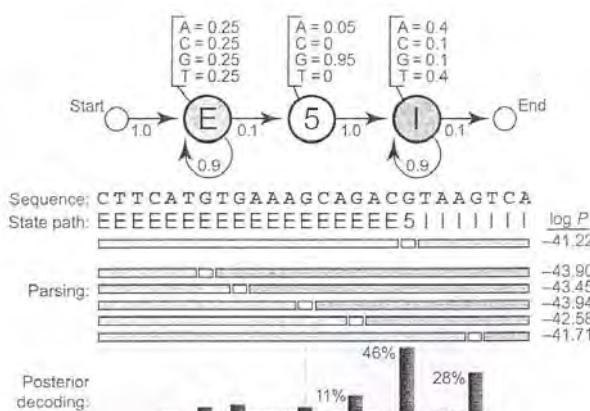


باشد. اجازه دهید اختلاف‌های ساده را در نظر بگیریم مثلاً در اگزون توزیع بازها به طور متوسط یکنواخت است (یعنی ۲۵٪ از هر کدام) در اینترون‌ها ($A/T = 0.5$) زیاد است (مثلاً ۴۰٪ درصد از هر یک از A, T و ۱۰٪ درصد از دو حرف دیگر C و G) و اجماع نوکلئوتیدی در محل SS' تقریباً همیشه G است (مثلاً ۹۵٪ اوقات G و ۵٪ اوقات A).

با شروع از این اطلاعات، می‌توانیم یک HMM رسم کنیم (شکل ۱). این HMM سه وضعیت را به وجود می‌آورد که هر یک از این‌ها را به عنوان برچسبی به هر نوکلئوتید قرار می‌دهیم (اگزون) $5'$ و SS' (۵) و I (اینترون). هر وضعیت، دارای احتمال‌های انتشار مختص خودش است (که در بالای وضعیتها نشان داده شده‌اند)، این احتمال‌ها ترکیب بازه‌ای اگزون، اینترون‌ها و اجماع G در محل SS' را الگو می‌کنند. همچنین هر وضعیت دارای احتمال‌های انتقال است (کمان‌ها) که احتمال‌های حرکت از هر وضعیت به وضعیت دیگر را نشان می‌دهد. احتمال‌های انتقال، رتبه خطی را که طی آن انتظار داریم وضعیتها رخ دهنند، توجیه می‌کنند یک یا چند E ، یک یا چند 5 ، یک یا چند I .



شکل ۱ یک HMM نمایشی برای تشخیص محل برش ۵.

پس چه چیز پنهان است؟

مفید است تصور کنیم که یک HMM یک توالی تولید می‌کند. وقتی وضعیتی را ملاقات می‌کنیم از توزیع احتمال انتشارها یک عضو منتشر می‌نماییم. سپس بر مبنای توزیع احتمال وضعیت‌ها، وضعیت بعدی را انتخاب می‌کنیم بنا بر این مدل، دو رشته از اطلاعات تولید می‌شود. یکی مسیر وضعیت‌ها (برچسب‌ها) است که ما از یک وضعیت به وضعیت دیگر منتقل می‌شویم. دیگری توالی مشاهدات (AND) است که از هر وضعیت داخل مسیر یک عضو منتشر می‌شود.

مسیر وضعیت، یک زنجیر مارکف است یعنی هر وضعیتی که قرار است برویم فقط به وضعیت فعلی مان بستگی دارد. از آنجایی که به ما فقط توالی مشاهدات را می‌دهند، این مسیر وضعیت، پنهان است - این‌ها برچسب‌های عضوها هستند که علاقمند به استنباط راجع به آن‌ها هستیم. مسیر وضعیت یک زنجیر مارکف

