

---

*EM* التَّحْرِيقُ

---

الگوریتم EM که توسط دمپستر، لیرد و رابین<sup>۴</sup> در سال ۱۹۷۷ و توسط مک لچلان و کریشنان<sup>۵</sup> در سال ۱۹۹۷ ارائه شد، ابزاری قدرتمند برای برآورد حداکثر درستنمایی با داده‌های ناقص می‌باشد. در اینجا معنی کلمه « ناقص » حالتی کلی دارد و در موقعیتهای متفاوت می‌تواند به معانی گوناگونی مانند: داده‌های گم شده، مؤلفه‌های مجهول، مشاهدات سانسور شده، متغیرهای نهفته و نظایر آنها اشاره داشته باشد. شرح مختصری از الگوریتم EM در زیر ارائه می‌شود.

یکی از روش های برآورد پارامتر روش ماکسیمم  
درست‌نمایی است

اگر بردار مشاهدات غیر کامل باشد که فقدان قسمتی  
از داده ها ممکن

است از مکانیزم سانسور کردن یا از مقادیر گم‌شده و  
..... باشد

---

استفاده نمود EM میتوان از الگوریتم

فرض کنید  $\mathbf{Y}$  بردار داده‌های مشاهده شده و  $\mathbf{X}$  بردار داده‌های نامعلوم باشند. همچنین فرض کنید  $\theta$  پارامتر مجهول است که می‌خواهیم برآورد شود و  $l_c(\theta; \mathbf{Y}, \mathbf{X})$  لگاریتم تابع درست‌نمایی بر اساس کلیه داده‌ها می‌باشد، که به ازای تمام مقادیر ممکن  $\theta$  در فضای پارامتر  $\Omega$ ، تعریف می‌شود.

الگوریتم با یک مقدار اولیه  $\theta^{(0)} \in \Omega$  آغاز می‌شود و دو مرحله زیر را تا رسیدن به همگرایی، تکرار می‌کند:

مرحله E :

محاسبه  $l^{(j)}(\theta) = E_{\mathbf{X}|\mathbf{Y}, \theta^{(j-1)}} [l_c(\theta; \mathbf{Y}, \mathbf{X})]$  که امید

ریاضی با توجه به داده‌های گم شده  $\mathbf{X}$  به شرط داده‌های مشاهده شده  $\mathbf{Y}$  گرفته می‌شود و باید توجه شود که مقدار  $\theta^{(j-1)}$  در این امید ریاضی جایگذاری می‌شود.

مرحله M : یافتن  $\theta^{(j)} \in \Omega$  بقسمی که  $l^{(j)}(\theta)$  بیشینه شود.

تکرار دو مرحله فوق به ازای  $j = 1, 2, \dots$  ، منجر به همگرایی دنباله  $\theta^{(1)}, \theta^{(2)}, \dots$  در ماکزیمم موضعی لگاریتم درستنمایی کلیه داده‌ها می‌شود. لازم به ذکر است که تحت شرایط خاصی این دنباله به همگرایی نمی‌رسد، برای جزئیات بیشتر به [۴] رجوع شود.

از جمله کاربردهای رایج این الگوریتم، حل مسائلی با داده‌های گم شده، یافتن نمای توزیع پسین در چهارچوب بیز [۳]، تشخیص مدل‌های آمیخته و کاربردهایی در داده‌های گروه‌بندی شده، سانسور شده و یا بریده شده، می‌باشند

مثال ۱) فرض کنید طول عمر لامپ های مربوط به یک کارخانه دارای توزیع نمایی با میانگین مجهول  $\theta$  باشد. از این کارخانه تعداد  $M+N$  لامپ را به تصادف انتخاب کرده و در دو آزمایش مستقل شرکت داده می شود. آزمایشگر تعداد  $N$  لامپ را در آزمایش اول قرار داده و طول عمر تک تک آنها را به صورت  $y_1, \dots, y_N$  ثبت می کند، در حالیکه در آزمایش دوم زمان  $t > 0$  را در نظر گرفته و همه  $M$  لامپ را داخل آزمایش می کند و فقط لامپ هایی را که تا

زمان  $t$  هنوز نسوخته‌اند مشخص می‌کند. بنابراین در آزمایش دوم طول عمر دقیق تک تک مشاهدات مشخص نیست و تنها مشخصه های  $E_1, E_2, \dots, E_M$  در دسترس هستند که در آن

$$E_i = \begin{cases} 1 & \text{اگر لامپ تا زمان } t \text{ سوخته نباشد} \\ 0 & \text{اگر لامپ قبل از زمان } t \text{ سوخته باشد} \end{cases}$$

$$, \quad i = 1, \dots, M$$

حال سؤال اینجاست که با این داده ها  $MLE(\theta)$  چقدر می‌شود؟



فرض کنید که  $X_1, \dots, X_M$  طول عمرهای (مشاهده نشده) لامپ های شرکت کننده در آزمایش دوم باشند و  $Z = \sum_{i=1}^M E_i$  تعداد لامپ هایی باشد که در آزمایش دوم تا زمان  $t$  هنوز نسوخته اند. بنابراین ترکیب داده های مشاهده شده در دو

