروموزوم‌ها) را تغییر می‌دهد، در حالی که مناطق مختلف ژنوم انسان ممکن است از نظر بنیادی در محتوای داخلی توالی آنها متفاوت باشد [16].

۳. روشهای قبلی صراحتاً روش تصادفی که جهش‌ها در ژنوم وارد می‌شوند را در نظر نگرفته اند. تصادفی بودن، همراه با ماهیت گسسته بودن جهش، منجر به الگوهای خاص noise در طیف‌های جهش مشاهده شده می‌شود. هنگامی که زیرمجموعه‌های کوچک‌تر ژنوم مورد تجزیه و تحلیل قرار می‌گیرند، این noise تقویت می‌شود، و این یک راه حل احتمالی صریح از داده‌ها ضروری است.

روش جدید:

براساس expectation-maximization (EM)

برای استنتاج تعداد فرآیندهای جهش‌زا ابتدایی و طیف‌های آنها از داده‌های توالی سرطان.

با امکان در نظر گرفتن bias ها در فرصت جهش‌زایی و داده‌های noise با استفاده از نتایج شبیه‌سازی، ما نشان می‌دهیم که اضافه کردن فرصت جهش‌زایی در شناسایی تعداد صحیح فرآیندهای جهش‌زا در این مدل بسیار مهم است. تعداد فرایندها با بیشترین احتمال با Bayesian information criterion (BIC) می‌تواند بدست بیاید. تخمین خطای معنی دار از امضاهای جهش‌زا می‌تواند به روش تحلیلی یا عددی با روش‌های Markov chain Monte Carlo (MCMC) استنباط شود. با توسعه منطقی الگوریتم یادگیری-پارامتر به زیر مجموعه‌های کوچک‌تر یک ژنوم سرطان، می‌توان فعالیت موضعی فرآیندهای مختلف را اندازه گرفت. **ما از داده تفسیر (annotation) از**

پروژه [20] ENCODE استفاده می کنیم تا تغییرات در فرآیندهای جهش را مربوط به حالت chromatin را بررسی کنیم. پروژه ENCODE بخش هایی از ژنوم انسان را در حالت های مختلف عملکردی کروماتین فراهم می کند [21]، که عمدتاً بر اساس توزیع علائم اصلاح histone است. در سرطان، یک مطالعه جدید ارتباطات قوی جهش های جسمی را با برخی از این تغییرات histone، به ویژه H3K9me3 مشخص کرده است [14]. ما در اینجا شواهد آماری را ارائه می دهیم که نشان می دهد فرآیندهای جهش را در مناطق با حالت chromatin مختلف در سرطان سینه متفاوت عمل می کنند. به طور خاص، یکی از فرآیندهای جهش، تولید کننده جهش های C > T عمدتاً در محل های CpG، در مناطق promoter به شدت کم رنگ است و در مناطق heterochromatin تقویت می شود. این نشان می دهد که این فرایند با متیلاسیون DNA ارتباط دارد، همانطور که در مورد حذف خود به خودی آمینواسید (deamination) وجود دارد [22].

