



داده کاوی با استفاده از استنتاج و یادگیری بیزی

زیر نظر آقای دکتر رهگذر

توسط مصطفی حقیر چهرقانی

چکیده :

استدلال بیزی یک روش احتمالی برای استنتاج فراهم می‌آورد. این روش بر اساس این فرض بنا نهاده شده است که مقادیر مورد توجه از یک توزیع احتمال پیروی می‌کنند و اینکه تصمیم‌های بهینه می‌توانند با استدلال بر روی این احتمالات به همراه داده‌های مشاهده شده اتخاذ شوند. بدلیل اینکه این روش یک راه کار کمی برای وزن دهی شواهدی که از فرض‌های مختلف پشتیبانی می‌کند، فراهم آورد، در مبحث یادگیری ماشین از اهمیت فراوانی برخوردار است. استدلال بیزی، روشی مستقیم برای کار با احتمالات برای الگوریتم‌های یادگیری فراهم آورد، و همچنین چارچوبی برای تحلیل عملکرد الگوریتم‌هایی که مستقیماً با احتمالات سروکار ندارند ایجاد می‌نماید.

۱ - مقدمه

اهمیت استدلال بیزی داده کاوی را می‌توان به دو دلیل عمده نسبت داد. اول اینکه، الگوریتم‌های یادگیری بیزی که به طور صریح بر روی احتمالات فرض‌های مختلف کار می‌کنند، مانند naïve Bayes classifier که از جمله کاراترین و عملی‌ترین الگوریتم‌های ممکن برای برخی مسائل یادگیری می‌باشد. به عنوان مثال Michie (۱۹۹۴) مقایسه کاملی بین این الگوریتم و سایر الگوریتم‌ها مانند درخت تصمیم و شبکه عصبی انجام داده است. این محقق نشان می‌دهد که الگوریتم naïve Bayes classifier قابل رقابت با سایر الگوریتم‌ها و در برخی موارد بهتر از آنها عمل می‌کند.



دلیل دوم این است که روش‌های استدلال بیزی چشم انداز مفیدی برای درک عملکرد الگوریتم‌های که مستقیماً بر روی احتمالات عمل نمی‌کنند ایجاد می‌کند. برای نشان دادن این موضوع یک دلیل روشن برای یکی از مسائل یادگیری در شبکه عصبی که همان انتخاب کمترین مجموع مربعات خطا می‌باشد بیان می‌کنیم. همچنین تابع خطای دیگری به نام cross entropy ارائه می‌شود که از تابع خطای کمترین مجموع مربعات خطا در یادگیری توابع هدفی که بیانگر احتمال هستند کارتر می‌باشد. به عنوان مثالی دیگر، اصل Minimum Description Length را با استفاده از اصول تئوری اطلاعات توضیح داده می‌شود.

از ویژگی‌های یادگیری بیز می‌توان موارد زیر را نام برد:

- هر نمونه آموزشی جدید که مشاهده می‌شود می‌تواند احتمال درستی یک فرض را افزایش یا کاهش دهد. به این خاطر از روش‌های که با ناسازگاری یک نمونه فرض را کلاً حذف می‌کنند منعطف‌تر می‌باشد.

- دانش پیشین می‌تواند با مشاهدات ترکیب شده تا دانش جدید یا به عبارت دیگر احتمال درستی فرضیات را به وجود آورد. دانش پیشین به وسیله (۱) در نظر گرفتن احتمال هر فرض و (۲) انتساب یک توزیع احتمال برای مشاهدات، ساخته می‌شود.

- روش‌های بیزی می‌توانند از فرض‌هایی که احتمال را پیشبینی می‌کنند بهره‌گیرند (به عنوان مثال « این مریض به احتمال ۹۳٪ شانس بهبودی کامل را دارد »)

- نمونه‌های جدید می‌توانند با استفاده از ترکیب وزنی نمونه‌های قبل متناسب با احتمال آنها تولید شود.

- حتی در مواردی که روش‌های بیزی از لحاظ پیچیدگی محاسبات غیر قابل استفاده باشند می‌توان از آنها به عنوان بهترین روش (gold standard) برای مقایسه سایر روشها استفاده کرد.



یک مشکل عملی در استفاده از روش‌های بیزی این است که آنها عموماً نیاز به دانستن احتمالات پیشین بسیاری دارند. وقتی این احتمالات از قبل معلوم نباشد آنها بر اساس داده‌های موجود و دانش پیشین و توزیع احتمالی که بر روی فرض‌ها وجود دارد تخمین زده می‌شود. یک مشکل دیگر که در عمل به وجود می‌آید، هزینه محاسباتی زیاد هنگام محاسبه بهترین فرض بیزی در حالت عمومی است (هزینه محاسبات به طور خطی با تعداد فرض‌ها افزایش می‌یابد، یا به طور نمایی با تعداد متغیرها). در موارد خاص این هزینه محاسبات می‌تواند به شدت کاهش یابد.

