



ژنومیک محاسباتی

مطهری و فروغمند
پاییز ۱۴۰۰

درخت تبارزایی (۴)

جلسه هفتم: کار با چند درخت

نگارنده: سهیل رستگار

۱ مروری بر مباحث گذشته

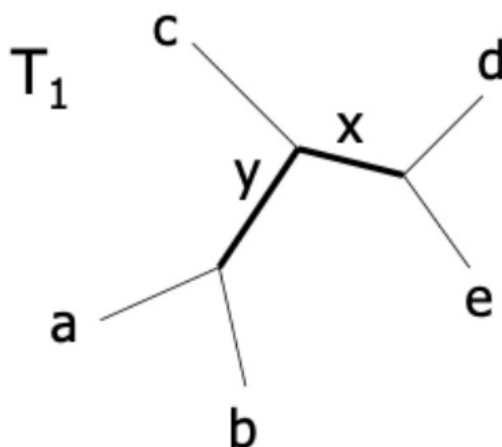
در جلسات گذشته‌ی حول این موضوع به درخت‌های تبارزایی^۱ پرداختیم و روش‌های تولید آن‌ها را از روی داده‌های خام، یا ماتریس‌های فاصله بررسی کردیم. دیدیم که اگر ساختار درخت را داشته باشیم می‌توان با پیچیدگی محاسباتی کم (خطی) و با استفاده از برنامه‌نویسی پویا مقادیر میانی درخت تبارزایی را از روی گونه‌های موجود در حال حاضر تولید کرد. همینطور دیدیم که عموماً ساختار درخت را نداریم و تولید یک درخت تبارزایی در این حالت بسیار پیچیده و در واقع یک الگوریتم np است.

در جلسه‌های بعد چند مدل احتمالاتی را برای ساخت درخت تبارزایی بررسی کردیم و در نهایت یک روش خوب که ساختن درخت بر اساس ماتریس فواصل بین گونه‌های مختلف بود را معرفی کرده و چند الگوریتم برای حالات مختلف ماتریس فاصله ارائه دادیم. در این جلسه قصد داریم که روش‌هایی را برای مقایسه‌ی درخت‌های مختلف ارائه دهیم تا بتوانیم مکانیزمی برای مقایسه‌ی و ارزیابی درخت‌های تولید شده و درخت‌های موجود داشته باشیم. همینطور در ادامه روش‌هایی را برای ساختن درخت‌های توافقی و اجماعی از روی چند درخت ارائه می‌دهیم.

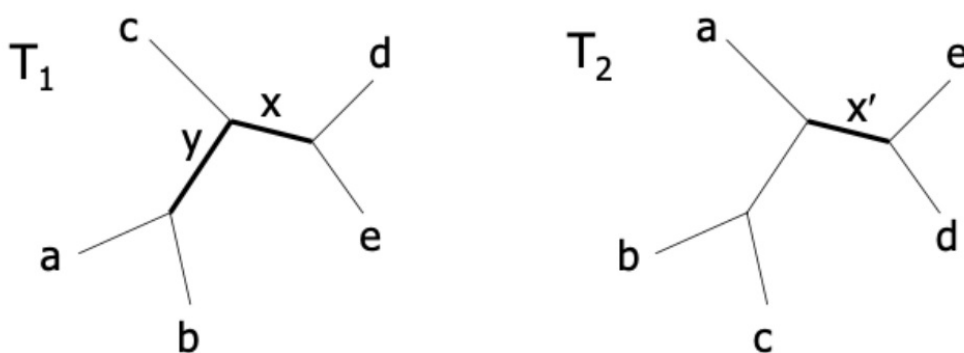
۲ مقایسه‌ی درخت‌ها

در این بخش قصد داریم روشی را برای محاسبه‌ی فاصله‌ی دو درخت ارائه دهیم تا بتوان میان درخت‌های مختلف مقایسه انجام داد و شباهت و تفاوت آن‌ها را ارزیابی کرد.

