

مثالهایی از نحوه استفاده از الگوریتم EM در محاسبه برآوردهای درستنمایی ماکزیمم

دکتر ناصر رضا ارقامی^۲

سید محمد ابراهیم حسینی نسب^۱

چکیده

گاهی اوقات، در بدست آوردن برآوردهای درستنمایی ماکزیمم، پس از مشتق‌گیری از تابع درستنمایی به معادلاتی می‌رسیم که نمی‌توان از آنها فرم بسته‌ای را برای برآوردهای درستنمایی ماکزیمم نتیجه گرفت. در این گونه موارد الگوریتم EM که یک فن تکراریست و برای اولین بار توسط Dempster و سایرین در سال ۱۹۹۷ ارائه گردید، یک راه عملی برای یافتن جوابهای درستنمایی ماکزیمم پیش پای ما می‌گذارد.

معمولاً، در منابع درسی هنگام معرفی الگوریتم EM کوششی در جهت روشن کردن منطقی که این الگوریتم بر آن استوار است و ارائه دلایلی که چرا این الگوریتم ما را به برآوردهای درستنمایی ماکزیمم می‌رساند به عمل نمی‌آید. در این مقاله سعی ما بر این است که علاوه بر توضیح موارد فوق، مثالهای ساده‌ای که به درک عمیقتر این الگوریتم کمک می‌کند ارائه نماییم.

۱ مقدمه

هر تکرار الگوریتم EM، شامل دو مرحله است: مرحله E و مرحله M. در مرحله E امید ریاضی شرطی، $Q(\theta|\theta^{(n)}) = E[\ln f(X|\theta)|y; \theta^{(n)}]$ تشکیل و در مرحله M این امید شرطی نسبت به θ ماکزیمم می‌شود. $\theta^{(n)}$ در عبارت امید ریاضی شرطی نشان‌دهنده مقدار برآورد بدست آمده از مرحله تکرار $\theta^{(n)}$ و $\theta^{(n+1)}$ مقداری از θ است که امید ریاضی شرطی فوق را با فرض ثابت بودن $\theta^{(n)}$ ماکزیمم می‌کند و در نتیجه

فرض کنید مشاهدات ما بردار Y با تابع چگالی احتمال $g(y|\theta)$ که θ پارامتر مورد نظر است باشد. در این الگوریتم، فرض بر این است که بردار تصادفی X با تابع چگالی احتمال $f(x|\theta)$ به گونه‌ای است که Y تابعی از X مانند $h(X)$ باشد در این صورت به X داده‌های کامل و به Y داده‌های ناقص گویند. معمولاً Y را داده‌های قابل مشاهده و X را داده‌های غیر قابل مشاهده نامند.

^۱ سید محمد ابراهیم حسینی نسب، دانشجوی کارشناسی ارشد آمار دانشگاه فردوسی مشهد
^۲ دکتر ناصر رضا ارقامی عضو هیأت علمی گروه آمار دانشگاه فردوسی مشهد

$$\begin{aligned}
 &= E_x \left\{ \ln \left[\frac{f(X|\underline{\theta})}{g(Y|\underline{\theta})} \right] \middle| y; \underline{\theta}^{(n)} \right\} \\
 &= \int_{\{x; h(x)=y\}} \ln \left[\frac{f(x|\underline{\theta})}{g(y|\underline{\theta})} \right] f(x|y; \underline{\theta}^{(n)}) dx \\
 &= \int_{\{x; h(x)=y\}} \ln \left[\frac{f(x|\underline{\theta})}{g(y|\underline{\theta})} \right] \left[\frac{f(x|\underline{\theta}^{(n)})}{g(y|\underline{\theta}^{(n)})} \right] dx
 \end{aligned}$$

