

## مثالهایی از نحوه استفاده از الگوریتم EM در محاسبه برآوردهای درستنایی ماکزیمم

سید محمد ابراهیم حسینی نسب<sup>۱</sup> دکتر ناصر رضا ارقامی<sup>۲</sup>

### چکیده

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

معمولًا، در منابع درسی هنگام معرفی الگوریتم EM کوششی در جهت روشن کردن منطقی که این الگوریتم بر آن استوار است و ارایه دلایلی که چرا این الگوریتم ما را به برآوردهای درستنایی ماکزیمم می‌رساند به عمل نمی‌آید. در این مقاله سعی می‌کنیم این مفهوم را بوضیح موارد فوق، مثالهای ساده‌ای که به درک عمیق‌تر این الگوریتم کمک می‌کند ارایه نماییم.

### ۱ مقدمه

هر تکرار الگوریتم EM، شامل دو مرحله است: مرحله E و مرحله M. در مرحله E امید ریاضی شرطی،  $E[\ln f(X|\theta)|y; \theta^{(n)}] = E[\ln f(X|\theta)|y; \theta^{(n)}]$  تشکیل و در مرحله M این امید شرطی نسبت به  $\theta$  ماکزیمم می‌شود. در عبارت امید ریاضی شرطی نشان دهنده مقدار برآورد بدست آمده از مرحله تکرار (n)ام و  $(n+1)$ مقداری از  $\theta$  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  $f(y|\theta)$  که  $\theta$  پارامتر مورد نظر است باشد. در این الگوریتم، فرض براین است که بردار تصادفی X با تابع چگالی احتمال  $f(x|\theta)$  به گونه‌ای است که Y تابعی از X مانند  $f(Y|X)$  باشد در این صورت به X داده‌های کامل و به Y داده‌های ناقص گویند. معمولًا Y را داده‌های قابل مشاهده و X را داده‌های غیرقابل مشاهده نامند.

<sup>۱</sup>سید محمد ابراهیم حسینی نسب، دانشجوی کارشناسی ارشد آمار دانشگاه فردوسی مشهد

<sup>۲</sup>دکتر ناصر رضا ارقامی عضو هیأت علمی گروه آمار دانشگاه فردوسی مشهد

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

$\underline{\theta}^{(n+1)}$  مقدار براورد به دست آمده از مرحله  $(n+1)$  است. توضیح این که چرا الگوریتم EM منجر به بیشینه شدن تابع درستنمایی می‌شود، با تعریف تابع  $H(\underline{\theta}|\underline{\theta}^{(n)})$  با ضابطه  $H(\underline{\theta}|\underline{\theta}^{(n)}) = E_x[\ln f(X|\underline{\theta})|y; \underline{\theta}^{(n)}] - \ln g(y|\underline{\theta})$

بنابراین ۱ می‌توان نوشت:

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

تسهیل می‌یابد و نقش اساسی این الگوریتم که همان افزایش درستنمایی از مرحله‌ای به مرحله بعد است بهتر روش می‌گردد. در این تابع مقادیر  $y$  و  $\underline{\theta}^{(n)}$  معلوم فرض می‌شود،  $X$  متغیر تصادفی و  $\underline{\theta}$  متغیر غیر تصادفی است. ابتدا ثابت می‌کنیم این تابع ماکزیمم مقدارش را در  $\underline{\theta}^{(n+1)} = \underline{\theta} = \underline{\theta}^{(n)}$  اختیار می‌کند، سپس نشان خواهیم داد که مقدار تابع درستنمایی  $H(\underline{\theta}|\underline{\theta}^{(n)})$  به ازای  $\underline{\theta} = \underline{\theta}^{(n+1)}$  بیشتر از مقدار  $H(\underline{\theta}|\underline{\theta}^{(n)})$  به ازای  $\underline{\theta} = \underline{\theta}^{(n)}$  است و در نتیجه در هر مرحله، مقدار تابع درستنمایی افزایش می‌یابد.