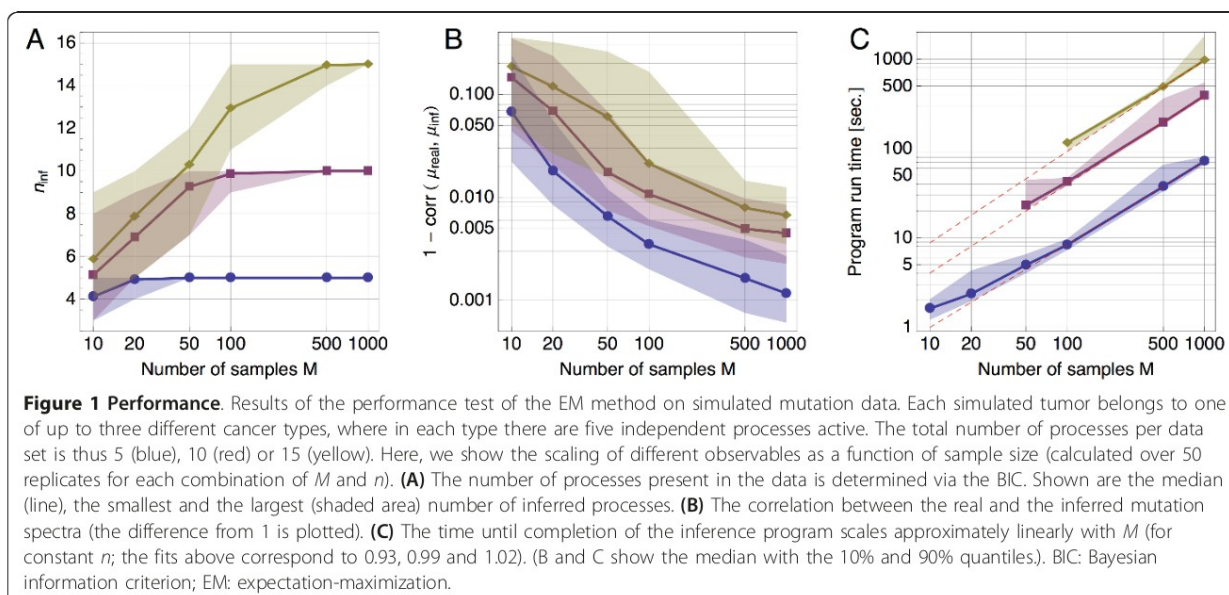
## نتایج

ارزیابی با داده های شبیه سازی شده

تعداد نمونه های سرطانی:  $M \in \{10, 20, 50, 100, 500, 1000\}$

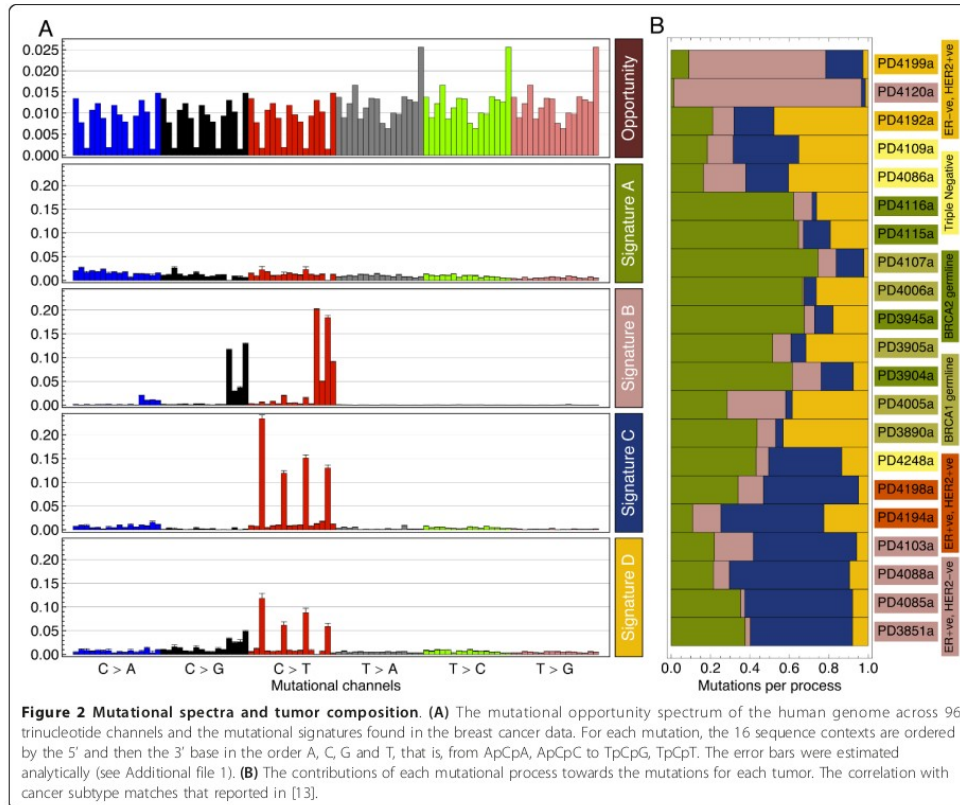
تعداد فرآیندهای جهش را:  $n \in \{5, 10, 15\}$

در هر حالت، افزایش اندازه نمونه، توانایی تمایز بین فرآیندهای جهش را بهبود بخشید. در اکثر موارد، پنج فرایند جهش از بیست نمونه تومور قابل تشخیص است. در فرآیندهای شناسایی شده، طیف جهش با میزان واقعی آن ارتباط قوی دارد. استفاده از BIC در انتخاب مدل، تعداد فرایندها هیچ وقت overestimated نشده است. هرچند با نادیده گرفتن فرصت جهش زایی باعث شد تعداد فرایند با BIC بیشتر از مقدار قابل انتظار شود.