بنابه لم ۱ می توان نوشت:

$$\begin{aligned}
 &\int_{\{x; h(x)=y\}} \ln \left[\frac{f(x|\underline{\theta})}{g(y|\underline{\theta})} \right] \left[\frac{f(x|\underline{\theta}^{(n)})}{g(y|\underline{\theta}^{(n)})} \right] dx \\
 &\leq \int_{\{x; h(x)=y\}} \ln \left[\frac{f(x|\underline{\theta}^{(n)})}{g(y|\underline{\theta}^{(n)})} \right] \left[\frac{f(x|\underline{\theta}^{(n)})}{g(y|\underline{\theta}^{(n)})} \right] dx \\
 &= \int_{\{x; h(x)=y\}} [\ln f(x|\underline{\theta}^{(n)}) - \ln g(y|\underline{\theta}^{(n)})] f(x|y; \underline{\theta}^{(n)}) dx \\
 &= E_x [\ln g(Y|\underline{\theta}^{(n)}) | y; \underline{\theta}^{(n)}] - \ln g(y|\underline{\theta}^{(n)})
 \end{aligned}$$

یعنی این که :

$$H(\underline{\theta}|\underline{\theta}^{(n)}) \leq H(\underline{\theta}^{(n)}|\underline{\theta}^{(n)})$$

با توجه به قضیه ۱ می توان نوشت:

$$\begin{aligned}
 H(\underline{\theta}^{(n)}|\underline{\theta}^{(n)}) &\geq H(\underline{\theta}^{(n+1)}|\underline{\theta}^{(n)}) \Rightarrow \\
 E_x [\ln f(X|\underline{\theta}^{(n)}) | y; \underline{\theta}^{(n)}] - \ln g(y|\underline{\theta}^{(n)}) &\geq \\
 E_x [\ln f(X|\underline{\theta}^{(n+1)}) | y; \underline{\theta}^{(n)}] - \ln g(y|\underline{\theta}^{(n+1)}) &\Rightarrow \\
 E_x [\ln f(X|\underline{\theta}^{(n)}) | y; \underline{\theta}^{(n)}] - E_x [\ln f(X|\underline{\theta}^{(n+1)}) | y; \underline{\theta}^{(n)}] & \\
 \geq \ln g(y|\underline{\theta}^{(n)}) - \ln g(y|\underline{\theta}^{(n+1)}) & \quad (1.1)
 \end{aligned}$$

اما چون در تکرار $(n+1)$ ام، $\underline{\theta}^{(n+1)}$ مقداری از $\underline{\theta}$ است که
 $E[\ln f(X|\underline{\theta}) | y; \underline{\theta}^{(n)}]$ را ماکزیمم می کند لذا:

$$E_x [\ln f(X|\underline{\theta}^{(n)}) | y; \underline{\theta}^{(n)}] \leq E_x [\ln f(X|\underline{\theta}^{(n+1)}) | y; \underline{\theta}^{(n)}]$$

$\underline{\theta}^{(n+1)}$ مقدار برآورد به دست آمده از مرحله $(n+1)$ ام است.
 توضیح این که چرا الگوریتم EM منجر به بیشینه شدن تابع درستنمایی می شود، با تعریف تابع $H(\underline{\theta}|\underline{\theta}^{(n)})$ با ضابطه

$$H(\underline{\theta}|\underline{\theta}^{(n)}) = E_x [\ln f(X|\underline{\theta}) | y; \underline{\theta}^{(n)}] - \ln g(y|\underline{\theta})$$

تسهیل می یابد و نقش اساسی این الگوریتم که همان افزایش درستنمایی از مرحله ای به مرحله بعد است بهتر روشن می گردد. در این تابع مقادیر y و $\underline{\theta}^{(n)}$ معلوم فرض می شود، X متغیر تصادفی و $\underline{\theta}$ متغیر غیر تصادفی است. ابتدا ثابت می کنیم این تابع ماکزیمم مقدارش را در $\underline{\theta} = \underline{\theta}^{(n)}$ اختیار می کند، سپس نشان خواهیم داد که مقدار تابع درستنمایی $g(y|\underline{\theta})$ به ازای $\underline{\theta} = \underline{\theta}^{(n+1)}$ بیشتر از مقدار $g(y|\underline{\theta})$ به ازای $\underline{\theta} = \underline{\theta}^{(n)}$ است و در نتیجه در هر مرحله، مقدار تابع درستنمایی افزایش می یابد.

