الگوریتم EM برای داده های از دست رفته (EXPECTATION-MAXIMIZATION)

استاد راهنما: دكتر وحيد رضايي تبار

ارائه دهندگان: حسام افشار – مهراب کیانی

فهرست مطالب

- مقدمه
- انواع داده های از دست رفته
- الگوريتم EM و نحوه عملكرد آن
 - مثال

مقدمه

بر اساس یک نظرسنجی علم داده که در سال 2020 توسط Anaconda صورت گرفت، متخصصان داده 45 درصد از زمان خود را صرف آماده سازی داده ها می کنند. یکی از مشکلات رایج که همه دانشمندان داده در این 45 درصد از زمان در دسترس خود با آن مقابله می کنند، رسیدگی به داده های از دست رفته در مجموعه داده خود است. و با این حال، هیچ راه حل یکسانی برای رسیدگی به این مشکل وجود ندارد.

دادههای از دست رفته یا مقادیر از دست رفته زمانی رخ میدهد که دادهای برای متغیرها یا شرکتکنندگان خاصی ذخیره نکنید. داده ها ممکن است به دلیل ورود ناقص داده ها، خرابی تجهیزات، فایل های از دست رفته و بسیاری دلایل دیگر از بین بروند و معمولا در مجموعه داده های بزرگ داده های از دست رفته وجود دارند.

انواع داده های از دست رفته

دلیل وجود داده های از دست رفته باید در نظر گرفته شود، زیرا به شما کمک می کند تا نوع داده های از دست رفته و اقداماتی که باید در مورد آن انجام دهید را تعیین کنید. در حالت کلی سه نوع اصلی از داده های از دست رفته وجود دارد که در جدول زیر آمده است:

مفهوم	انواع داده از دست رفته
داده های از دست رفته به طور تصادفی در بین متغیرها توزیع می شوند و با سایر متغیرها ارتباطی ندارند و هیچ ساختار خاصی برای توصیف این کمبود وجود ندارد.	به طور کاملاً تصادفی (MCAR)
داده های از دست رفته به طور تصادفی توزیع نمی شوند، اما آنها توسط سایر متغیر های مشاهده شده محاسبه می شوند.	به طور تصادفی (MAR)
داده های از دست رفته به طور سیستماتیک با مقادیر مشاهده شده متفاوت است.	به طور غیر تصادفی (MNAR)

انواع داده های از دست رفته

حذف، انتساب میانه یا میانگین، انتساب چندگانه یا حداکثر درستنمایی برخی از رایج ترین روش ها برای رسیدگی به داده های از دست رفته هستند اما آنها فقط برای مجموعههای دادهای قابل اجرا هستند که MCAR هستند و کاملاً MNAR نیستند.

در موارد خاص، ما بیشتر علاقه مند خواهیم بود که بدانیم آیا الگو یا دلیلی برای داده های از دست رفته وجود دارد یا خیر و با مراقبت از آن، آن را نسبت دهیم.

در چنین مواردی، روش هایی مانند الگوریتم EM مناسب تر است.

با این فرض که توزیع مشترک داده های گمشده و داده های مشاهده شده معلوم است، هدف الگوریتم EM یافتن تخمین پارامتری است که لگاریتم احتمال داده های مشاهده شده را به حداکثر می رساند، یعنی تابع چگالی احتمال مشاهدات و از آنجا مقادیر از دست رفته را تخمین می زند.

این الگوریتم دارای دو مرحله است: مرحله E-step) expectation) و مرحله maximization) که به صورت تکرار شونده انجام می شوند. M-step) که به صورت تکرار شونده انجام می شوند.

مراحل اجرای این الگوریتم عبارتند از:

- 1) با توجه به مجموعه داده ناقص یک سری پارامتر اولیه در نظر میگیریم.
- 2) (E-step): با استفاده از داده های موجود، مشاهدات گمشده را برآورد میکنیم.
- 3) (M-step): با استفاده از داده های کامل تولید شده در مرحله 2، پارامترهای مدل را به حداکثر میرسانیم.
 - 4) مرحله 2 و 3 را تا زمان همگرایی تکرار میکنیم، یعنی برآورد پارامتر بین تکرارها تغییر زیادی نکند.

یکی از مزایای اصلی الگوریتم EM این است که بایاس پارامترهای برآورد شده بسیار کمتر است.

فرض کنید که مشاهدات x_1, x_2, \dots, x_n را با x_1, x_2, \dots, x_n را با x_1, x_2, \dots, x_n را با x_1, x_2, \dots, x_n و همهٔ پار امتر های توزیع را نیز با x_1, x_2, \dots, x_n داده ها (پنهان و نمایان = مشاهدات) بر ابر خواهد بود با:

 $I(\theta) = \log p(d|\theta) = \log \sum_{h} p(d,h|\theta)$

از آنجا که لگاریتم تابع اکیداً صعودی است، میتوان لگاریتم درست نمایی کل داده ها را نسبت به θ بیشینه کرد. ولی آرگومان لگاریتم یک مجموع است و نمیتوان به سادگی پاسخ تحلیلی برای θ افت. از این رو، الگوریتم با ترفندی را برای بیشینه کردن حد پایین لگاریتم درست نمایی بکار میبرد. این حد پایین از نابرابری پنسن بدست میآید.