آزمایش به صورت زیر است:  $\mathbf{Y} = (Y_1, \dots, Y_N, E_1, \dots, E_M)$

و داده های مشاهده نشده عبارتند از:  $\mathbf{X} = (X_1, \dots, X_M)$

همچنین تابع درستنمایی بر اساس کلیه داده ها به صورت

$$L(\theta; \mathbf{Y}, \mathbf{X}) = \frac{1}{\theta^{M+N}} \exp\left(\frac{-(\sum_{i=1}^N Y_i + \sum_{i=1}^M X_i)}{\theta}\right)$$

بنابراین خواهیم داشت:

$$l_c(\theta; \mathbf{Y}, \mathbf{X}) = \ln L(\theta)$$

$$= -(M + N) \log \theta - \frac{1}{\theta} \left( \sum_{i=1}^N Y_i + \sum_{i=1}^M X_i \right) \quad (1)$$

حال  $E(X_i | \mathbf{Y})$  را در دو حالت  $E_i = 0$  و  $E_i = 1$

محاسبه می‌کنیم. با توجه به این موضوع که پیشامد  $E_i = 1$  معادل

$X_i > t$  می‌باشد، می‌توان نوشت:

$$f(x_i | x_i > t) = \frac{f(x_i)}{1 - F(t)} = \frac{1}{\theta} e^{-\frac{t-x_i}{\theta}}$$

$$E(X_i | E_i = 1) = E(X_i | X_i > t)$$

$$= \int_t^{\infty} x_i f(x_i | x_i > t) dx_i = t + \theta$$

بطور مشابه در حالت  $E_i = 0$  خواهیم داشت:

$$E(X_i | E_i = 0) = \int_0^t x_i f(x_i | x_i \leq t) dx_i = \theta - \frac{te^{t/\theta}}{1 - e^{t/\theta}}$$

در نتیجه می توان نوشت:

$$E(X_i | \mathbf{Y}) = E(X_i | E_i) = \begin{cases} t + \theta & E_i = 1 \\ \theta - \frac{te^{t/\theta}}{1 - e^{t/\theta}} & E_i = 0 \end{cases} \quad (2)$$

حال می توان گفت که  $Z$  آمین مرحله از الگوریتم شامل

جایگذاری  $E(X_i | E_i)$  از رابطه (2)، بجای  $X_i$  در رابطه (1)

می باشد. لذا خواهیم داشت:

$$l^{(j)}(\theta) = -(N + M) \log \theta - \frac{1}{\theta} [N\bar{Y} + Z(t + \theta^{(j-1)}) + (M - Z)(\theta^{(j-1)} - t\theta^{(j-1)})] \quad (3)$$

که در آن  $p^{(j)} = \frac{e^{-t/\theta^{(j)}}}{1 - e^{-t/\theta^{(j)}}}$  و منظور از  $\theta^{(j)}$ ،  $\theta$ ی بدست

آمده در  $j$ امین تکرار می‌باشد. در  $j$ امین تکرار مرحله  $M$ ،  $\theta^{(j)}$  را چنان برآورد می‌کنیم که رابطه (۳) را بیشینه می‌سازد، به عبارت

دیگر با فرض ثابت بودن  $\theta^{(j+1)}$ ، برآورد بیشترین درست‌نمایی  $\theta^{(j)}$  را بدست می‌آوریم:

$$\frac{\partial l^{(j)}(\theta)}{\partial \theta} = 0 \quad \Rightarrow \quad (4)$$

$$\theta^{(j)} = f(\theta^{(j-1)})$$

$$\equiv \frac{N\bar{Y} + Z(t + \theta^{(j-1)}) + (M - Z)(\theta^{(j-1)} - tp^{(j-1)})}{N + M}$$

حال می‌توان با انتخاب یک  $\theta^{(0)} > 0$  معادله (4) را تا رسیدن به همگرایی تکرار کرد. خاصیت خودسازگاری الگوریتم وقتی آشکار می‌شود که  $Z = M$  باشد (یعنی کلیه لامپ‌هایی که در

آزمایش دوّم شرکت داده‌ایم تا زمان  $t$  سالم باشند)، در این حالت به  
جواب معروف زیر می‌رسیم:

$$\hat{\theta} = \frac{N\bar{Y} + Mt}{N}$$

### مراجع

- [1] Dempster, A. P., Laird, N. M., and Rubin, D. B., 1977, Maximum Likelihood from Incomplete Data via the EM Algorithm, (with discussion), *Journal of the Royal Statistical Society, Ser. B.* 39, 1-38.
- [2] McLachlan, G., and Krishnan, T., 1997, *The EM-algorithm and Extensions*, New York: Wiley.
- [3] Tanner, M. A., 1996, *Tools for Statistical Inference*, (3rd ed.), New York: Springer-Verlag.
- [4] Wu, C. F. J., 1983, On the Convergence Properties of the EM Algorithm, *The Annals of Statistics*, 11, 95-103.