۲- تئوری بیز

در بسیاری از موارد به دنبال پیدا کردن بهترین فرض در فضای مفروضات H ، با در اختیار داشتن داده‌های آموزشی D هستیم. یک روش برای بیان بهترین فرض این است که بگوییم ما به دنبال محتمل‌ترین فرض، با داشتن داده D به علاوه دانش اولیه در مورد احتمالات پیشین فرض‌های H ، هستیم. قضیه بیز روش مستقیم برای محاسبه این احتمالات فراهم می‌آورد.

برای تعریف قضیه بیز ابتدا کمی نماد گذاری لازم است. ما از $P(h)$ برای بیان احتمال اولیه‌ای که فرض h درست است استفاده می‌کنیم، پیش از آنکه داده‌های آموزشی را دیده باشیم. $P(h)$ را عموماً احتمال پیشین می‌نامند و بیانگر هر دانش پیشینی می‌باشد که در مورد شانس درستی فرض h سخن می‌گوید. اگر هیچ دانش اولیه‌ای از مفروضات نداشته باشیم می‌توانیم یک احتمال یکسان به کل فضای مفروضات H اختصاص دهیم. به طور مشابه از $P(D)$ برای بیان احتمال پیشین که داده‌های D مشاهده می‌شوند استفاده می‌کنیم (به عبارت دیگر احتمال مشاهده D به شرط اینکه هیچ دانشی در مورد درستی مفروضات موجود نباشد). همچنین از $P(D|h)$ برای بیان احتمال D در دنیایی که فرض h صادق است استفاده می‌کنیم. در یادگیری ماشین ما به دنبال $P(h|D)$



هستیم ، یعنی احتمال درستی فرض h به شرط مشاهده داده‌های آموزشی D . $P(h|D)$ احتمال پسین h نام دارد ، بدین علت که بیانگر اطمینان ما از فرض h پس از مشاهده داده‌های D می‌باشد

قضیه بیز اصلی‌ترین سنگ بنای یادگیری بیزی می‌باشد، زیرا روشی برای محاسبه احتمال پسین $P(h|D)$ را از احتمال پیشین $P(h)$ به همراه $P(D)$ و $P(D|h)$ فراهم می‌آورد .

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)} \quad (۲,۱)$$

همانطور که انتظار می‌رود ، می‌توان مشاهده کرد که $P(h|D)$ با افزایش $P(h)$ و همچنین $P(D|h)$ افزایش می‌یابد . به همین ترتیب منطقی به نظر می‌رسد که $P(D|h)$ با افزایش $P(D)$ کاهش یابد ، زیرا با افزایش احتمال وقوع $P(D)$ که مستقل از h می‌باشد ، شواهد کمتری در D برای پشتیبانی از h وجود خواهد داشت .

در بسیاری از سناریوهای یادگیری، یادگیرنده مجموعه‌ای از فرض‌ها H را در نظر می‌گیرد و علاقمند به یافتن فرضی $h \in H$ می‌باشد که محتمل‌ترین باشد (یا حداقل یکی از محتمل‌ترین مفروضات اگر چندین تا وجود دارد) . هر فرضی که دارای این خصوصیت باشد به فرض (MAP) Maximum a posteriori نام دارد . ما می‌توانیم فرض MAP را با استفاده از قضیه بیز برای محاسبه احتمال پسین هر کاندیدا بیابیم . به عبارت دقیق‌تر h_{MAP} فرضی است که،

$$\begin{aligned} h_{MAP} &= \arg \max_{h \in H} P(h|D) \\ &= \arg \max_{h \in H} \frac{P(D|h)P(h)}{P(D)} \\ &= \arg \max_{h \in H} P(D|h)P(h) \end{aligned} \quad (۲,۱)$$



توجه کنید که در مرحله آخر بالا $P(D)$ را حذف کردیم زیرا محاسبه آن مستقل از h می باشد و همیشه یک عدد ثابت است. در برخی موارد، فرض می کنیم که تمام مفروضات $h \in H$ احتمال وقوع یکسانی دارند (یعنی $(\forall h_i, h_j \in H; P(h_i) = P(h_j))$). در این صورت یک ساده سازی دیگر نیز در فرمول (۲,۲) می توان انجام داد. به عبارت دیگر می توان فرضی را که $P(D|h)$ را ماکزیمم می کند در نظر گرفت (Maximum Likelihood).

$$h_{ML} = \arg \max_{h \in H} P(D | h) \quad (۲,۲)$$

در اینجا، از داده های D به عنوان نمونه های آموزشی برای یک تابع هدف و از H به عنوان مجموعه ای از توابع هدف ممکن یاد می کنیم. اما در حقیقت قضیه بیز کلی تر از این بحث است. یعنی می توان آن را به گونه ای مشابه برای مجموعه ای از هر نوع مفروضات دو به دو مستقل از هم که مجموع احتمالات آنها یک می شود استفاده کرد.