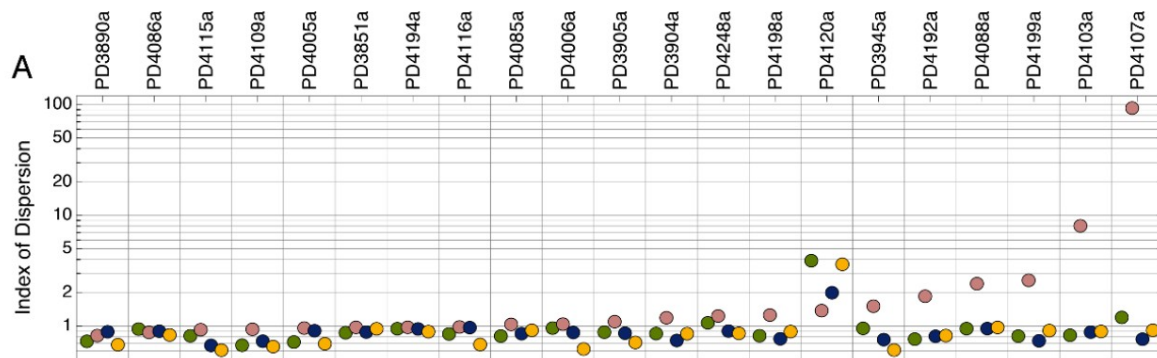
## فرایندهای جهش‌زا در سرطان سینه:

در سراسر ژنوم انسان، فرصت جهش‌زایی بین کانالها به اندازه 17-fold تغییر می‌کند، و این موضوع را به یک عامل مهم در استخراج طیف‌های جهش‌زا تبدیل می‌کند. در مجموعه داده‌های سرطان سینه، این با وجود تغییرات copy number در مقیاس بزرگ همراه است. این منجر به طیفی از فرصت جهش‌زایی می‌شود که به نوع تومور بستگی دارد.



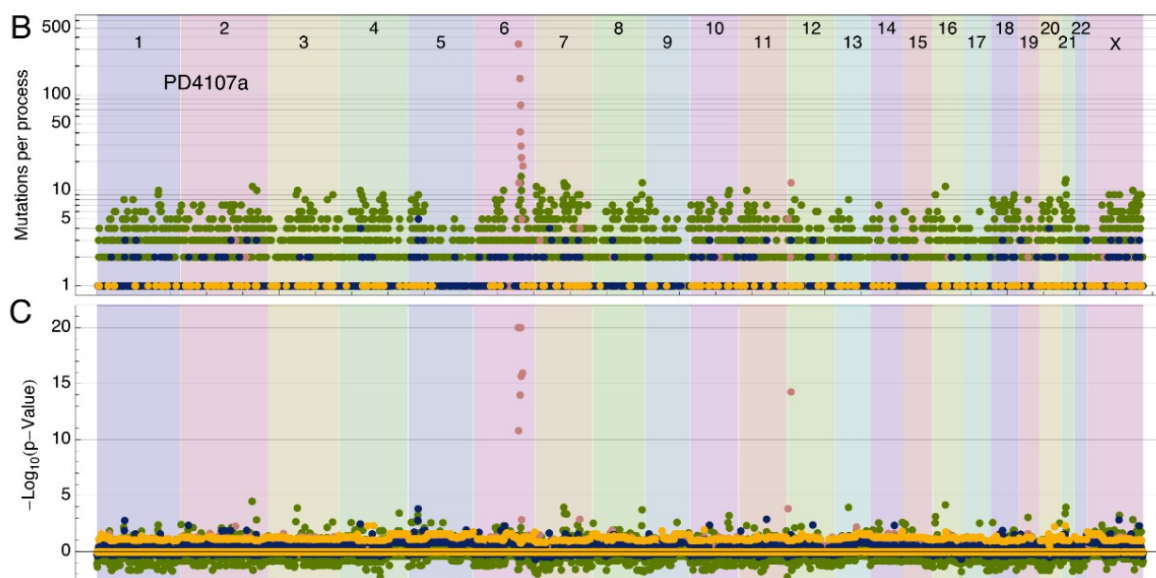
تجزیه و تحلیل موضعی، مناطقی را نشان می‌دهد که به شدت توسط فرایندهای تکی مورد هدف قرار می‌گیرند:

برای موضعی کردن مناطق ژنوم که در آن فرایندهای مختلف جهش‌زایی فعال هستند، از امضاهای جهش‌زابه دست آمده برای تخمین فعالیت هر فرآیند برای هر پنجره توالی 1 Mb استفاده کردیم. تنوع در فعالیت موضعی فرایندهای جهش‌زا تکی در هر ژنوم سرطان به ویژه در فرآیند B مشاهده شد. برای به دست آوردن میزان ناهمگنی جهش‌زایی، برای هر فرآیند در هر نمونه سرطانی یک شاخص پراکندگی (نسبت واریانس به میانگین جهش‌ها در هر پنجره)، معادل فرصت جهش‌زایی، جایی که مقدار 1 مطابق توزیع unbiased و همگن جهش‌ها در سراسر ژنوم است.



خوشه بندی تعداد زیادی جهش در ناحیه کوچکی از ژنوم که kataegis نامیده می‌شود ، قبلاً در ژنوم‌های سرطان سینه در نظر گرفته شده است که در 13 مورد از 21 تومور وجود دارد. در کار قبلی ، رویدادهای kataegis تا حدود زیادی جهش‌های  $C > T$  و  $C > G$  را نشان میدادند. اکنون می بینیم که kataegis به شدت با امضای جهش‌ها B مرتبط است.

یک دید بهتر از تغییرات جهش‌یافته مناطق خوشه بندی جهش را شناسایی کرد. برای هر منطقه 1 مگابایت ژنوم ، تعداد جهش‌های شناسایی شده از هر فرآیند در برابر این فرض صفر ارزیابی شد که هر فرآیند جهش‌های خود را بطور تصادفی در ژنوم سرطان توزیع می‌کند ، با توجه به یک فعالیت یکنواخت در سرتاسر ژنوم ، تنها با فرصت موضعی تنظیم شده. بررسی P-value به عنوان تابعی از موقعیت ژنوم در نمونه تومور PD4107a بسیار واضح است که رویداد kataegis قبلاً در کروموزوم 6 گزارش شده است ، اما علاوه بر این ، یک واقعه کوچکتر اما بسیار مهم در ابتدای کروموزوم 12 پیدا کرده است.



## فرآیندهای جهشی و حالت chromatin:

تجزیه و تحلیل موضعی را می‌توان برای هر زیر مجموعه تعریف شده ژنوم انجام داد. تقسیم بندی ژنوم به حالت‌های chromatin از لحاظ عملکردی مرتبط با پروژه ENCODE مورد توجه ویژه ای است [20,21]. در حقیقت، تراکم منطقه‌ای جهش‌ها در حالت chromatin در human mammary epithelial cells (HMECs) متفاوت است. مقایسه فرصت جهش‌زایی در هر حالت chromatin با جهش‌های واقعی مشاهده شده نشان داد که 7٪ جهش بیشتر از آنچه انتظار می‌رفت در بخش‌های heterochromatin وجود دارد، اما تقریباً 50٪ جهش کمتری در مناطق promoter داشت. با استفاده از روش استنباط موضعی ما حساسیت فرآیندهای جهش‌زا به حالت chromatin را بررسی کردیم. برای این منظور، میزانی که یک فعالیت موضعی فرآیند جهش‌زا به حالت chromatin یک بخش ژنومی بستگی دارد، اندازه‌گیری می‌شود. فرایند جهش‌زا  $C > T$  در محل‌های XpCpG یک تغییر بسیار مهم در فعالیت هر دو منطقه heterochromatin و promoter را نشان می‌دهد. این سیگنال به طور مداوم در تمام 21 سرطان مشاهده می‌شود. برای تعیین میزان این اثر، تعداد جهش‌هایی را که به یک فرآیند خاص در هر حالت کروماتین اختصاص داده شده را با تعداد جهش‌هایی که انتظار می‌رود در صورت عملکرد یکنواخت فرایند در سراسر ژنوم دیده شود، مقایسه کردیم. نسبت این دو عدد را برای هر فرایند در هر حالت کروماتین و برای هر ۲۱ تومور نشان دادیم. در مناطق heterochromatin فعالیت فرایند C در همه تومورها افزایش می‌یابد، به طور متوسط تا ۴۰٪، نسبت به به اندازه آن بخش ژنومی اثر بزرگی دارد. در مناطق promoter، فرایند C به طور مداوم کاهش می‌یابد، به طور متوسط 80٪. این مشاهده این احتمال را ایجاد می‌کند که فرایند C به حالت DNA methylation وابسته باشد. فرایند C عمدتاً از جهش‌های  $C > T$  در محل‌های CpG تشکیل شده است. چنین جهش‌هایی مشاهده شده است که بیشتر در سایت‌های CpG هنگام متیله شدن cytosine اتفاق می‌افتد، در نتیجه یک deamination خود به خودی [22,23]. در مناطق promoter ژن‌های فعال، تعداد محل‌های CpG بیشتر از حد معمول وجود دارد (افزایش فرصت برای فرایند C). اما به طور عمده آنها متیله نشدند (کاهش فعالیت مشاهده شده فرایند C).

