



ژنومیک محاسباتی

مطهری و فروغمند

پاییز ۱۴۰۰

بازسازی درخت تبارزایی

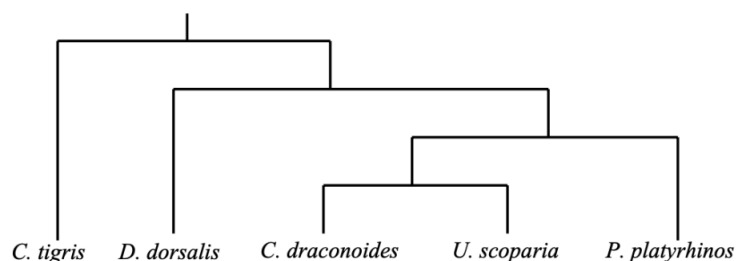
جلسه چهارم

نگارنده: سید پوریا لقایی

۱ مروری بر مباحث گذشته

در جلسه‌های گذشته مباحث هم‌ترازی رشته‌ها مطرح شد. در این جلسه مباحث بازسازی درخت‌های تبارزایی مطرح می‌شود.

۲ نمونه‌ای از درخت‌های تبارزایی

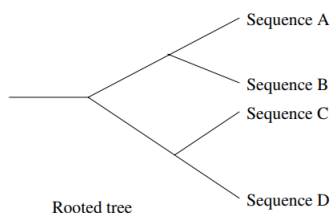


شکل ۱: درخت تبارزایی

در مدل‌های تکاملی سعی می‌کنیم با داشتن اطلاعاتی از گونه‌ها درخت تبارزایی آن را بسازیم. به عنوان مثال ممکن است برای یک ویروس درخت تبارزایی بسازیم تا در نهایت منشأ آن را پیدا کنیم.

۱.۲ ساخت درخت تبارزایی با ریشه

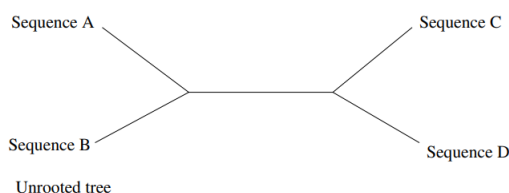
درخت تبارزایی در اصل یک درخت با ریشه است اما در بسیاری از موارد توانایی ساختن درخت ریشه دار را نداریم. معمولاً درخت های بدون ریشه ساخته می شوند و با استفاده از تکنیک هایی، از روی آن درخت ریشه دار را استخراج می کنیم.



شکل ۲: درخت تبارزایی ریشه دار

۲.۲ درخت تبارزایی بدون ریشه

اطلاعاتی که درخت بی ریشه به ما می دهد در مورد شباهت بین گونه ها است. اگر بتوانیم جای ریشه را در درخت بی ریشه پیدا کنیم به همان درخت ریشه دار (که در اصل در مسائل هم به دنبال آن هستیم) تبدیل می شود. نکته دیگر آنکه همه اطلاعات موجود در درخت ریشه دار در درخت بی ریشه نیز موجود است ولی فقط جای ریشه را نداریم. با توجه به اینکه ساختن درخت ریشه دار مسئله سخت است، فرض کنید یک الگوریتمی در زمان چند جمله ای برای ایجاد درخت بی ریشه داشته باشیم. در این صورت یک مسئله داریم با این عنوان که چطور از روی درخت بی ریشه، درخت ریشه دار را ایجاد کنیم.