۳- درست نمایی ماکزیمم (Maximum Likelihood) و فرض کمترین مربعات

خطا

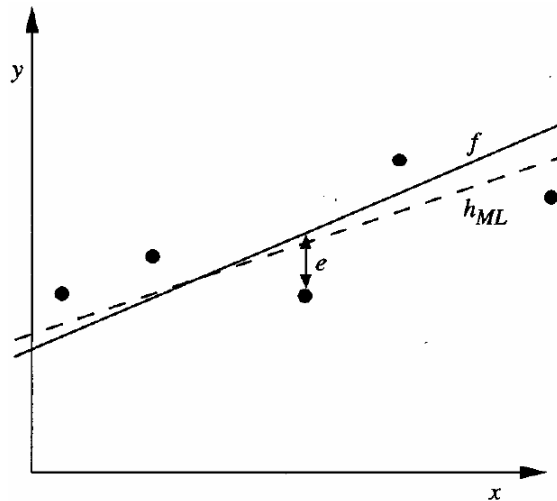
در این بخش، مسأله تابع هدف پیوسته را در نظر می گیریم، مسأله ای که در بسیاری از روش های یادگیری مانند شبکه عصبی، رگرسیون خطی و برازش منحنی های چند جمله ای وجود دارد. یک تحلیل بیزی ساده نشان می دهد که « در صورت برقرار بودن یک سری شرایط هر الگوریتم داده کاوی که مربعات خطای بین خروجی پیش بینی شده فرض و داده های آموزشی را کمینه کند، منجر به خروجی یک فرض درست نمایی ماکزیمم می شود ». اهمیت این نتیجه در آن است که صحت کارایی بسیاری از الگوریتم هایی که خطای مربعات را کمینه می کنند را نشان می دهد.



مسئله‌ای با شرایط زیر را در نظر گیرد. الگوریتم یادگیری L در فضای نمونه X و فضای مفروضات H که شامل دسته‌ای از توابع حقیقی بر روی X (یعنی هر $h \in H$ تابعی به شکل $h: X \rightarrow \mathbb{R}$ که مجموعه اعداد حقیقی می‌باشد) است.

مسئله‌ای که الگوریتم L باید به آن پردازد یافتن تابع هدفی مانند $f: X \rightarrow \mathbb{R}$ است که در فضای H می‌باشد. مجموعه‌ای به اندازه m از نمونه‌های آموزشی موجود می‌باشد که مقدار خروجی هر نمونه دارای مقداری نویز با توزیع نرمال می‌باشد. به عبارت دقیق‌تر هر نمونه آموزشی از یک جفت $\langle x_i, d_i \rangle$ تشکیل شده که در آن $d_i = f(x_i) + e_i$. در اینجا $f(x_i)$ مقدار بدون نویز تابع e_i و یک متغیر تصادفی است که نمایانگر نویز می‌باشد. فرض بر آن است که مقادیر e_i به صورت مستقل تولید شده که توزیع احتمال آن از یک توزیع نرمال با میانگین صفر پیروی می‌کند. وظیفه L بدست آوردن یک فرض درست‌نمایی ماکزیمم، یا به طور مشابه فرض MAP هنگامی که تمام مفروضات دارای احتمال یکسان هستند می‌باشد.

یک مثال ساده، یادگیری یک تابع خطی می‌باشد، اگر چه تحلیل که در اینجا بیان می‌شود برای هر تابعی می‌تواند استفاده شود. شکل (۱،۳) یک تابع خطی f را با خط ممتد به همراه داده‌های آموزشی دارای نویز نشان می‌دهد. خط نقطه چین شده بیانگر یک فرض h_{ML} می‌باشد که دارای کمترین مربعات خطا است. توجه داشته باشید که فرض ML همیشه برابر فرض اصلی، f ، نمی‌باشد، زیرا از تعداد محدودی نمونه آموزشی خطا دار منتج شده است.