^۱ phylogenetic tree



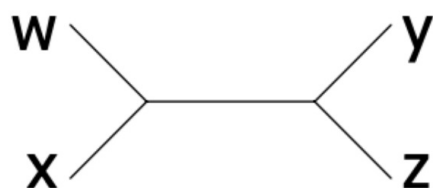
شکل ۱: یک نمونه درخت



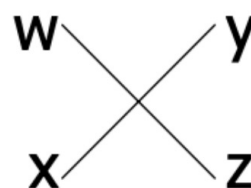
شکل ۲: دو درخت دلخواه

۱.۲ تقسیم (Split)

می‌دانیم که درخت گراف بدون دور است. پس اگر یک یال آن را قطع کنیم به دو بخش مجزا می‌رسیم که یالی بین هم‌دیگر ندارند. از همین قضیه برای تعریف تقسیم روی درخت استفاده می‌کنیم. دو مجموعه‌ی مجزا از برگ‌های درخت را یک تقسیم روی یال x می‌نامیم، اگر با حذف یال x از درخت اصلی، برگ‌های درخت به دو مجموعه‌ی مجزای داده شده افزاش شوند. برای مثال در شکل ۱ در صورتی که تقسیم روی یال x انجام گیرد، به دو مجموعه مجزای $\{abc | de\}$ می‌رسیم و اگر تقسیم روی یال y انجام شود به دو بخش $\{ab | cde\}$.



Butterfly quartet



Star quartet

شکل ۳: حالت‌های بخش‌بندی

۲.۲ فاصله‌ی Robinson-Foulds

برای محاسبه‌ی این معیار فاصله بین دو درخت، ابتدا تمام تقسیم‌های درخت اول و درخت دوم را به صورت جدا حساب می‌کنیم، در واقع به دو مجموعه از تقسیم‌های مختلف می‌رسیم. سپس عناصر غیرمشترک بین این دو مجموعه، یا به عبارت دیگر اجتماع دو مجموعه منهای اشتراک دو مجموعه را حساب کرده و تقسیم بر ۲ می‌کنیم. عدد بدست آمده، فاصله‌ی Robinson-Foulds بین دو درخت است. برای مثال جهت محاسبه‌ی فاصله‌ی بین دو درخت در شکل ۲ داریم:

$$Splits(T_1) = \{a|bcde, b|acde, c|abde, d|abce, e|abcd, ab|cde, abc|de\}$$

$$Splits(T_2) = \{a|bcde, b|acde, c|abde, d|abce, e|abcd, bc|ade, abc|de\}$$

$$D = Splits(T_1) \cup Splits(T_2) - Splits(T_1) \cap Splits(T_2) = \{ab|cde, bc|ade\}$$

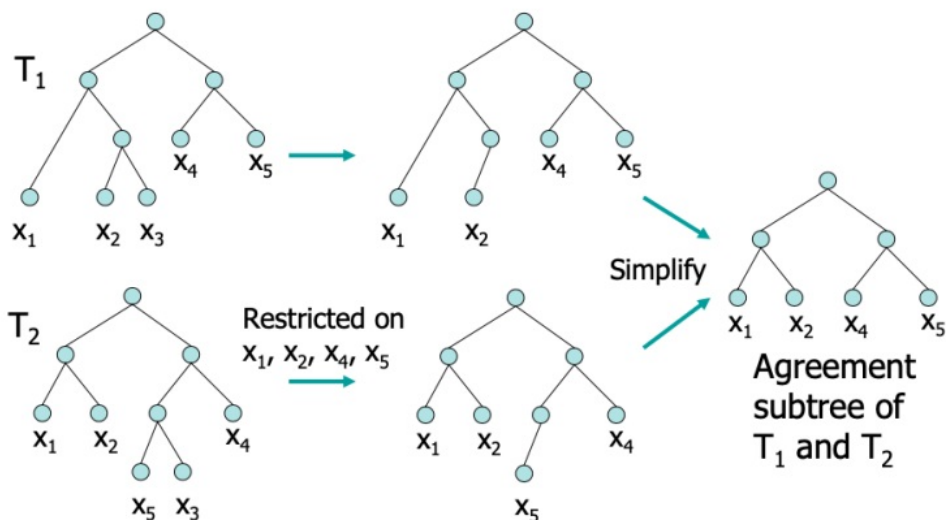
$$\text{Robinson-Foulds}(T_1, T_2) = \frac{|D|}{2} = 1$$

از تعریف بعضی از خواص فاصله‌ی Robinson-Foulds مشخص می‌شود. برای مثال فاصله‌ی هر درخت با خودش ۰ است و فاصله‌ی درخت a با درخت b برابر فاصله‌ی درخت b تا درخت a است.