پیشنهادهای:

تفسیر حالت کروماتین که در اینجا مورد استفاده قرار گرفت ، از آزمایشات انجام شده بر روی بافت سالم بدست آمد. یک توسعه جالب برای این کار ، بررسی ارتباط بین فرآیندهای جهش‌زا و تفسیر عملکردی و اپی‌ژنتیکی ژنوم‌های سرطانی واقعی است که در آنها یافت شده است. این نوع داده‌ها هنوز در دسترس نیست اما مطمئناً در دست فناوری‌های فعلی است.

روش استنباط ارائه شده در این مطالعه محدود به ۹۶ کانال جهش سه‌نوکلئوتیدی نمی باشد. علاوه بر این، می‌توان اطلاعاتی از قبیل رشته رونویسی بودن جهش‌ها یا اینکه آیا جهش دیرتر یا زودتر در زمان سرطان رخ داده است [19].

به دلیل بازآرایی‌های کروموزومی (chromosomal rearrangements) و تغییرات در تعداد نسخه‌ها (copy number changes) در مقیاس بزرگ در طی تکامل جسمی ، این فرصت جهش‌زایی در واقع یک مقدار پویا است [19]. به عنوان یک تقریب ، ما برای محاسبه فرصت جهش‌زایی از copy number state ژنوم‌های سرطان در زمان توالی‌یابی استفاده می‌کنیم. به عنوان مثال ، ممکن است در اواخر مرحله تکامل سرطان یک واقعه کپی بر روی یک نسخه والدین از یک کروموزوم اتفاق افتاده باشد. بنابراین فرصت جهش‌زایی برای این قطعه کپی 50٪- (از دو تا سه نسخه) برای زمان باقی مانده تا توالی‌یابی افزایش می‌یابد. در این کار ، ما فرض کردیم که این فرصت افزایش جهش‌زایی در همه فرآیندهای جهش‌زا در سراسر ژنوم وجود دارد. با اطلاعات دقیق‌تر در مورد زمان‌بندی نسبی تغییرات در تعداد نسخه‌ها ، می‌توان یک فرصت جهش‌زایی مؤثر را ایجاد کرد.

کاربردهای آینده این روش شامل فهرست‌نویسی امضاهای جهش‌زا در انواع مختلف سرطان ، شناسایی فرآیندهای جهش‌زا مشترک و خاص از نوع سرطان (و از این رو ، طبقه‌بندی بهتر سرطان‌ها براساس ترکیب فرآیند آنها) و سرانجام یک مدل ابتدایی وابسته- به نوع - سرطان (cancer-type-dependent) برای شناسایی انواع گذرگر (driver) مسبب سرطان، از طریق سیگنال‌های انتخاب ، می‌شود.



