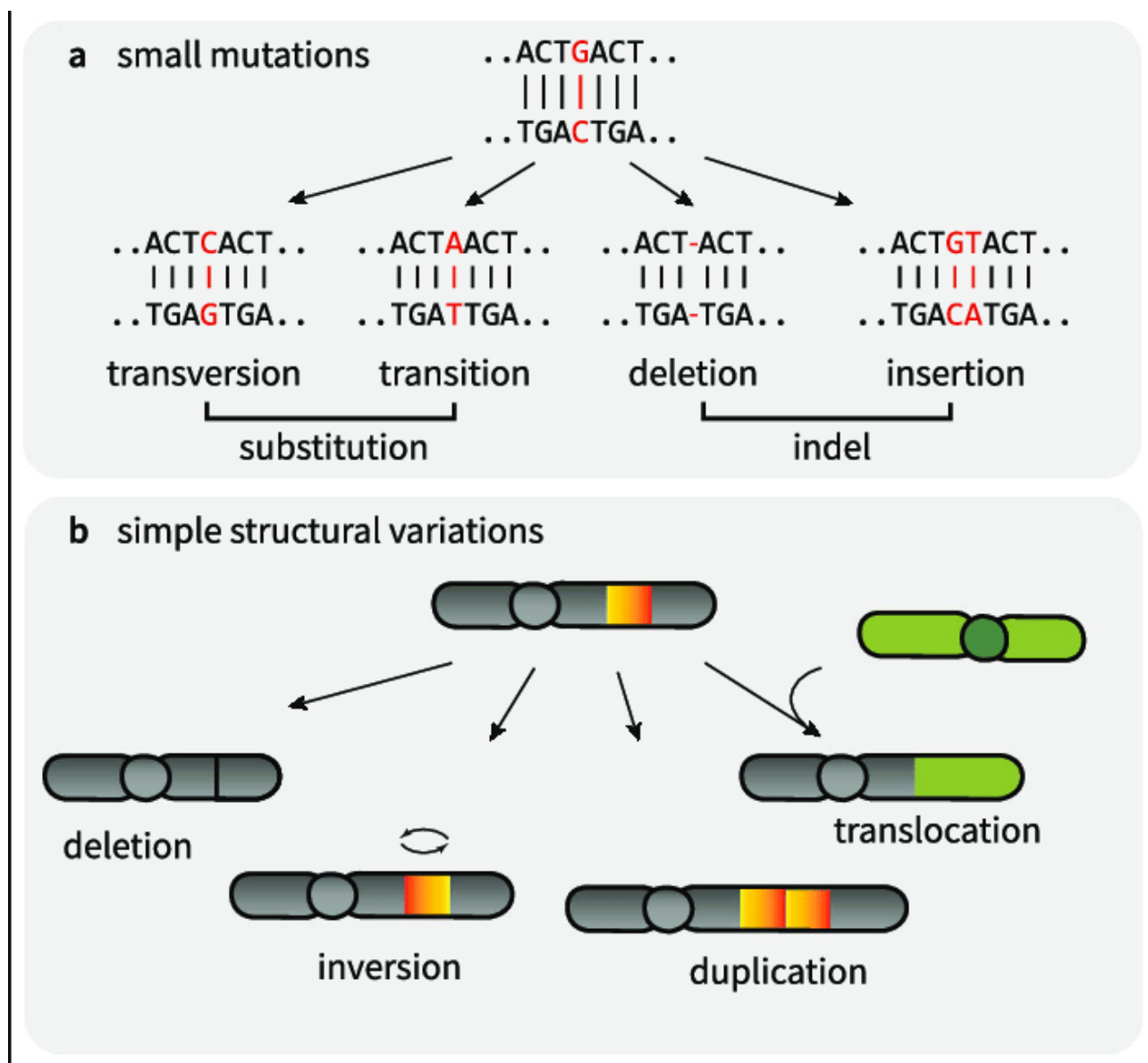


جهش های جسمی معمولاً به چهار کلاس گروه بندی می شوند. تعویض نوکلئوتید ، ایندل های کوچک ، بازآرایی مجدد و تغییر تعداد کپی ها.



امضای یک فرآیند جهش به عنوان یک تابع چگالی احتمال گسسته با دامنه ای از ویژگی های جهش از پیش انتخاب شده نشان داده می شود.