بر اساس نابرابری ینسن برای هر دسته k تایی از t_i ها و w_i ها اگر $w_i = 1$ برای هر دسته t_i خواهیم داشت:

 $\sum_{i=1}^{k} w_i \log t_i \le \log \sum_{i=1}^{k} w_i t_i$

اکنون (θ) را به صورت زیر باز نویسی میکنیم:

$$\log \sum_{h} q(h) \frac{p(d,h|\theta)}{q(h)} \ge \sum_{h} q(h) \log \frac{p(d,h|\theta)}{q(h)} = J(q,\theta)$$

که با گزینش $q(h)=p(h|d,\theta)$ داریم:

$$\begin{split} \log p(d|\theta) &= \sum_{h} p(h|d,\theta^{(t)}) \log p(d,h|\theta) - \sum_{h} p(h|d,\theta^{(t)}) \log p(h|d,\theta) \\ &= Q(\theta|\theta^{(t)}) + H(\theta|\theta^{(t)}) \\ \log p(d|\theta^{(t)}) &= \sum_{h} p(h|d,\theta^{(t)}) \log p(d,h|\theta^{(t)}) - \sum_{h} p(h|d,\theta^{(t)}) \log p(h|d,\theta^{(t)}) \\ &= Q(\theta^{(t)}|\theta^{(t)}) + H(\theta^{(t)}|\theta^{(t)}) \\ \log p(d|\theta) - \log p(d|\theta^{(t)}) &= Q(\theta|\theta^{(t)}) - Q(\theta^{(t)}|\theta^{(t)}) + H(\theta|\theta^{(t)}) - H(\theta^{(t)}|\theta^{(t)}) \end{split}$$

که با توجه به نابرابری گیبز
$$H(\theta|\theta^{(t)}) \ge H(\theta^{(t)}|\theta^{(t)})$$
 داریم $H(\theta^{(t)}|\theta^{(t)}) \ge H(\theta^{(t)}|\theta^{(t)})$ در نتیجه میتوان نوشت:

$$\log p(d|\theta) - \log p(d|\theta^{(t)}) \ge Q(\theta|\theta^{(t)}) - Q(\theta^{(t)}|\theta^{(t)})$$

در واقع یعنی انتخاب θ به گونه ای که $Q(\theta|\theta^{(t)})$ را بهبود دهد باعث میشود $\log p(d|\theta)$ نیز حداقل همان اندازه بهبود یابد

$$Q(\theta|\theta^{(t)}) = E_{h|d,\theta}[\log L(\theta|d,h)] = \sum_{h} p(h|d,\theta^{(t)}) \log p(d,h|\theta)$$
$$= \sum_{h} q(h) \log p(d,h|\theta)$$

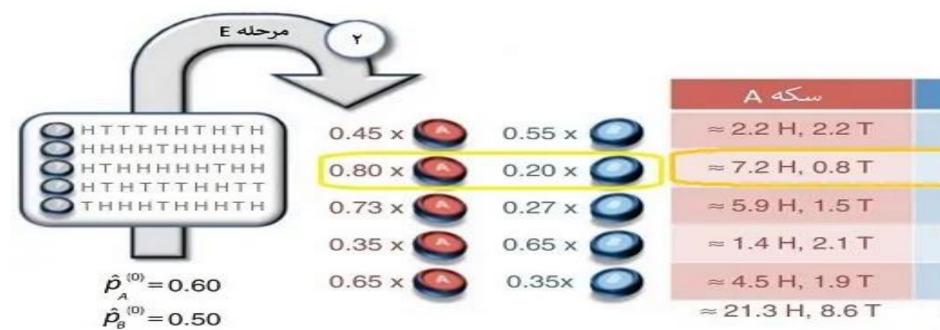
در نتیجه روش کار الگوریتم امید ریاضی-بیشینه کردن به صورت زیر است:

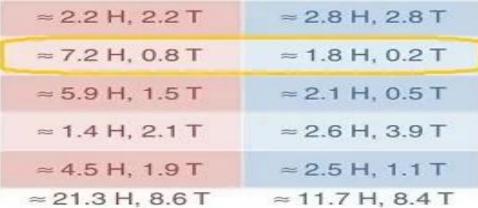
1- پارامترها را مقدار آغازین $\theta^{(0)}$ میدهیم.

2-تا رسیدن به همگرایی دو گام زیر را انجام میدهیم:

- $q^{(t)}$ = $armax Q(\theta|\theta^{(t)})$:(E-step) مید ریاضی گام امید ریاضی
- $\theta^{(t)}$ = $ar\max_{\theta} Q(\theta|\theta^{(t)})$:(M-step) گام بیشینه کردن

سکه B







$$\hat{p}_{A}^{(1)} \approx \frac{21.3}{21.3 + 8.6} \approx 0.71$$

$$\hat{p}_{\rm g}^{(1)} \approx \frac{11.7}{11.7 + 8.4} \approx 0.58$$