ابتدا به بیان یک لم و اثبات یک قضیه می پردازیم:

لم ۱ (نامساوی آنتروپی): فرض کنید f و g دو تابع چگالی احتمالی، نسبت به اندازه μ و همچنین هر دو تقریباً همه جا (a.e) نسبت به μ مثبت باشند. اگر امید ریاضی نسبت به اندازه احتمال $f d\mu$ باشد، آنگاه $E_f(\ln f) \geq E_f(\ln g)$ و تساوی برقرار است اگر و تنها اگر $f = ga.e\mu$ (اثبات در مرجع [۴])

قضیه ۱ (مرجع [۳]): $H(\underline{\theta}|\underline{\theta}^{(n)})$ در نقطه $\underline{\theta} = \underline{\theta}^{(n)}$ ماکزیمم می شود.

اثبات :

$$\begin{aligned}
 E_x [\ln g(Y|\underline{\theta}) | y; \underline{\theta}^{(n)}] &= \int_{\{x; h(x)=y\}} \ln g(y|\underline{\theta}) f(x|y; \underline{\theta}^{(n)}) dx \\
 &= \ln g(y|\underline{\theta}) \int_{\{x; h(x)=y\}} f(x|y; \underline{\theta}^{(n)}) dx \\
 &= \ln g(y|\underline{\theta})
 \end{aligned}$$

بنابراین

$$H(\underline{\theta}|\underline{\theta}^{(n)}) = E_x [\ln f(X|\underline{\theta}) | y; \underline{\theta}^{(n)}] - E_x [\ln g(Y|\underline{\theta}) | y; \underline{\theta}^{(n)}]$$

مطالب گفته شده بالا را بطور خلاصه نشان می‌دهد.

آنتی ژن	گونه‌های ژنی
A	A A و A O
B	B B و B O
AB	A B
O	O O

در این مسأله، تعداد افراد مشاهده شده از دسته‌های سمت چپ جدول، داده‌های Y را تشکیل می‌دهند در صورتی که تعداد افرادی که هریک از ۶ گونه ژنی سمت راست را دارند و غیر مشخص و نامعلوم هستند داده‌های X را مشخص می‌کنند. فرض کنید تعداد افرادی که دارای گونه‌های $A|A, A|O, B|B, B|O, A|B, O|O$ هستند به ترتیب با x_1, x_2, \dots, x_6 نشان دهیم. n حجم نمونه گرفته شده ما و y_1, y_2, y_3, y_4 به ترتیب تعداد افراد مشاهده شده برای دسته‌های A, B, AB, O هستند که $n = y_1 + y_2 + y_3 + y_4$. توجه کنید که در این جا $y_1 = x_1 + x_2$ ، $y_2 = x_3 + x_4$ ، $y_3 = x_5$ و $y_4 = x_6$ می‌باشد.

بنابراین داریم:

$$f(x|p) = \binom{n}{x_1, x_2, \dots, x_6} (P_A^y)^{x_1} (2P_{AP}O)^{x_2} \times (P_B^y)^{x_3} (2P_{BP}O)^{x_4} (2P_{APB})^{x_5} (P_O^y)^{x_6}$$

که $P_A^y, 2P_{APB}, 2P_{BP}O, P_B^y, 2P_{AP}O, P_A^y$ به ترتیب احتمالات متناظر با گونه‌های ژنی بالا است. در گام E از الگوریتم EM ما باید $E[\ln f(X|P)|y; P^{(m)}]$ را تشکیل دهیم که $P^{(m)} = (P_A^{(m)}, P_B^{(m)}, P_O^{(m)})^T$ بردار جاری P می‌باشد.

$$\begin{aligned} Q(P|P^{(m)}) &= E(\ln f(X|P)|y; P^{(m)}) \\ &= x_1^{(m)} \ln(P_A^y) + x_2^{(m)} \ln(2P_{AP}O) \\ &+ x_3^{(m)} \ln(P_B^y) + x_4^{(m)} \ln(2P_{BP}O) \\ &+ y_3 \ln(2P_{APB}) + y_4 \ln(P_O^y) \end{aligned}$$

لذا از نامساوی (۱.۱) داریم:

$$\ln g(y|\underline{\theta}^{(n)}) \leq \ln g(y|\underline{\theta}^{(n+1)})$$

به عبارت دیگر، الگوریتم EM باعث افزایش لگاریتم درست‌نمایی از مرحله‌ای به مرحله بعد از آن می‌شود و این روند تا همگرا شدن الگوریتم به جواب، که ثابت می‌شود تحت شرایطی همان برآورد درست‌نمایی ماکزیمم است ادامه خواهد داشت. لازم به ذکر است که همگرایی، به انتخاب نخستین $\underline{\theta}$ (یعنی $\underline{\theta}^{(0)}$) برای شروع بستگی ندارد.