شکل ۳: درخت تبارزایی بدون ریشه

۳ روش های تبدیل درخت بدون ریشه به ریشه دار

از راه حل های موجود برای حل این مسئله می توان از موارد زیر نام برد:

- اضافه کردن یک گونه متفاوت و دور نسبت به گونه های موجود در درخت بی ریشه (حالتی از تبدیل درخت بی ریشه به ریشه دار که در آن اطلاعاتی به مجموعه اضافه کرده ایم)

- وسط بزرگترین یال را به عنوان ریشه در نظر بگیریم

برای ساختن درخت تبارزایی باید ارتباط گونه ها با هم را داشته باشیم. به دو صورت ارتباط بین گونه ها مدل می شود:

- ماتریس ویژگی گونه ها

- ماتریس فاصله بین گونه ها

در بعضی از شرایط ها که تعداد ویژگی ها زیاد است ساخت درخت تبارزایی از روی ماتریس ویژگی کار مناسبی نیست. در حال حاضر درخت ساختن برای ویروسی مثل کرونا که داده های آن بسیار زیاد شده، کار بسیار سخت است. معمولاً با تغییر درخت های قدیمی ویروس کرونا، درخت جدید آن را می سازند.

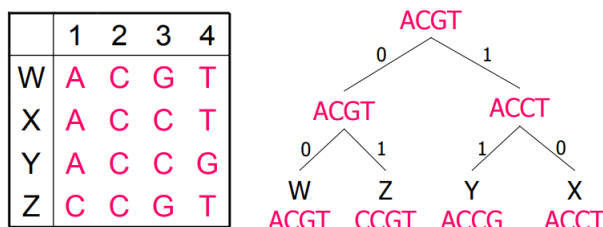
۱.۳ معیارهای سنجش درخت

حال که ساختن درخت ها را بررسی کردیم باید معیاری داشته باشیم که بتوانیم بهترین درخت را تشخیص دهیم. معیار های موجود شامل دو دسته می شوند :

- معیار بیشینه صرفه جویی
- معیار های احتمالاتی (بررسی تطابق درخت با داده های زیست شناسی)

۴ بازسازی ویژگی مبانی درخت تبارزایی

در این روش ورودی، ماتریس ویژگی گونه ها بوده. به عنوان نمونه در ژنوم انسان، ابتدا یک هم ترازی بین رشته ها ایجاد می شود و سپس از روی آن ماتریس ویژگی را استخراج می کنند. همانطور که بالاتر هم گفته شد، یک معیار برای تشخیص بهترین درخت، بیشینه صرفه جویی است. اگر درختی داشته باشیم که توپولوژی، برگ ها و نودهای میانی آن مشخص باشد و همینطور جمع میزان تغییرات بین همه نودها و فرزندان آنها، به عنوان هزینه درخت محسوب شود، مسئله ای تعریف می کنیم تحت این عنوان که بین درخت های ممکن، درختی را بیابیم که عدد جمع میزان تغییرات بین نودها و فرزندان آنها، کمینه باشد.



شکل ۴: مثال بیشینه صرفه جویی

۱.۴ حل مسئله بیشینه صرفه جویی (sankoff)

همانطور که بالاتر هم گفته شد اگر یک ماتریس ویژگی و یک توپولوژی برای درخت به عنوان ورودی داده شود (حتی جای گره های برگ را نیز داشته باشیم)، و از ما خواسته شود با این شرایط صرفه جویانه ترین درخت را پیدا کنیم راه حلی چند جمله ای برای این مسئله وجود دارد. الگوریتمی به نام sankoff که برای هر ویژگی یک بار جداگانه اجرا می شود. یک امتیاز تحت عنوان $A[v, c]$ تعریف می کنیم که به معنی بهترین امتیاز زیر درخت v به ازای انتصاب حرف c به آن است. بر همین اساس یک رابطه بازگشتی می نویسیم :

$$A[v, c] = \sum_{u: v \rightarrow u} \text{minscore}(c, c') + A[u, c']$$

فرمول بالا میزان امتیاز برای هر نود میانی است. برای برگ ها هم رابطه زیر را داریم :

$$A[v, c] = \begin{cases} 0 & \text{if } c = a \\ \infty & \text{else} \end{cases}$$