۳.۲ فاصله‌ی چهارتایی (Quartet Distance)

فرض می‌کنیم ۴ برگ از یک درخت را به نام‌های w, x, y, z داشته باشیم. این ۴ برگ بر اساس بخش‌بندی ۴ حالت دارند، اینکه x با w جد مشترک نزدیک‌تری داشته باشند و y با z، اینکه x با y و w با z، یا اینکه x با w و y با z، سه حالت را می‌سازند که در این سه حالت یعنی یالی وجود دارد که در صورت قطع شدن آن، برگ‌های که جد مشترک نزدیکتر دارند در یک درخت و دو برگ دیگر که باز جد مشترک نزدیکتر دارند در درخت دیگر قرار می‌گیرند. به این ۳ حالت، Butterfly quartet گفته می‌شود. یک حالت نیز وجود دارد که جد مشترک هر ۴ برگ یکسان باشد که به آن Star quartet گفته می‌شود (شکل ۳). پس در کل ۴ حالت بخش‌بندی وجود دارد.

برای محاسبه‌ی فاصله‌ی چهارتایی، تمام ۴ تایی‌های ممکن را در نظر می‌گیریم و بخش‌بندی آن‌ها را حساب می‌کنیم. اگر بخش‌بندی این ۴ تایی در دو درخت یکسان نبود، آن حالت را به عنوان یک واحد فاصله‌ساز در نظر می‌گیریم. در نهایت تعداد ۴ تایی‌هایی که در دو درخت بخش‌بندی متفاوتی دارند را حساب کرده و به عنوان فاصله‌ی دو درخت در نظر می‌گیریم.



شکل ۴: مراحل رسیدن به زیردرخت توافقی

۳ درخت توافقی

بزرگترین زیردرخت توافقی بین دو درخت، درختی است که با حذف کمترین تعداد گونه از هر دو درخت و ساده‌سازی (یعنی حذف راس‌های درجه دو به جز ریشه، یعنی راس‌هایی که فرزند خود را به والد متصل می‌کنند، و اتصال مستقیم فرزند به والد) بتوان به آن رسید. شکل ۴ مثالی از زیردرخت توافقی را نشان می‌دهد. واضح است که با انتخاب دو گونه‌ی دلخواه یکسان از هر دو درخت و حذف بقیه‌ی برگ‌ها، به دو درخت یکسان که در آن این دو گونه به ریشه متصل‌اند می‌رسیم، پس در حالت حداقلی درختی به اندازه ۲ وجود دارد.

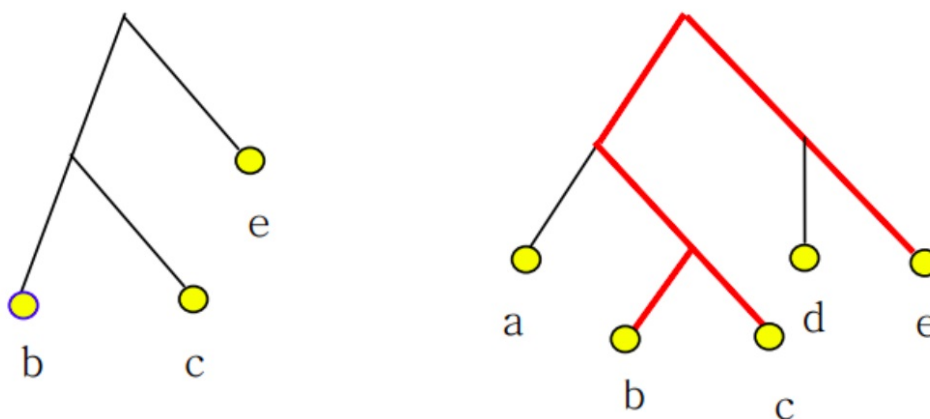
الگوریتم بدست آوردن بزرگترین زیردرخت توافقی از برنامه‌نویسی پویا کمک می‌گیرد. به این صورت که به ازای هر راس u از درخت اول و هر راس v از درخت دوم، امتیاز بزرگترین زیردرخت بین زیردرخت‌هایی که این رئوس ریشه‌ی آن هستند را حساب می‌کند، و برای اینکار حالات مختلف انتصاب زیردرخت هرکدام از فرزندان این راس‌ها به همدیگر را بررسی کرده و مقدار بیشینه را حساب می‌کند. در رابطه‌ی بازگشتی زیر، این موضوع نشان داده شده. با این نکته که $MAST(T_1^u, T_2^v)$ یعنی اندازه‌ی بزرگترین زیردرخت توافقی بین زیردرخت نوادگان u از T_1 و زیردرخت نوادگان v از T_2 . حالت پایه نیز زمانی است که u و v جزو برگ‌ها باشند که اگر یکسان باشند جواب ۱ و در غیراین صورت جواب ۰ است.