تذکر ۱: هنر استفاده از الگوریتم EM در انتخاب و تشخیص داده‌های کامل X است. اگر چه روشهای زیادی برای محاط کردن Y در یک فضای نمونه بزرگتر وجود دارد اما اغلب طبیعت مسئله یا ملاحظات دیگر به یک تعریف روشن و صریح از X منجر می‌شود.

مثالهای زیر ما را در جهت استفاده صحیح از این الگوریتم با در نظر گرفتن نکات و ریزه‌کاری‌های آن یاری می‌دهد.

مثال ۱: فرض کنید هریک از والدین دارای سه نوع ژن A, B, O با احتمالات P_A و P_B و P_O که $P_A + P_B + P_O = 1$ باشند. بر اساس اصول ژنتیک، گونه‌های ژنی بوجود آمده توسط آنها عبارتند از: $A|A, A|B, A|O, B|B, B|O, O|O$ است. به عنوان مثال گونه ژنی $A|B$ نشان دهنده آنست که ژن دریافتی از مادر A و ژن دریافتی از پدر B می‌باشد یا بالعکس.

بر اساس نمونه‌های خون گزفته شده از افراد، ما قادر به مشاهده فراوانی گونه‌های ژنی بالا نمی‌باشیم اما با تحت تأثیر قرار دادن نمونه‌های خونی بوسیله آنتی‌ژنهای A و B می‌توان اطلاعاتی درباره گونه‌های ژنی بدست آورد. به این ترتیب که در هر نمونه اگر آنتی‌ژن A به تنهایی نمایان شود یکی از گونه‌های $A|A$ یا $A|O$ ، اگر آنتی‌ژن B به تنهایی ظاهر شود یکی از گونه‌های $B|B$ یا $B|O$ ، اگر آنتی‌ژنهای A و B هر دو با هم نمایان شوند گونه $A|B$ و اگر هیچ کدام از آنتی‌ژنهای A یا B نمایان نگردد گونه $O|O$ مشخص خواهد شد. جدول زیر

آمیخته پواسن، تابع درستنمایی مشاهدات به صورت زیر است:

$$L(y; \theta) = \prod_{i=0}^9 [\alpha e^{-\mu_1} \frac{\mu_1^i}{i!} + (1-\alpha) \frac{e^{-\mu_2} \mu_2^i}{i!}]^{y_i}$$

که α پارامتر ترکیب و μ_1 و μ_2 میانگین‌های دو توزیع پواسن و $\theta = (\alpha, \mu_1, \mu_2)^T$ هستند.

تعداد مرگ و میر (i)	فراوانی (y _i)	تعداد مرگ و میر (i)	فراوانی (y _i)
۰	۱۶۲	۵	۶۱
۱	۲۶۷	۶	۲۷
۲	۲۷۱	۷	۸
۳	۱۸۵	۸	۳
۴	۱۱۱	۹	۱

اگر فرض کنیم $Z_i(\theta)$ احتمال پسین آن که یک روز با i مرگ و میر به جامعه پواسن ۱ متعلق باشد، آنگاه $1 - Z_i(\theta)$ احتمال پسین تعلق یک روز با i مرگ و میر به جامعه پواسن ۲ خواهد بود و به صورت زیر بدست می‌آید:

$$Z_i(\theta) = \frac{\alpha e^{-\mu_1} \mu_1^i}{\alpha e^{-\mu_1} \mu_1^i + (1-\alpha) e^{-\mu_2} \mu_2^i} \quad i = 0, 1, 2, \dots, 9$$