شکل ۱،۳ یادگیری یک تابع حقیقی. خط ممتد متناظر با تابع هدف f نقاط متناظر با نمونه‌های آموزشی می‌باشد. خط نقطه‌چین شده متناظر با فرض درست‌نمایی ماکزیمم که بر اساس کمترین مجموع مربعات خطا آموزش دیده می‌باشد.

در اینجا به نشان دادن مسأله اصلی، یعنی اینکه فرض کمترین مربعات خطا در حقیقت فرض درست‌نمایی ماکزیمم می‌باشد می‌پردازیم. خاطر نشان می‌شود که به جای $P(D|h)$ از $p(D|h)$ که نشان دهنده تابع چگالی احتمال می‌باشد استفاده شده است زیرا تابع مورد نظر پیوسته است و نمی‌توان از تابع توزیع احتمال برای آن استفاده کرد،

$$h_{ML} = \arg \max_{h \in H} p(D|h)$$

همانند قبل مجموعه‌ای از نمونه‌های آموزشی $\langle x_1, \dots, x_m \rangle$ را در نظر می‌گیریم و از اینجا مجموعه متناظر مقادیر هدف $D = \langle d_1, d_2, \dots, d_m \rangle$ وجود خواهد داشت. همچنین $d_i = f(x_i) + e_i$ با فرض استقلال دو به دو نمونه‌ها به شرط h می‌توان نوشت،

$$h_{ML} = \arg \max_{h \in H} \prod_{i=1}^m p(d_i | H)$$

بفرض اینکه نویز e_i از یک توزیع نرمال با میانگین صفر و واریانس نامعلوم σ^2 پیروی می‌کند، هر d_i نیز باید از یک توزیع نرمال با واریانس σ^2 که دارای میانگین $f(x_i)$ می‌باشد پیروی کند. در نتیجه $p(d_i|h)$ را می‌توان با یک توزیع نرمال با واریانس σ^2 و میانگین $\mu = f(x_i)$ نشان داد.



به علت آنکه ما عبارت احتمال d_i را با فرض درست بودن h برای تعریف تابع هدف f می نویسیم ،
در این فرمول نیز قرار می دهیم $\mu = f(x_i) = h(x_i)$ در نتیجه،

$$\begin{aligned} h_{ML} &= \arg \max_{h \in H} \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(d_i - \mu)^2} \\ &= \arg \max_{h \in H} \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(d_i - h(x_i))^2} \end{aligned} \quad (3,1)$$

حال در اینجا یک تبدیل که در محاسبات درست‌نمایی ماکزیم معمول می‌باشد انجام می‌دهیم . به جای ماکزیم کردن عبارت پیچیده (3,1) از لگاریتم این عبارت که دارای پیچیدگی کمتری می‌باشد استفاده می‌کنیم . این عمل امکان پذیر است زیرا تابع لگاریتم یک تابع اکیداً یکنوا می‌باشد . بنابراین با ماکزیم کردن $\ln(p)$ به جای p داریم،

$$h_{ML} = \arg \max_{h \in H} \sum_{i=1}^m \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} (d_i - h(x_i))^2 \quad (3,2)$$

جمله اول در عبارت (3,2) که ثابتی مستقل از h می‌باشد و می‌توان آن را نادیده گرفت . در نتیجه،

$$h_{ML} = \arg \max_{h \in H} \sum_{i=1}^m -\frac{1}{2\sigma^2} (d_i - h(x_i))^2 \quad (3,3)$$

ماکزیم کردن عبارت (3,3) معادل مینیمم کردن منفی آن است،

$$h_{ML} = \arg \max_{h \in H} \sum_{i=1}^m \frac{1}{2\sigma^2} (d_i - h(x_i))^2$$

نهایتاً ، می‌توان ضرایب ثابتی که مستقل از h هستند را نادیده گرفت،

$$h_{ML} = \arg \max_{h \in H} \sum_{i=1}^m (d_i - h(x_i))^2 \quad (3,4)$$

بنابراین عبارت (3,4) نشان می‌دهد که فرض درست‌نمایی ماکزیم h_{ML} ، فرضی است که مجموع مربعات خطا بین داده‌های آموزشی مشاهده شده d_i و پیشبینی خروجی از فرض $h(x_i)$ را مینیمم



کند. این مسأله در صورتی صادق است که داده‌های آموزشی به وسیله اضافه کردن نویز نرمال تصادفی با میانگین صفر و به طور مستقل تولید شده باشند. مشاهده می‌شود که عبارت مربع خطا مستقیماً از نمای توزیع نرمال ناشی می‌شود. می‌توان عبارات مشابهی را با فرض‌های دیگری در مورد تابع توزیع خطا بدست آورد.