$$MAST(T_1^u, T_2^v) = \max \begin{cases} MAST(T_1^{u_1}, T_2^{v_1}) + MAST(T_1^{u_2}, T_2^{v_2}) \\ MAST(T_1^{u_1}, T_2^{v_1}) + MAST(T_1^{u_2}, T_2^{v_2}) \\ MAST(T_1^{u_1}, T_2^v) \\ MAST(T_1^{u_2}, T_2^v) \\ MAST(T_1^u, T_2^{v_1}) \\ MAST(T_1^u, T_2^{v_2}) \end{cases}$$

५

الگوریتم ۱ ساخت درخت اجماعی اکثریت

- ۱: به ازای $i = 1$ تا m :
- ۲: به ازای هر یال c در درخت T_i :
- ۳: یال q را یال والد c در T_i در نظر می‌گیریم
- ۴: اگر $P[c] = \emptyset$ یا $|B_P[c]| < |B_q|$ آنگاه:
- ۵: $P[c] = q$
- ۶: پایان اگر
- ۷: پایان به ازای
- ۸: پایان به ازای



شکل ۵: ساخت ۳ تایی از روی یک درخت ریشه‌دار

۳.۴ درخت R^*

درخت R^* نیز یک مدل درخت اجماعی است که بر اساس سه‌تایی‌های درخت‌های موجود ساخته می‌شود.

۱.۳.۴ سه‌تایی

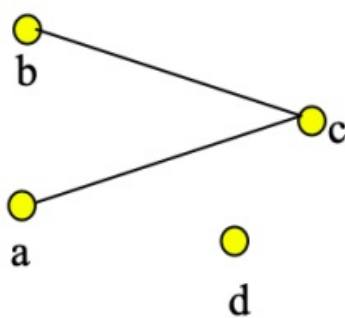
در یک درخت ریشه‌دار، با نگهداری ۳ برگ و حذف بقیه برگ‌ها و منقبض سازی درخت، به یک درخت ۳ تایی می‌رسیم. در این ۳ تایی والد دو برگ، فرزند والد برگ دیگر خواهد بود، یا به عبارت دیگر دو برگ به هم نزدیک‌تر هستند تا برگ سوم. برای نمایش این موضوع از نماد $ab|c$ استفاده می‌کنیم که به این معنی است که برگ c در سه‌تایی، از برگ‌های a و b دورتر است. با توجه به این موضوع، برای هر ۳ تایی ۳ حالت مختلف داریم، براساس اینکه کدام برگ دورتر باشد. برای مثال سه‌تایی شکل ۵ به صورت $bc|e$ قابل نمایش است.

۲.۳.۴ گراف باقی‌مانده

الگوریتم ساخت گراف باقی‌مانده، سه‌تایی‌های درخت را دریافت کرده و به ازای هر $ab|c$ یک یال بین a و b می‌گذارد. در مرحله‌ی نهایی، یک گراف غیرهمبند خواهیم داشت که هر بخش یکی از زیرشاخه‌های ریشه‌ی اصلی درخت هستند. به بیان دیگر این الگوریتم با حذف ریشه، مجموعه‌های باقی‌مانده را برمی‌گرداند. مثالی از روش ساخت گراف باقی‌مانده را در شکل ۶، و چند نمونه از مجموعه‌های باقی‌مانده چند درخت را در شکل ۷ می‌توان مشاهده کرد.