ابتدا به بیان یک لم و اثبات یک قضیه می‌پردازیم:

لم ۱ (نامساوی آنتروپی): فرض کنید  $f$  و  $g$  دو تابع چگالی احتمالی، نسبت به اندازه  $\mu$  و همچنین هر دو تقریباً همه جا (a.e) نسبت به  $\mu$  مثبت باشند. اگر  $E_f(\ln f) \geq E_f(\ln g)$  نسبت به اندازه احتمال  $\mu$  باشد، آنگاه  $f d\mu \geq g d\mu$  و تساوی برقرار است اگر و تنها اگر  $f = g a.e^\mu$  (اثبات در مرجع [۳])

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

اثبات :

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

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

بنابراین

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

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

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

مطلوب 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  $O|O, A|B, B|O, B|B, A|O, A|A$  نشان دهیم.  $n$  حجم نمونه گرفته شده ما و  $y_1, y_2, y_3, y_4, y_5, y_6$  به ترتیب تعداد افراد مشاهده شده برای دسته‌های  $O, AB, B, A$  هستند که  $n = y_1 + y_2 + y_3 + y_4 + y_5 + y_6$ . توجه کنید که در اینجا  $y_1 = x_1 + x_2$ ،  $y_2 = x_2 + x_4$ ،  $y_3 = x_1 + x_4$ ،  $y_4 = x_5$  و  $y_5 = x_6$  می‌باشد.

بنابراین داریم:

$$f(x|p) = \binom{n}{x_1, x_2, \dots, x_6} (P_A)^{x_1} (2P_AP_O)^{x_2} \\ \times (P_B)^{x_3} (2P_BP_O)^{x_4} (2P_AP_B)^{x_5} (P_O)^{x_6}$$

که  $P_A, P_B, P_O$  به ترتیب احتمالات متناظر با گونه‌های ژنی بالا است. در گام 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 می‌باشد.

$$Q(P|P^{(m)}) = E[\ln f(X|P)|y; P^{(m)}] \\ = x_1^{(m)} \ln(P_A) + x_2^{(m)} \ln(2P_AP_O) \\ + x_3^{(m)} \ln(P_B) + x_4^{(m)} \ln(2P_BP_O) \\ + y_1 \ln(2P_AP_B) + y_2 \ln(P_O)$$

لذا از نامساوی (۱.۱) داریم:

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

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  $\theta$  (یعنی  $\theta^{(0)}$ ) برای شروع بستگی ندارد.

نذکر 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 درنظر گرفتن نکات و ریزه‌کاری‌های آن پاری می‌دهد.

مثال ۱: فرض کنید هریک از والدین دارای سه نوع ژن با احتمالات  $P_A, P_B$  و  $P_O$  که  $P_A + P_B + P_O = 1$  باشند. بر اساس اصول ژنتیک، گونه‌های ژنی بوجود آمده توسط آنها عبارتند از:  $O|O, B|B, B|O, A|O, A|B, A|A$  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 گرفته شده از افراد، ما قادر به مشاهده فراوانی گونه‌های ژنی بالا نمی‌باشیم اما با تحت تأثیر قرار دادن نمونه‌های خونی بوسیله آنتی ژنهای  $A$  و  $B$  می‌توان اطلاعاتی درباره گونه‌های ژنی بدست آورد. به این ترتیب که در هر نمونه اگر آنتی ژن  $A$  به تنها یی نمایان شود یکی از گونه‌های  $A|O$  یا  $A|A$ ، اگر آنتی ژن  $B$  به تنها یی نمایان شود یکی از گونه‌های  $B|B$  یا  $B|O$ ، اگر آنتی ژنهای  $A$  و  $B$  هردو با هم نمایان شوند گونه  $A|B$  و اگر هیچ کدام از آنتی ژنهای  $A$  یا  $B$  نمایان نگردد گونه  $O|O$  مشخص خواهد شد. جدول زیر

آمیخته پواسن،تابع درستنمایی مشاهدات به صورت زیر است:

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

که  $\alpha$  پارامتر ترکیب و  $\mu_1$  و  $\mu_2$  میانگین های دو توزیع پواسن و  $\underline{\theta} = (\alpha, \mu_1, \mu_2)^T$  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  $\theta$  مرگ و میر به جامعه پواسن ۱ متعلق باشد، آنگاه  $Z_i(\theta) = Z_i - 1$  احتمال پسین تعلق یک روز با  $\theta$  مرگ و میر به جامعه پواسن ۲ خواهد بود و به صورت زیر بدست می آید:

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

با تعریف  $H$  به صورت:

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

$$P(H = 1) = \alpha$$

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

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

که

که در آن  $x_i^{(m)} = E(X_i|y; P^{(m)})$ ،  $i = 1, 2, 3, 4, 5$  و  $x_6^{(m)} = y_4$ ،  $x_5^{(m)} = y_3$  است.