## داده جهش:

183,916 جهش جسمی نقطه‌ای

M = ۲۱ نمونه سرطان سینه

N<sub>c</sub> = ۹۶ کانال جهش سه نوکلئوتیدی

۶ جفت تغییر نوکلئوتیدی ممکن، C:G > A:T, C:G > G:C, C:G > T:A, T:A > A:T, T:A > C:G, T:A > G:C

به همراه ۱۶ حالت مختلف از دو نوکلئوتید مجاور هر جهش

← مجموعه مشاهدات:

X<sub>j</sub><sup>m</sup> (j = 1, ..., N<sub>c</sub>, m = 1, ..., M) of channel- and tumor-specific mutation counts

## فرایندهای ابتدایی جهش‌زا:

n فرایند جهش‌زا ابتدایی مجزا

امضاهای جهش‌زا:

{μ<sub>aj</sub>} (a = 1, ..., n, j = 1, ..., N<sub>c</sub>), where  $\sum_{j=1}^{N_c} \mu_{aj} = 1$

x<sub>a</sub><sup>m</sup> فعالیت فرایند a در تومور m، میزان عملکرد یک فرایند در تومور مشخص (سهم

فرایند در تومور)

توجه به این نکته ضروری است که استخراج میزان جهش واقعی یک فرایند امکان پذیر نیست. مجموع فعالیت یک فرایند در یک سرطان معین، **حالت ضرب میزان جهش آن، زمانی که فرایند فعال بوده است و نسبت ژنومی که در واقع مورد حمله قرار گرفته است.** در حالی که این مورد آخر را می توان با یک تجزیه و تحلیل مکانی حل کرد، دو مورد اول بدون توالی‌یابی time-resolved تومور غیرقابل حل است.

با توجه به مطالب فوق، با استفاده از توزیع پواسون نتیجه ایجاد تصادفی جهش‌ها در ژنوم را توصیف می‌کنیم. با فرض اینکه فرایندهای جهش‌زا مستقل از یکدیگر باشند، احتمال مشاهده بردار x<sup>m</sup> از تعداد جهش در تومور m در کانال j:

$$P(X^m | x^m, \omega^m, \mu) \equiv \prod_{j=1}^{N_c} \text{Pois} \left( X_j^m \mid \sum_{a=1}^n x_a^m \mu_{aj} \omega_j^m \right), \quad (1)$$

ω<sub>j</sub><sup>m</sup> فرصت جهش‌ها در کانال j که در تومور m رخ داده است

در داده‌های مورد بررسی، تنوع تعداد کپی (copy number variation) در فرصت جهش‌زایی بین تومورها متفاوت است. ما در اینجا در اطلاعات ploidy موجود [19] به



شرح زیر استدلال کرده ایم: هر جفت نوکلئوتید بسته به محتوای توالی مجاور نوکلئوتید جهش یافته می‌تواند به سه روش مختلف متناظر با سه کانال 96 جهش‌زا، جهش یابد. بردار مجموع فرصت جهش‌زایی  $\omega^m$ ، مجموع تجمعی در کل نوکلئوتیدها است، که در آن هر نوکلئوتید سه کانال نه با یک بلکه با تعداد کپی آن نوکلئوتید در ژنوم  $m$  افزایش می‌یابد. در روشی که فرصت جهش‌زایی را نادیده گرفته است، مقادیر  $\omega^m$  به 1 تنظیم شده است. توجه داشته باشیم که، از پارامترهای موجود در معادله بالا، فعالیت‌های  $x$  و طیف  $\mu$  ناشناخته هستند. تخمین حداکثر احتمال این مقادیر با استفاده از الگوریتم EM انجام شد.

### پیدا کردن طیف‌های جهش‌زایی (امضاهای جهش‌زا):

الگوریتم EM تخمین حداکثر درست‌نمایی (ML) پارامترهای نامعلوم مدل (در اینجا طیف  $\mu$ ) و داده‌های پنهان (در اینجا فعالیت  $x$ ) را می‌یابد. پس از یک حدس اولیه،  $\mu(0)$  ساخته شد، دو مرحله به صورت iterative تکرار می‌شوند. در مرحله اول، با توجه به بهترین حدس فعلی،  $\mu(k)$  و داده‌های مشاهده شده،  $X$ ، تخمینی برای  $x$  یافت می‌شود. در مرحله دوم، این برآورد،  $\hat{x}$  از داده‌های پنهان برای به دست آوردن برآورد پارامتر به روز شده،  $\mu(k+1)$  استفاده می‌شود. این دو مرحله تا زمانی رسیدن همگرایی به حداکثر (محلی) درست‌نمایی داده،  $P(X|\mu)$ ، تکرار می‌شود [17]. به صورت دقیق‌تر روال به شرح زیر است:

0. (Initialize) باتوجه به نرمال سازی  $\sum_{j=1}^{N_c} \mu_{aj} = 1$  مقداردهی اولیه  $\mu$

1. (E-step) باتوجه به تخمین حاضر،  $\mu^k$ ، برآورد حداکثر درست‌نمایی را برای همه فعالیت‌های پنهان پیدا کنید:

$$\hat{x}^m \equiv \operatorname{argmax}_{x \in \mathbb{R}_+^n} \log P(X^m | x, \omega^m, \mu^{(k)}) = \operatorname{argmin}_{x \in \mathbb{R}_+^n} \sum_{j=1}^{N_c} \left[ \sum_{a=1}^n x_a \mu_{aj}^{(k)} \omega_j^m - X_j^m \log \left( \sum_{a=1}^n x_a \mu_{aj}^{(k)} \omega_j^m \right) \right]$$

2. (M-step) با استفاده از تخمین‌های داده‌های پنهان، طیف جهش را به روزرسانی

$$\mu^{(k+1)} \equiv \operatorname{argmax}_{\mu \in \mathbb{R}^{n \times N_c}} \sum_{m=1}^M \log P(X^m | \hat{x}^m, \omega^m, \mu), \text{ with constraint } \sum_{j=1}^{N_c} \mu_{aj}^{(k+1)} = 1.$$

کنید:

3. (Finish) تست برای همگرایی. در صورت نیاز به مرحله E-step و تکرار برگردید.

### محاسبه مجموع داده‌های درست‌نمایی:

پس از همگرایی الگوریتم EM به حداکثر تخمین درست‌نمایی،  $\hat{\mu}$ ، از  $n$  طیف، مقدار عددی درست‌نمایی با ادغام متغیرهای نهفته،  $x$  محاسبه شد. برای جلوگیری از هزینه ادغام عددی، این کار با استفاده از تقریب نقطه زینی انجام شد:

$$\log P(X|\hat{\mu}) = \sum_{m=1}^M \log \int d^n x P(X^m|x, \omega^m, \hat{\mu}) \approx \sum_{m=1}^M \left[ \frac{n}{2} \log 2\pi - L(\hat{x}^m) - \frac{1}{2} \log \det H(L)(\hat{x}^m) \right]$$

که در آن  $L(x) = -\log P(X|x, \omega, \mu)$  لگاریتم درست‌نمایی داده شرطی است و  $H(L)(x)$  ماتریس Hessian مشتق دوم آن است.

### پیدا کردن تعداد فرایندهای جهش‌زا:

مقایسه درست‌نمایی داده‌های به دست آمده با مقادیر مختلف  $n$ ، محتمل‌ترین تعداد فرایندهای جهش‌زا را نشان می‌دهد. با افزایش مقدار  $n$  تعداد پارامترهای مدل موجود برای fitting افزایش می‌یابد، و به طور کلی توضیحی بهتر برای داده‌ها،  $X$  ایجاد می‌کند. این به معنای افزایش ارزش درست‌نمایی،  $P(X|\hat{\mu})$  است. برای جلوگیری از overfitting داده‌ها، از Bayesian Information Criterion (BIC) برای اصلاح پیچیدگی مدل استفاده شد. در نهایت مدل با بیشترین مقدار BIC انتخاب شد:

$$BIC = 2 \log P(X|\hat{\mu}) - n(N_c - 1) \log M.$$

### برآورد فعالیت‌های فرآیند و اختصاص جهش‌ها به فرآیندها:

#### تجزیه و تحلیل سراسری:

تکمیل الگوریتم EM تخمین‌های حداکثر درست‌نمایی،  $\hat{\mu}$  از امضاها، جهش‌زا  $n$  فرآیند را فراهم می‌کند. به همین ترتیب، E-step تخمین‌های،  $\{\hat{x}_a^m\}$ ، از فعالیت‌های فرآیند برای هر نمونه را بدست می‌آورد. در مجموع، می‌توان از این برای تخمین تعداد جهش‌هایی استفاده کرد  $\hat{X}^m = \sum_{a=1}^n \hat{x}_a^m$ . تمایل زیاد توسط هر فرآیند در مجموع همه