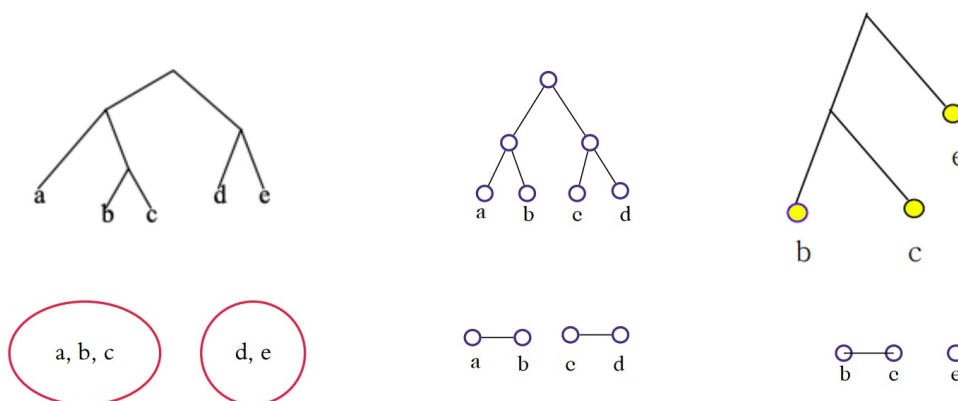
۳.۳.۴ ساخت درخت از روی سه‌تایی‌ها

در این بخش قصد داریم الگوریتمی ارائه دهیم که با داشتن سه‌تایی‌ها، درختی که آن سه‌تایی‌ها را داشته باشد بسازد. الگوریتم ۲ بدین منظور ارائه شده است. یک مثال از اجرای این الگوریتم برای ورودی $bc|a$ و $ac|d$ و $b|de$ را می‌توان در شکل ۸ مشاهده کرد.