با تعریف H به صورت:

$$H = \begin{cases} 1 & \text{اگر روزی با هر تعداد مرگ و میر} \\ & \text{از جامعه پواسن ۱ باشد.} \\ 0 & \text{در غیر اینصورت} \end{cases}$$

$$P(H=1) = \alpha \text{ داریم}$$

در این مثال $X = (U, H)$ ، که متغیر تصادفی U تعداد مرگ و میر است و مقادیر $0, 1, 2, \dots, 9$ را می‌گیرد.

$$\begin{aligned} f(x|\theta) &= f_{\theta}(u, h) \\ &= (\alpha g_1(u))^h [(1-\alpha)g_2(u)]^{1-h} \quad h = 0, 1 \end{aligned}$$

که

که در آن $x_i^{(m)} = E(X_i|y; P^{(m)})$ ، $i = 1, 2, 3, 4, 5, 6$ و $x_4^{(m)} = y_4, x_5^{(m)} = y_5$ است.

اما در گام M از الگوریتم EM تابع $Q(P|P^{(m)})$ ماکزیمم می‌شود، که در این جا بدلیل قید $P_A + P_B + P_O = 1$ این کار با استفاده از روش لاگرانژ انجام می‌گردد.

$$H(P, \lambda) = Q(P|P^{(m)}) + \lambda(P_A + P_B + P_O - 1)$$

اگر نسبت به P_A, P_B, P_O و λ مشتق گرفته و مساوی صفر قرار دهیم، پس از ساده کردن داریم:

$$\begin{aligned} P_A^{(m+1)} &= \frac{2x_1^{(m)} + x_2^{(m)} + y_3}{2n} \\ P_B^{(m+1)} &= \frac{2x_3^{(m)} + x_4^{(m)} + y_3}{2n} \\ P_O^{(m+1)} &= \frac{x_5^{(m)} + x_6^{(m)} + 2y_4}{2n} \end{aligned}$$

بر اساس داده‌های $y_1 = 186, y_2 = 38, y_3 = 13$ و $y_4 = 284$ و با بهره‌گیری از رایانه، برآوردهای زیر حاصل می‌شوند:

$$\hat{P}_A = 0/2136, \hat{P}_B = 0/0501, \hat{P}_O = 0/7363$$

(این مثال در مرجع [۴] آمده است).

مثال ۲: داده‌های جدول زیر مربوط به سالهای ۱۹۱۰ تا ۱۹۱۲ شهر لندن است. ستون اول این جدول تعداد مرگ و میر در بین زنان ۸۰ و بیشتر از ۸۰ سال را نشان می‌دهد، ستون دوم y ها یعنی فراوانی روزهایی را که در آن روزها i مرگ و میر داشته‌ایم مشخص می‌سازد. بدلیل الگوهای متفاوت مرگ و میر در زمستان و تابستان، یک توزیع پواسن به تنهایی نمی‌تواند برازش خوبی بر داده‌ها داشته باشد، لذا به نظر می‌رسد ترکیبی از دو توزیع پواسن، برازش بهتری را ایجاد کند. تحت مدل

ساده کردن و در نهایت مشتق‌گیری داریم:

$$\alpha^{(m+1)} = \frac{\sum_{i=0}^9 y_i Z_i(\underline{\theta}^{(m)})}{\sum_{i=0}^9 y_i}$$

$$\mu_1^{(m+1)} = \frac{\sum_{i=0}^9 y_i Z_i(\underline{\theta}^{(m)}) i}{\sum_{i=0}^9 y_i Z_i(\underline{\theta}^{(m)})}$$

$$\mu_2^{(m+1)} = \frac{\sum_{i=0}^9 y_i (1 - Z_i(\underline{\theta}^{(m)})) i}{\sum_{i=0}^9 y_i (1 - Z_i(\underline{\theta}^{(m)}))}$$

$$g_j(u) = \frac{e^{-\mu_j} (\mu_j)^u}{u!} \quad j = 1, 2$$

$$A(u) = \frac{g_1(u)}{g_2(u)} \quad u = 0, 1, \dots, 9$$

لذا:

$$\begin{aligned} f(x_1, \dots, x_n) &= \prod_{r=1}^n f(x_r | \underline{\theta}) \\ &= (1 - \alpha)^{\sum_{i=0}^9 y_i} \left(\frac{\alpha}{1 - \alpha}\right)^{\sum_{i=0}^9 q_i} \\ &\quad \times \prod_{i=0}^9 [g_2(i)]^{y_i} [A(i)]^{q_i} \end{aligned}$$