جهش‌های مشاهده شده شرکت می‌کند: . برای بهبود برآورد فعالیت ، یک تخمین ساده از حداکثر احتمال لگاریتم درست‌نمایی از

$$\hat{x}^m \equiv \operatorname{argmax}_{x \in \mathbb{R}_+^n} \log P(X^m | x, \omega^m, \mu^{(k)}) = \operatorname{argmin}_{x \in \mathbb{R}_+^n} \sum_{j=1}^{N_c} \left[ \sum_{a=1}^n x_a \mu_{aj}^{(k)} \omega_j^m - X_j^m \log \left( \sum_{a=1}^n x_a \mu_{aj}^{(k)} \omega_j^m \right) \right]$$

را تعمیم می‌دهیم تا اطلاعات پیشین را به صورت pseudocounts درج کنیم. همانطور که در زیر نشان می‌دهیم ، این pseudocounts به طور کلی در کانال‌ها یکنواخت نیستند. فعالیت تخمین زده می‌شود ،  $\{\hat{x}_a^m\}$  ، و دارای وضعیت بیشینه احتمال پسین است. سرانجام ، تعداد جهش ها ،  $\hat{x}_a^m$  ، در سرطان m که از طریق فرایند a در کانال z منتشر شده‌اند به شرح زیر است:

$$\hat{x}_{a,j}^{m,g} = \frac{|X^m| \hat{x}_a^{m,g} \hat{\mu}_{a,j} \omega_j^m}{\sum_{b=1}^n \hat{x}_b^{m,g} (\hat{\mu} \omega^m)_b}, \text{ where } |X^m| \equiv \sum_{j=1}^{N_c} X_j^m, \text{ and } (\hat{\mu} \omega^m)_a \equiv \sum_{j=1}^{N_c} \mu_{aj} \omega_j^m.$$

در اینجا ، اندیس بالا g مقدار سراسری بدست آمده را نشان می‌دهد. قاعده تخصیص را می‌توان به شرح زیر درک کرد: کسری از جهش‌های مشاهده شده که توسط یک فرآیند خاص تولید شده‌اند باید متناسب با فعالیت آن ، طیف جهش آن و فرصت جهش در آن کانال باشد. به هر مشاهده ،  $X^m$  ، در مجموع n تعداد pseudocounts اضافه شد. این شمارش دقیقاً به روش فوق انجام شد ، اما بدون هیچ گونه bias ناشی از فعالیت‌های (هنوز ناشناخته):

$$\tilde{X}_{aj}^m \equiv \frac{n \hat{\mu}_{aj} \omega_j^m}{\sum_{b=1}^n (\hat{\mu} \omega^m)_b}$$

بنابراین ، آنها نمایانگر مشاهده‌ای هستند که فقط براساس اطلاعات شناخته شده پیشین از برآورد فعالیت انجام می‌شود. pseudocounts با یک عبارت اضافی وارد تابع لگاریتم درست‌نمایی می‌شود:

$$L(x; X^m, \omega^m, \mu) \rightarrow L(x; X^m, \omega^m, \mu) - \sum_{j=1}^{N_c} \sum_{a=1}^n \tilde{X}_{aj}^m \log(x_a \mu_{aj} \omega_j^m)$$

### تجزیه و تحلیل موضعی:

با همان منطق مشابه در بالا ، ما تخمین‌هایی از فعالیت موضعی هر فرآیند در هر سرطان گرفتیم و هر ژنوم را به پنجره های غیر همپوشانی با طول 1Mb تقسیم می‌کنیم. با این حال ، ما اکنون می‌توانیم از فعالیت‌های سراسری قبلاً محاسبه شده ،  $\{x_a^{m,g}\}$  ، به عنوان اطلاعات پیشین آگاه‌تر برای استنباط فعالیت‌های فرآیند در هر پنجره megabase ، به طول 1 استفاده کنیم:

$$\tilde{X}_{aj}^{m,l} = \frac{n \hat{x}_a^{m,g} \hat{\mu}_{aj} \omega_j^{m,l}}{\sum_{b=1}^n \hat{x}_a^{m,g} (\hat{\mu} \omega^{m,l})_b} \xrightarrow[\text{using } X^{m,l}, \omega^{m,l}]{\text{E - step with eq. 8}} \hat{x}_a^{m,l}.$$