الگوی مارکف پنهان چیست؟

شون آر. ادی

بخش ژنتیک دانشگاه واشنگتن

الگوهای آماری بنام الگوهای مارکف پنهان در زیست‌شناسی محاسباتی به دفعات مورد ارجاع قرار می‌گیرند. این مدل‌ها چه هستند و چرا این قدر برای مسائل متعدد مفیدند؟ تجزیه و تحلیل توالی‌های زیستی اغلب به عنوان قرار دادن برچسب درست روی هر عضو توالی تعریف می‌شود (که آن عضو را معمولاً با نام residue می‌شناسند). در شناسایی زن می‌خواهیم برچسب روی نوکلئوتیدها را تحت عنوان اگزون‌ها، اینترون‌ها و یا توالی مابین ژنی قرار دهیم. دربحث هم ردیفی توالی‌ها به دنبال نسبت دادن اعضای یک دنباله ناشناش به اعضای همسان آن در یک مجموعه هدف از دنباله‌ها هستیم. ما همیشه می‌توانیم یک برنامه خاص برای یک مسئله داده شده بنویسیم اما غالباً با مشکلات مشابهی مواجه می‌شویم. یکی از مشکلات عبارت است از به کارگیری منابع ناهمگون اطلاعات. برای مثال یک زن باب باید نواحی مشترک برای برش، اربیسی کودونی، ترجیه طولی اگزون اینترون و آنالیز چارچوب باز را در یک سیستم نمره‌دهی ترکیب کند. چطور این پارامترها باید تنظیم شوند؟ به اطلاعات گوناگون چگونه باید وزن داد؟ مشکل دیگر تعبیر و تفسیر نتایج به صورت احتمالی است. یافتن بهترین نمره (امتیاز) پاسخ یک هدف است ولی معنی نمره چیست و چقدر مطمئنیم که بهترین نمره درست را یافته‌ایم؟ سومین مسئله، تعیین پذیری است. هم زمان با کامل ساختن زن باب ویژه آرزو می‌کنیم ای کاش اجماع اولیه انتقالی، برش‌های دیگر و علامت پلیا دنیلی را هم مدل می‌کردیم.

اغلب اوقات با افزودن حقایق بیشتر به یک برنامه ویژه ظریف موجودات فروبریزی آن را زیر بار وزن سنجینش پذیرد می‌آوریم. الگوهای مارکف پنهان (HMM) مبانی رسمی الگوهای احتمالی برچسب‌گذاری بر توالی‌های خطی را ارائه می‌دهند. آن‌ها ابزار مفهومی برای ساخت مدل‌های پیچیده با رسم تصویری ذهنی پذیرد می‌آورند. آن‌ها در قلب دامنه وسیعی از برنامه‌ها مثل، زن باب‌ها، جست و جوگرهای پروفایلی، هم‌ردیف‌باب‌های چندگانه، توالی‌ها و تعیین کننده‌های محل‌های تنظیم ژنی قرار دارند. HMM‌ها نمادهای تجزیه و تحلیل محاسباتی توالی‌ها هستند.

یک HMM نمایشی تشخیص دهنده محل برش ۵

به عنوان یک مثال ساده مسئله تشخیص محل برش ۵ در نمودار زیر را در نظر بگیرید. فرض کنید یک توالی DNA داریم که از یک اگزون شروع شده و یک محل برش ۵ دارد و به یک اینترون ختم می‌شود. مسئله تعیین محل برش ۵ است. یعنی جایی که اگزون به اینترون تغییر پیدا می‌کند - یا به عبارتی مسئله تعیین SS' است. برای ما که می‌خواهیم حدسی هوشمندانه بزنیم توالی‌های اگزون، محل‌های برش و اینترون‌ها باید خواص آماری مختلفی داشته