با استفاده از رایانه جوابهای عددی بدست آمده عبارتند از:

$$\hat{\alpha} = 0/3599, \quad \hat{\mu}_1 = 1/2561, \quad \hat{\mu}_2 = 2/6634$$

تذکر ۲: در این دو مثال، بوضوح می‌توان دید که Y تابعی از X است. مجدداً تأکید می‌شود که مشاهدات X در دسترس نیستند و تنها مشاهدات Y است که در اختیار ما قرار دارد.

تذکر ۳: گاهی اوقات ممکن است بنا به دلایلی اجرای گام M از الگوریتم EM امکان‌پذیر نباشد، در این گونه مواقع می‌توان از الگوریتم GEM (الگوریتم EM تعمیم یافته) که آن نیز برای نخستین بار توسط Dempster و سایرین در سال ۱۹۷۷ (در [۱]) ارائه شده، استفاده نمود. روند این الگوریتم طوری است که:

$$Q(\underline{\theta}^{(n+1)} | \underline{\theta}^{(n)}) \geq Q(\underline{\theta}^{(n)} | \underline{\theta}^{(n)})$$

بنابراین در این جا نیز از مرحله‌ای به مرحله بعد افزایش مقدار تابع درستنمایی را خواهیم داشت و ملاحظه می‌شود که الگوریتم EM حالت خاصی از الگوریتم GEM می‌باشد.

تذکر ۴: الگوریتم EM و GEM موارد استفاده فراوانی دارند که به عنوان مثال می‌توان به استفاده از آنها در بازسازی تصویر درسی‌تی‌اسکن، مسایل ژنتیک و نیز مسایلی که در آن داده گمشده وجود دارد، اشاره نمود.

و در آن $i = 0, 1, 2, \dots, 9$ ، $S_i = \{r | u_r = i\}$ و $N(S_i) = y_i$ می‌باشند، واضح است که $q_i = \sum_{r \in S_i} H_r$

$$\begin{aligned} E[\ln f(X | \underline{\theta}) ; \underline{\theta}^{(m)}] &= \sum_{i=0}^9 y_i [\ln(1 - \alpha) + \ln g_2(i)] \\ &\quad + \sum_{i=0}^9 [\ln\left(\frac{\alpha}{1 - \alpha}\right) + \ln A(i)] E[q_i | y; \underline{\theta}^{(m)}] \end{aligned}$$

و اما:

$$\begin{aligned} E[q_i | y; \underline{\theta}^{(m)}] &= \sum_{r \in S_i} E[H_r | y; \underline{\theta}^{(m)}] \\ &= \sum_{r \in S_i} Z_i(\underline{\theta}^{(m)}) \\ &= y_i Z_i(\underline{\theta}^{(m)}) \end{aligned}$$

که $\underline{\theta}^{(m)} = (\alpha^{(m)}, \mu_1^{(m)}, \mu_2^{(m)})^T$ پارامترهای بدست آمده از تکرار m الگوریتم (پارامتر جاری) است. با مقدارگذاری به جای $A(i)$ و $g_2(i)$ در امید شرطی بالا،

- [1] Dempster AP , Laird NM, Rubin DB. Maximum likelihood from incomplete data via the EM algorithm. *JR Stat Soc Series B* 1977; 39: 1-38.
 - [2] Wu CF. On the convergence properties of the EM algorithm. *Ann stat* 1983; 11: 95-103.
 - [3] K.Lange and R.Carson, EM Reconstruction Algorithms for Emission and Transmission Tomography. *J. of Computer Assisted Tomography*, Vol. 8, No. 2, 1984.
 - [4] K. Lange. *Mathematical and Statistical Methods for Genetic Analysis*. Springer, 1997, p 22-32.
-