اما در گام  $M$  از الگوریتم EM تابع  $Q(P|P^{(m)})$  ماکریم می شود، که در اینجا بدلیل قید  $P_A + P_B + P_O = 1$  کار با استفاده از روش لاگرانژ انجام می گردد.

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

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

$$\begin{aligned} P_A^{(m+1)} &= \frac{2x_1^{(m)} + x_2^{(m)} + y_3}{2n} \\ P_B^{(m+1)} &= \frac{2x_3^{(m)} + x_4^{(m)} + y_2}{2n} \\ P_O^{(m+1)} &= \frac{x_2^{(m)} + x_5^{(m)} + 2y_4}{2n} \end{aligned}$$

بر اساس داده های  $y_1 = 186, y_2 = 13, y_3 = 38, y_4 = 284$  و با بهره گیری از ریاضی، برآوردهای زیر حاصل می شوند:

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

(این مثال در مرجع [۴] آمده است).

مثال ۲: داده های جدول زیر مربوط به سالهای ۱۹۱۰ تا ۱۹۱۲ شهر لندن است. ستون اول این جدول تعداد مرگ و میر در بین زنان  $80$  و بیشتر از  $80$  سال را نشان می دهد، ستون دوم، یعنی فراوانی روزهایی را که در آن روزها مرگ و میر داشتایم مشخص می سازد. بدلیل الگوهای متفاوت مرگ و میر در زمستان و تابستان، یک توزیع پواسن به تنها بی نمی تواند برآش خوبی بر داده ها داشته باشد، لذا به نظر می رسد ترکیبی از دو توزیع پواسن، برآش بهتری را ایجاد کند. تحت مدل

ساده کردن و در نهایت مشتق‌گیری داریم:

$$\begin{aligned}\alpha^{(m+1)} &= \frac{\sum_{i=0}^9 y_i Z_i(\underline{\theta}^{(m)})}{\sum_{i=0}^9 y_i} & g_j(u) &= \frac{e^{-\mu_j} (\mu_j)^u}{u!} & j = 1, 2 \\ \mu_1^{(m+1)} &= \frac{\sum_{i=0}^9 y_i Z_i(\underline{\theta}^{(m)}) i}{\sum_{i=0}^9 y_i Z_i(\underline{\theta}^{(m)})} & A(u) &= \frac{g_1(u)}{g_2(u)} & u = 0, 1, \dots, 9 \\ \mu_2^{(m+1)} &= \frac{\sum_{i=0}^9 y_i (1 - Z_i(\underline{\theta}^{(m)})) i}{\sum_{i=0}^9 y_i (1 - Z_i(\underline{\theta}^{(m)}))}\end{aligned}$$

لذا:

با استفاده از رایانه جوابهای عددی بدست آمده عبارتند از:

$$\hat{\alpha} = 0/3599, \hat{\mu}_1 = 1/2561, \hat{\mu}_2 = 2/6624$$

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 موقع می‌توان از الگوریتم GEM (الگوریتم EM تعمیم یافته) که آن نیز برای نخستین بار توسط Dempster و سایرین در سال ۱۹۷۷ (در [۱]) ارایه 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 مرحلهٔ 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 درسی تی‌اسکن، مسایل ژنتیک و نیز مسایلی که در آن داده گمشده وجود دارد، اشاره نمود.

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

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

واما:

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

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  $(i)$  و  $(j)$  در امید شرطی بالا،

## مراجع

- [1] Dempster AP , Laird NM, Rubin DB. Maximum likelihood from incomplete data via the FM 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 Tomograohy. *J. of Computer Assisted Tomograohy*, Vol. 8, No. 2, 1984.
  - [4] K. Lange. Mathematical and Statistical Methods for Genetic Analysis. Springer, 1997, p 22-32.
-