$$bc \mid a, ac \mid d$$

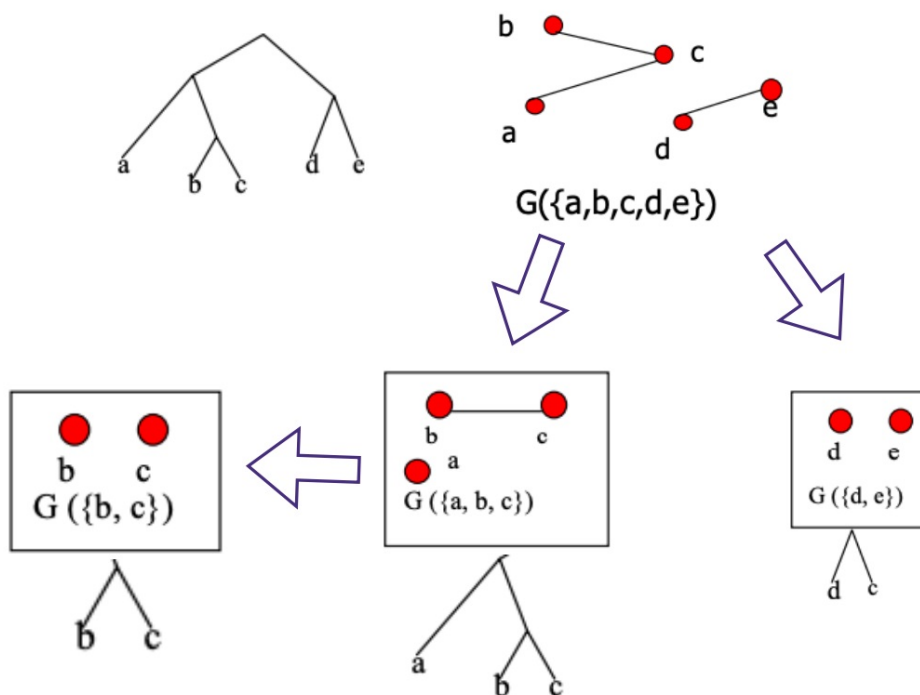
شکل ۶: یک نمونه گراف باقی مانده



شکل ۷: چند مجموعه‌ی باقی مانده

الگوریتم ۲ ساخت درخت از روی سه‌تایی‌ها

- ۱: گراف باقی مانده را بساز
- ۲: اگر ناهمبند بود:
- ۳: برای هر زیرمجموعه درخت را بساز (بازگشتی)
- ۴: یک راس ریشه‌ی جدید اضافه کن و همه‌ی درخت‌ها را زیردرخت آن بگذار
- ۵: این درخت را برگردان
- ۶: اگر نه:
- ۷: باید حتماً تک راس باقی مانده باشد، وگرنه نمی‌توان درخت ساخت
- ۸: آن تک راس را به عنوان درخت بازسازی شده برگردان



شکل ۸: مثال از اجرای الگوریتم ۲

۴.۳.۴ ساخت درخت R^*

برای ساخت درخت T از روی درخت‌های $\{T_1, T_2, \dots, T_m\}$ به روش R^* ، از هر درخت مجموعه‌ی ۳ تایی‌ها را در یک ماتریس علامت می‌زنیم. در آن ماتریس به ازای هر ۳ برگ مختلف، شیوه‌های ساخت ۳ تایی که ۳ حالت می‌تواند داشته باشد، در نظر گرفته می‌شود و اینکه از هر شیوه چه تعداد درخت، ۳ تایی را به آن صورت دارند هم ثبت می‌شود. در نهایت شیوه‌ای که بیشترین تعداد را داشته باشد، به عنوان ۳ تایی نهایی آن ۳ برگ برای درخت T انتخاب می‌شود. به این صورت از $TL(T_1), \dots, TL(T_m)$ ها به $TL(T)$ می‌رسیم (TL مجموعه‌ی سه‌تایی‌هاست). حال با داشتن مجموعه‌ی سه‌تایی‌های درخت T ، با استفاده از الگوریتم ۲ درخت را می‌سازیم.

۵ روش‌های نوشتن درخت

در این بخش روش‌هایی را برای تبدیل ساختار یک درخت به متن جهت انتقال راحت ارائه داده می‌دهیم.

۱.۵ قالب Newick

این قالب حالت‌های مختلف را پوشش می‌دهد اما قواعد کلی آن به این صورت است که هر پرانتز بیانگر یک راس میانی است و در آن فرزندان با ، از همدیگر جدا شده‌اند. اگر راسی اسم داشته باشد در سمت راست پرانتز مربوط به آن راس نوشته می‌شود، یا اگر برگ باشد به صورت مستقیم نوشته می‌شود. اگر یک راس فاصله تا ریشه داشته باشد، بعد از رعایت قواعد قبلی نام گذاری، در آخر در سمت راست اطلاعات هر راس : گذاشته شده و سمت راست آن میزان فاصله نوشته می‌شود. چند مثال این قالب در شکل ۹ نمایش داده شده است.

۲.۵ قالب NEXUS

این قالب علاوه بر ساختار درخت، اطلاعات دیگری را نگه می‌دارد، برای مثال ویژگی‌های گونه‌های موجود و قواعد تبدیل نام آن‌ها به اعداد برای استفاده در ساختار درخت. یا مثلاً ماتریس هم‌ردیفی بین گونه‌های موجود. در نهایت ساختار درخت طبق ساختار Newick داده می‌شود. مثالی از قالب NEXUS در شکل ۱۰ قابل مشاهده است.

<code>(,,(,));</code>	<i>no nodes are named</i>
<code>(A,B,(C,D));</code>	<i>leaf nodes are named</i>
<code>(A,B,(C,D)E)F;</code>	<i>all nodes are named</i>
<code>(:0.1,:0.2,(0.3,0.4):0.5);</code>	<i>all but root node have a distance to parent</i>
<code>(:0.1,:0.2,(0.3,0.4):0.5):0.0;</code>	<i>all have a distance to parent</i>
<code>(A:0.1,B:0.2,(C:0.3,D:0.4):0.5);</code>	<i>distances and leaf names (popular)</i>
<code>(A:0.1,B:0.2,(C:0.3,D:0.4)E:0.5)F;</code>	<i>distances and all names</i>
<code>((B:0.2,(C:0.3,D:0.4)E:0.5)F:0.1)A;</code>	<i>a tree rooted on a leaf node (rare)</i>