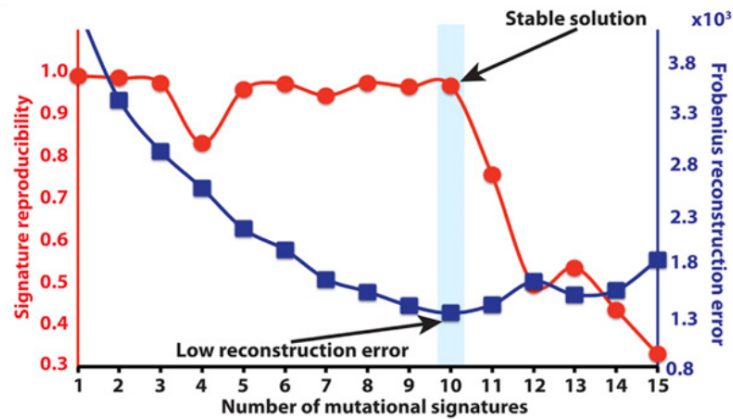
به طور ریاضی، می‌توان ویژگی‌های جهش را به صورت مجموعه الفبای متناهی Ξ با K حرف (هر حرف متناظر با یک ویژگی جهش) بیان کرد و طبق تعریف، فرایند جهش P_1 یک K -tuple با ترتیب lexicographical است؛ $P_1 = [p_1^1, p_1^2, \dots, p_1^K]^T$ ، که p_1^i احتمال فرایند p_1 که ویژگی جهش متناظر با i ام حرف از الفبای Ξ باعث می‌شود را نشان می‌دهد و چون p^i احتمال هستند:

$$p = [p_1^1, p_1^2, \dots, p_1^K]^T, \sum_{k=1}^K p_1^k = 1$$

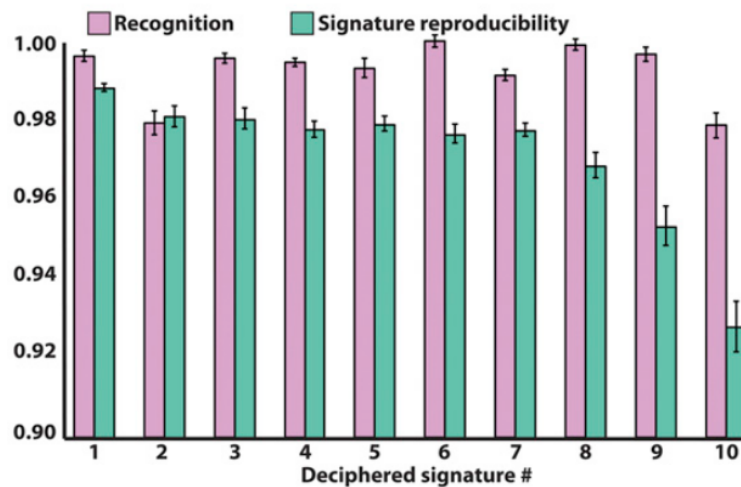
ژنوم‌های سرطانی متفاوت می‌توانند در معرض فرایند جهش‌زا با شدت‌های متفاوت قرار گیرند. به عنوان مثال، یک فرایند جهش‌زا می‌تواند باعث ایجاد ۱۰۰۰ جهش در یک ژنوم سرطان شود در حالی که باعث ایجاد ۲۰,۰۰۰ جهش در دیگری می‌شود. بنابراین یک فرایند جهش‌زا با امضای P_1 دارای شدت (تعداد جهش‌هایی که سبب می‌شود)، e_g^1 ، در ژنوم سرطانی g است. توجه داشته باشید که اندیس پایین امضا P_1 با اندیس بالا e_g^1 مطابقت دارد. بنابراین شدت e_g^1 با امضا P_1 مطابقت دارد.

روش NMF بر روی ۱۰۰ فهرست جهش‌زا ژنوم سرطانی شبیه‌سازی شده پیاده‌سازی شده است. مشابه بسیاری از ژنوم‌های سرطانی انسان، هر ژنوم سرطانی شامل ۵۰۰ تا ۵۰۰۰۰ جهش جایگزینی است. جهش‌های شبیه‌سازی شده با استفاده از ۱۰ فرایند جهش‌زا با امضا مجزا هر کدام با ۹۶ نوع جهش (۶ نوع جهش جایگزین و دو همسایه مجاور) تولید شده‌اند.