صورت کسر است که به مجموع ۱۴ مسیر وضعیت در مخرج کسر تقسیم شده است. ما برای حالتی که بهترین امتیاز پنجمین G درست باشد، احتمالی برابر ۴۶ به دست می آوریم، حال آنکه برای حالتی که ششمین G درست باشد احتمال برابر با ۲۸ را خواهیم داشت (شکل ۱، پایین). این را رمزگشایی پسین (posterior decoding) می نامند. برای مسائل بزرگتر، رمزگشایی پسین از دو الگوریتم به نام های پیشرو و پسرو استفاده می کنند. این ها اساساً شبیه الگوریتم ویتری بیشتر فقط به جای پیدا کردن محتمل ترین مسیر به جمع بندی روی تمام مسیرهای ممکن توجه می کنند. ساختن الگوهای واقع گرایانه تر ساختن یک HMM یعنی مشخص نمودن چهار چیز الفبای علائم K علامت مختلف (مثل ACGT)، $k = 4$ (ii) تعداد وضعیت های مدل M : (iii) احتمال های انتشار $e_{i(x)}$ برای هر وضعیت i که جمع روی K علامت x برابر ۱ می شود یعنی $\sum e_{i(x)} = 1$ و (vi) احتمال های انتقال t_{ij} برای هر وضعیت i که به وضعیت مانند j (شامل خودش) می رود و جمع روی M وضعیت j برابر ۱ می شود یعنی $\sum t_{ij} = 1$. هر مدلی که این خواص را داشته باشد یک HMM است.

یعنی هرکس می تواند با رسم تصویر از مسئله ی مورد نظرش شبیه به مسئله شکل ۱ یک HMM جدید بسازد. این سادگی گرافیکی اجازه می دهد هرکس به تعریف بیولوژیک خاصی از یک مسئله متتمرکز شود. برای مثال در الگوی نمایشی ما یعنی الگوی محل های برش ممکن است از قدرت ممیزی مدل راضی نباشیم ممکن است بخواهیم یک اجماع واقع گرایانه تر شش نوکلئوتیدی GTRAGT را به محل برش ۵' اضافه کیم می توانیم سطربی با شش وضعیت HMM به جای وضعیت ۵ قرار دهیم تا یک اجماع شش باز را با پارامتری کردن احتمال های انتشار در محل های برش ۳' الگو کنیم. ممکن است بخواهیم یک مدل کامل اینترونی شامل یک محل برش ۳' را مدل کنیم و یک ۳' به عنوان وضعیت اگزونی به مدل اضافه کنیم تا اجازه دهد یک توالی مشاهده شده به جای یک اینترون به یک اگزون ختم شود. بنابراین ممکن است علاقمند باشیم یک مدل کامل ژنی بسازیم. هر چه که اضافه می کنیم فقط کافی است به صورت ترسیمی مشخص شود.

جالب توجه HMM ها به وابستگی بین اعضاء توالی مشاهدات به نحو مناسب نمی پردازند زیرا آنها فرض می کنند به شرط وضعیت داده شده مشاهدات از هم مستقلند. مثالی که در آن مدل های HMM اغلب نامناسبند، تجزیه تحلیل ساختار دوم RNA است. جفت بازهای حفظ شده RNA همبستگی های جفتی با طول بلند ایجاد می کنند کی از موقعیت ها می توانند هر یک از اعضای توالی باشد اما جفت باز شریک باید مکمل آن باشد. یک وضعیت مسیر HMM هیچ راهی به جز به یاد آوردن این مطلب که کدام وضعیت دور تولید شده است، ندارد. گاهی اوقات می توانیم قوانین HMM را بدون آن که الگوریتمها را خراب کنیم، دور بزنیم. برای مثال ممکن است در بحث ژن یابی بخواهیم به جای انتشار سه عضو توالی، یک سه تایی وابسته از کodon ها را انتشار دهیم الگوریتم های HMM بدون هیچ مشکلی قابل تعمیم به انتشار وضعیت های

پنهان است.