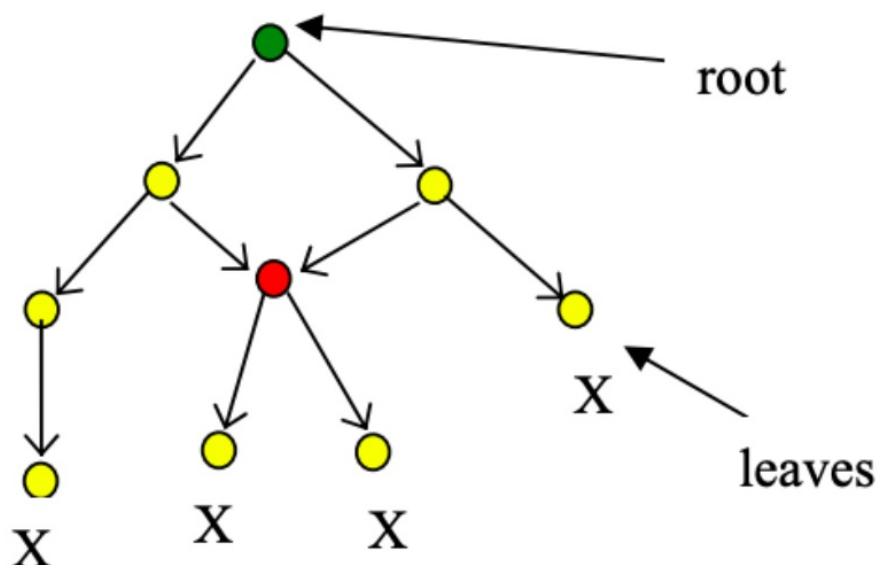
شکل ۹: قالب Newick

```
#NEXUS
Begin TAXA;
  Dimensions ntax=4;
  TaxLabels SpaceDog SpaceCat SpaceOrc SpaceElf
End;

Begin data;
  Dimensions nchar=15;
  Format datatype=dna missing=? gap=- matchchar=.;
  Matrix
    [ When a position is a "matchchar", it means that it is the same as the first entry
    at the same position. ]
    SpaceDog   atgctagctagctcg
    SpaceCat   .....??...-a.
    SpaceOrc   ...t.....-g. [ same as atgttagctag-tgg ]
    SpaceElf   ...t.....-a.
  ;
End;

BEGIN TREES;
  Tree tree1 = (((SpaceDog,SpaceCat),SpaceOrc,SpaceElf));
END;
```

شکل ۱۰: قالب NEXUS



شکل ۱۱: یک شبکه‌ی تبارزایی

۶ ابردخت (Supertree)

ابر درخت یک درخت تبارزایی است که به جای اینکه مستقیم از روی گونه‌های مختلف ساخته شود، از روی درخت‌های تبارزایی روی زیرمجموعه‌های گونه‌ها (مثلاً گونه‌های شبیه بهم از یک رده) ساخته می‌شود. این کار باعث می‌شود که بتوان روش‌های درخت اجماعی را روی زیرمجموعه‌های شبیه و کوچک از گونه‌ها اجرا کرد و از پیچیدگی زیاد مسئله جلوگیری نمود.

۷ شبکه‌های تبارزایی (Phylogenetic Networks)

شبکه‌های تبارزایی مانند درخت‌های تبارزایی هستند، با این تفاوت که هر موجود می‌تواند بیش از یک والد داشته باشد. کاربرد اصلی آن در باکتری‌ها و فرایندهای انتقال ژن افقی آنان است که در آن باکتری‌های یک نسل می‌توانند با یکدیگر تبادل ژنتیکی داشته باشند و فرزندی تولید کنند که از ژنتیک چندین گونه به ارث می‌برد. شکل ۱۱ مثالی از این شبکه‌های تبارزایی است.

مراجع

[For21] Hadi Foroughmand. Computational genomics course slides. Sharif University, 2021.

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