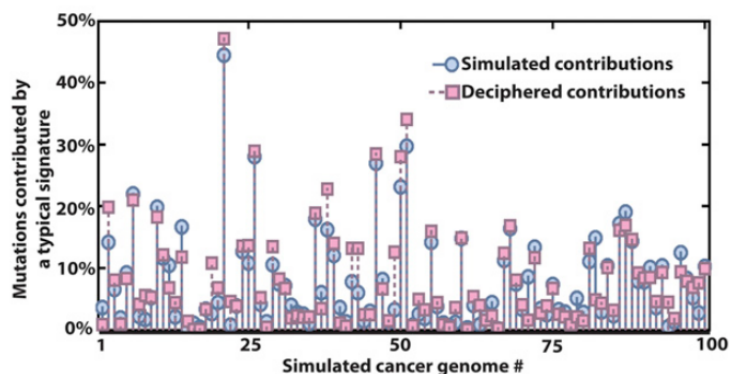
شناسایی تعداد، N ، فرایندهای جهش‌زا در مجموعه ژنوم‌های سرطانی به یک دانش پیشین را برای کشف امضاهایشان نیاز دارد. در روش NMF مقدار N را با استفاده از مقادیر مختلف شناسایی می‌کنیم. برای هر N ، شباهت بین فرایندهای به دست آمده در هر تکرار (قابلیت بازتولید فرایند) را ارزیابی می‌کنیم. علاوه بر این، برای هر N ، متوسط خطای بازسازی (reconstruction error) فروبنیوس از متوسط امضاهای کشف شده \bar{P} و شدت آن‌ها \bar{E} ($\|M - \bar{P}\bar{E}\|_F^2$) تعیین می‌شود. خطای کم بازسازی نشانگر توصیف دقیق از فهرست‌های ژنوم سرطان اصلی است. مقدار N ای را انتخاب می‌کنیم که فرایندهای به دست آمده قابلیت بازتولید و خطای بازسازی کمی داشته باشند. برای ۱۰۰ ژنوم سرطانی شبیه‌سازی شده، می‌توانیم جواب‌های قابل بازتولید را برای N با مقادیر بین ۲ و ۱۰، شناسایی کنیم (میزان شباهت بین تکرارهای یک فرایند به عبارت دیگر احتمال اینکه فرایندی که در هر تکرار تولید می‌شود، در یک خوشه قرار گیرد). افزایش تعداد امضاها از ۲ به ۱۰ به طور قابل ملاحظه‌ای خطای بازسازی را کاهش می‌دهد، افزایش بیش از ۱۰ امضا بیش از این خطا را کاهش نمی‌دهد. این نشان می‌دهد که رویکرد ما می‌تواند "به صورت بهینه" امضاهای ده فرایند جهش را تشخیص دهد، دقیقاً عددی که در ابتدا برای شبیه‌سازی فهرست‌های جهش یافته ۱۰۰ ژنوم سرطان استفاده شده است.



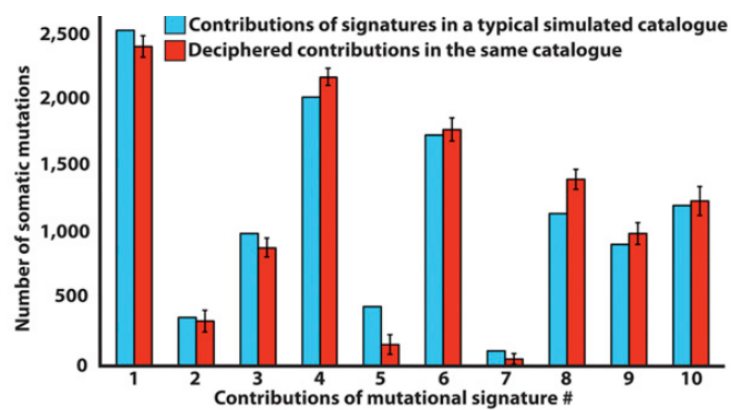
۱۰ امضای کشف شده قابلیت بازتولیدی بالایی دارند (شباهت امضاهایی تولید شده در هر تکرار در یک خوشه قرار می‌گیرند $0.96 < \text{silhouette width}$) و همچنین شباهت بالایی به امضاهایی دارند که برای تولید ۱۰۰ فهرست جهش‌زا مورد استفاده قرار گرفتند (متوسط شباهت کسینوسی < 0.98).



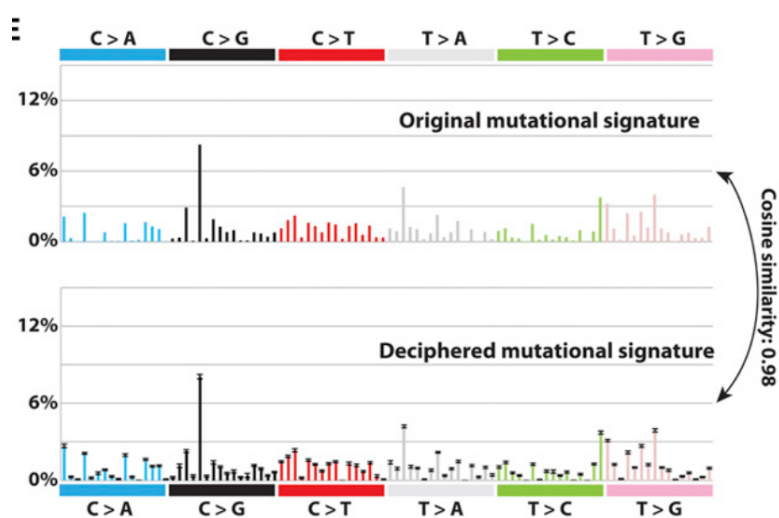
علاوه بر این، روش NMF قادر به شناسایی دقیق تعداد جهش‌های انجام شده توسط هر یک از ده فرآیند در هر یک از ژنوم‌ها است. مقایسه بین سهم تعداد جهش‌ها در یکی از امضاهای اصلی و امضا کشف شده‌اش در کلیه ژنوم‌ها در شکل زیر نشان داده شده است.



مقایسه سهم تعداد جهش هر ۱۰ امضا در یک ژنوم واحد در شکل زیر نشان داده شده است.



مقایسه بین یک امضاء اصلی و کشف شده در شکل زیر نشان داده شده است.



مقایسه بین یک فهرست جهش‌ها اصلی و بازسازی شده یک ژنوم در شکل زیر نشان داده شده است.