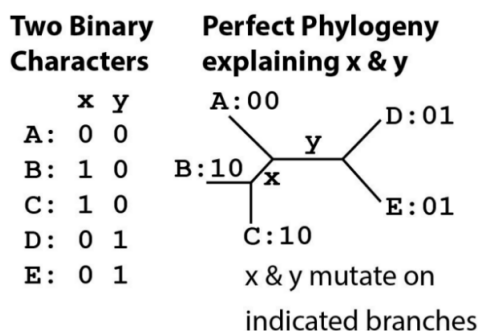
جواب نهایی برابر با حرفی است که کمترین مقدار امتیاز در ریشه را دارد. درخت موجود برای این مسئله تعداد $2m - 1$ راس خواهد داشت (m برگ و $m - 1$ نود میانی). برای هر ویژگی به تعداد گونه ها (m) ماتریس A پر می شود. تعداد کل ویژگی ها هم n تاست. پس پیچیدگی زمانی کلی الگوریتم (mn) است. پیچیدگی مکانی هم به دلیل اینکه به ماتریس ویژگی های پر شده قبلی نیاز نداریم و می توانیم آنها را دور بریزیم، (m) است.

۵ یافتن بهترین توپولوژی درخت

مسئله بعدی این است که بخواهیم بهترین توپولوژی درخت را نیز به دست آوریم. در این شرایط باید همه درخت ها را در نظر بگیریم. مشخص است که این مسئله راه حل چند جمله ای ندارد. اگر بخواهیم یک نود را به درخت اضافه کنیم باید روی یک یال اضافه شود. تعداد یال های درخت بدون ریشه برای n گونه و راس میانی درجه سه برابر است با : $(n - 2)!!$. در کل مسائل یافتن درخت تبارزایی صرفه جویانه، قابل حل در زمان چند جمله ای نیست. به همین دلیل روی این مسئله الگوریتم های هیوریستیک و تقریبی ارائه می شود

۶ درخت بی نقص

درخت بی نقص درختی است که مقدار هزینه آن دقیقاً برابر با تعداد ستون‌ها (ویژگی‌ها) باشد. نکته این درخت این است که به ازای هر دو ویژگی که در نظر بگیریم همه حالات را در سطرهای آن پیدا نکنیم.



شکل ۵: درخت بی نقص

۷ روش برآورد درست نمایی بیشینه

فرض کنید یک سکه داریم که احتمال شیر و خط آمدن آن را نمی‌دانیم. طبق مشاهدات می‌توانیم احتمال دیدن شیر و خط را تقریب بزنیم. هر چقدر میزان مشاهدات ما بیشتر بود تقریبی که داریم بهتر است. در روش برآورد درست نمایی بیشینه ما مدل احتمالاتی را انتخاب می‌کنیم که احتمال مشاهده مد نظر ما را بیشینه کند. مدل احتمالاتی از درخت و احتمال روی یال‌ها تشکیل شده است. اگر درخت و احتمالات را داشته باشیم و بدانیم روی هر راسی چه رشته‌هایی وجود دارد، می‌توانیم احتمال مشاهده را حساب کنیم. مسئله دیگر این است که فرض کنیم تنها برگ‌ها موجود باشد و توپولوژی درخت را نیز بدانیم. حال آیا می‌توانیم مسئله را با روش برآورد درست نمایی بیشینه حل کنیم؟ در این صورت باید همه حالات را در نظر بگیریم. این محاسبات بسیار طولانی بوده و ممکن است در زمان چند جمله‌ای قابل اجرا نباشد. اما با یک رابطه بازگشتی می‌توان این مسئله را در زمان بهتری حل کرد:

$$\prod_i \left(\frac{1}{2} L_i(r, 0) + \frac{1}{2} L_i(r, 1) \right)$$

حالت پایه برای برگ V نیز برابر است با:

$$L_i(v, s) = \mathbb{1}_{s=M[v,i]}$$

و همینطور یک رابطه بازگشتی برای راس میانی u با فرزندان v و w:

$$L_i(u, s) = \left[\sum_{y \in \{0,1\}} L_i(v, y) Pr(v_i = y | u_i = s) \right] \left[\sum_{x \in \{0,1\}} L_i(w, x) Pr(w_i = x | u_i = s) \right]$$

مراجع

[Sun09] Wing-Kin Sung. *Algorithms in bioinformatics: A practical introduction*. CRC Press, 2009.

[Sun09].