احتمال $P(S, \pi | HMM, \Theta)$ که HMM با پارامترهای Θ مسیر وضعیت و توالی مشاهدات S را تولید کند برابر است با حاصل ضرب همه احتمال های انتشارها در احتمال های انتقال های به وقوع پیوسته. برای مثال توالی ۲۶-نوکلئوتیدی به همراه مسیر وضعیت را در وسط شکل ۱ در نظر بگیرید که در آن در مجموع ۲۷ انتقال و ۲۶ انتشار وجود دارد. همه ۵۳ احتمال را در هم ضرب کنید (و لگاریتم گیری انجام دهید، زیرا که این ها اعداد کوچکی هستند) شما عدد $\log P(S, \pi | HMM, \Theta) = 41.22$ را به دست خواهید آورد. یک HMM یک الگو کامل احتمالی است - پارامترهای مدل و امتیاز کل توالی همگی احتمال ها هستند. بنابراین می توانیم از نظریه احتمال بیزی برای دست کاری این اعداد با استفاده از روش های توانمند استاندارد شامل بهینه سازی پارامترها و تفسیر معنی داری امتیازها، بهره جوئیم.

یافتن بهترین مسیر وضعیت

در یک مسئله تجزیه تحلیلی به ما توالی را داده اند و ما می خواهیم مسیر وضعیت پنهان را پیدا کیم. مسیرهای زیادی پتانسیل آن را دارند که یک توالی را تولید کنند. ما به دنبال توالی هستیم که بیشترین احتمال را داراست.

برای مثال، اگر به ما HMM و توالی ۲۶-نوکلئوتیدی شکل ۱ را بدهند، ۱۴ مسیر ممکن با احتمال غیر صفر وجود دارد، زیرا ۵'SS ۵' باید روی یکی از ۱۴ تا A یا G داخلی بیفتند. شکل ۱ شش تا از مسیرهای با امتیاز بالا را نشان می دهد (مسیرهایی با G در ۵'SS ۵'). بهترین مسیر دارای لگاریتم احتمال ۲۲.۴۱ است که نتیجه می دهد که محتملترین وضعیت $SS 5'$ در پنجمین G است.

در اکثر مسائل تعداد توالی های ممکن از وضعیت ها آنقدر زیاد است که قادر به نمایش همه آنها نیستیم. الگوریتم کارآی و پتری ضمانت می کند که با در دست داشتن توالی مشاهدات و MMH محتملترین مسیر وضعیت را پیدا کند. الگوریتم پتری یک الگوریتم برنامه ریزی پویا بسیار شبیه به آن الگوریتم های استانداردی است که در هم رده توالی ها به کار گرفته می شوند.

بعد از هم رده هایی بهترین امتیاز

شکل ۱ نمایانگر آن است که یک مسیر وضعیت مختلف دارای کمی اختلاف در امتیاز با حالتی است که $5'SS$ در پنجمین G قرار می گیرد (لگاریتم احتمال ۴۱.۲۱ - در مقابل ۴۱.۲۲). چقدر اطمینان داریم که پنجمین G، انتخاب درستی است؟ این از ویژگی های الگوهای احتمالی است ما می توانیم اطمینان مان را مستقیم محاسبه کنیم. احتمال این که یک عضو توالی مثل π توسط مسیرهای وضعیت که از وضعیت K برای انتشار استفاده می کنند (یعنی $\pi_i = K$ در مسیر وضعیت) تقسیم بر مجموع همه مسیرهای وضعیت ممکن. در مثال نمایشی ما، این یک مسیر خاص در

مصطفی‌احبیه

مصطفی‌احبیه با لاسلوب لواس



لاسلوب لواس در حال حاضر رئیس اتحادیه بین‌المللی ریاضیات (IMU) است. او در سال ۱۹۴۸ در بوداپست مجارستان به دنیا آمده و تمام تحصیلات خود را در مجارستان گذرانده است. لواس یکی از افراد سرشناس در ترکیبیات می‌باشد و به خاطر کارهای عمیقی که انجام داده، جوایز ول芙 و کنوت را در سال ۱۹۹۹ و جایزه بولیابی را در سال ۲۰۰۷ از آن خود کرده است. لواس در طول تحصیل در دیبرستان ۳ بار برندۀ مدار طلا از المپیاد جهانی ریاضیات شده است. پسرو از نیز در سال ۲۰۰۸ موفق به کسب این مدار گردید.

لواس در بهار ۱۳۸۶ به دعوت پژوهشکده ریاضیات پژوهشگاه دانش‌های بنیادی در کنفرانس نظریه جبری گراف شرکت کرد و چند سخنرانی ارائه داد. در خلال این کنفرانس از ایشان برای یک مصاحبه اختصاصی با خبرنامه دعوت به عمل آمد که با گشاده‌رویی دعوت را پذیرفتند.

این مصاحبه توسط آقایان دکتر: مهدی بهرزاد، سعید اکبری، منوچهر ذاکر و رشید زارع نهنده انجام شد. از پژوهشکده ریاضیات که امکانات این مصاحبه را فراهم کرد و از خانم عاطفه پارسا به خاطر ضبط مصاحبه و عکس‌ها و آقای فرزین منیعی و خانم سیده صادقه حق شناس که در پیاده‌سازی قسمتهایی از مصاحبه کمک کردند، سپس گزاری می‌شود. پیاده‌سازی نهایی و ترجمه مصاحبه توسط منوچهر ذاکر و رشید زارع نهنده انجام شده است.

• اکبری: نخست از این که دعوت ما را برای مصاحبه پذیرفتید، متشرکریم. لطفاً اگر خاطراتی از دوران نوجوانی خود و نحوه آشنایی با افرادی مانند اردوش (P. Erdős)، پلیکان (J. Pelikan) و سایر ترکیبیات‌دانان زیده دارید، بفرمایید.

اولین قدم من در ورود به ریاضیات زمانی بود که در دیبرستان بودم. آن زمان در مجارستان کلاس‌های ویژه‌ای برای دانش‌آموزان مستعد برپا کرده بودند و من هم وارد یکی از آنها شدم. دانش‌آموزان بسیار خوبی در آن کلاس بودند و من خیلی خوش‌شانس بودم که فرست هم کلاسی با آنان را داشتم. استادانی از دانشگاه به این کلاس‌ها می‌آمدند و درس‌های ویژه‌ای ارائه می‌کردند و موضوعات

سه‌تایی می‌باشند. به هر حال، تا الان HMM‌های پایه‌ای به کار گرفته شده‌اند. کلاس‌های الگوهای احتمالی پرتوانتری (اگرچه با کارایی کمتر) از HMM‌ها برای تجزیه و تحلیل توالی‌ها وجود دارند.

مراجع

1. Rabiner, L.R. A Tutorial on hidden Markov models and selected applications in speech recognition. Proc. IEEE 77, 257-286(1989)
2. Durbin, R., Eddy, S.R., Krogh, A. and Mitchinson, G.J. Biological sequence analysis: probabilistic models of proteins And nucleic acids (Cambridge University press, Cambridge UK. 1998).

این نوشه برگردانی از مقاله زیر است:

Eddy, S., R. What is a Hidden Markov Model? Nature Biotechnology Vol. 22 No. 10 October 2004.

مترجم: حمید پژشک

بخش آمار دانشگاه تهران

و هسته بیوانفورماتیک پژوهشگاه دانش‌های بنیادی



حق عضویت در انجمن ریاضی ایران

دوره ۸۷ - ۸۸

- اعضای پیوسته انجمن ۵۰۰ / ۰۰۰ ریال.
- اعضای وابسته با دریافت هر سه نشریه ۱۲۰ / ۰۰۰ ریال.
- اعضای وابسته با دریافت بولتن و خبرنامه ۱۰۰ / ۰۰۰ ریال.
- اعضای وابسته با دریافت خبرنامه و فرهنگ و اندیشه ریاضی ۱۰۰ / ۰۰۰ ریال.
- اعضای وابسته با دریافت فقط خبرنامه ۶۰ / ۰۰۰ ریال. کلیه دانش‌آموزان، دانشجویان، معلمان سطوح مختلف آموزش و پرورش و اعضای انجمن آمار ایران، انجمن ریاضی فرانسه و انجمن ریاضی آمریکا می‌توانند از تخفیف ۵۰٪ استفاده کنند.