توجه کنید که در روش بدست آوردن عبارت (۳,۴) از ماکزیمم کردن لگاریتم درست‌نمایی ($\ln p(D|h)$) استفاده شده است. این روش منجر به رسیدن به جواب یکسان اما با عبارت ریاضی ساده‌تر می‌شود.

به چه علت فرض نرمال بودن تابع توزیع احتمال نویز معقول به نظر می‌رسد؟ یک دلیل این است که موجب تحلیل ریاضی ساده‌ای می‌شود. دلیل دوم این است که شکل توزیع نرم و زنگوله‌ای این تابع، تقریب مناسبی از نویز سیستم‌های واقعی می‌باشد. در حقیقت قضیه حد مرکزی بیانگر این است که مجموع تعداد کافی از متغیرهای تصادفی مستقل با توزیع یکسان (independent identical distributed (iid)) از یک توزیع نرمال پیروی می‌کند مستقل از اینکه توزیع تک تک متغیرهای تصادفی چه باشند. البته مسلماً ممکن است در عمل فاکتورهای مختلف تولید نویز، دارای توزیع یکسانی نباشند که در نتیجه قضیه صادق نخواهد بود.

خاطر نشان می‌شود که نویز در تحلیل فوق بر روی خروجی فرض تاثیر می‌گذاشت. اگر مایل به بررسی شرایطی باشیم که نویز بر روی خصوصیات (پارامترهای ورودی) نیز تاثیر داشت، تحلیل‌های بسیار پیچیده تری مورد نیاز بود.

۴ - اصل طول توصیف مینیمم



در این بخش ما یک دیدگاه بیزی در اصل Ockham's razor و رابطه آن با اصل طول توصیف مینیمم بیان می‌کنیم. اصل Ockham's razor یک اصل معروف برای بایاس استنتاج به گونه‌ای است که :

«خلاصه‌ترین توضیح را برای بیان مشاهدات انتخاب کن.»

اصل طول توصیف مینیمم با بررسی تعریف h_{MAP} با استفاده از مفاهیم تئوری اطلاعات بخوبی بیان می‌شود.

تعریف آشنای h_{MAP} را در نظر بگیرید،

$$h_{MAP} = \arg \max_{h \in H} P(D|h)P(h)$$

که می‌توان آنرا به صورت زیر نیز نوشت،

$$h_{MAP} = \arg \max_{h \in H} \log_2 P(D|h) + \log_2 P(h)$$

یا به عبارت دیگر،

$$h_{MAP} = \arg \min_{h \in H} -\log_2 P(D|h) - \log_2 P(h) \quad (۴,۱)$$

با توجه به فرمول (۴,۱) می‌تواند این‌گونه توجیح کرد که فرض‌های کوتاه‌تر با در نظر گرفتن یک ارائه خاص و روشی برای کدگذاری فرض و داده بر دیگر فرض‌ها ترجیح دارند. برای توضیح این مسأله باید ابتدا با مفاهیمی از تئوری اطلاعات آشنا شویم.

فرض کنید باید رمزی برای ارسال یک سری پیام که به طور تصادفی انتخاب می‌گردد طراحی کنیم. به عبارت دیگر احتمال انتخاب پیام i ، p_i می‌باشد، البته ما به فشرده‌ترین رمز علاقه‌مندیم؛ یعنی رمزی که امید تعداد بیت‌های ارسالی برای فرستادن پیام را مینیمم کند. به وضوح می‌توان دریافت که برای مینیمم کردن امید طول رمز باید تعداد بیت کمتری به پیام‌هایی که با احتمال بیشتر ظاهر می‌شوند اختصاص داد. Weaver و Shanon (۱۹۴۹) نشان دادند که کد بهینه



$-\log_2 p_i$ بیت به پیام i ، اختصاص می‌دهد. از این پس به تعداد بیت‌های لازم برای رمز گذاری پیام i به روش کد گذاری C طول توصیف پیام i در روش C می‌گوییم و آن را با $L_C(i)$ نشان می‌دهیم.

حال با تحلیل معادله (۴,۱) با استفاده از نتایج بدست آمده از تئوری اطلاعات داریم:

• $-\log_2 P(h)$ طول توصیف بهینه برای h در فضای فرض‌های H می‌باشد. به عبارت دیگر این مقدار، اندازه توصیف بهینه برای فرض h با استفاده از بهترین ارائه ممکن می‌باشد. در نماد گذاری ما،

$$L_{C_H} = -\log_2 P(h)$$

که در آن C_H که بهینه برای فضای فرض H می‌باشد.

• $-\log_2 P(D|h)$ طول توصیف داده‌های آموزشی D به شرط درستی فرض h ، در کد گذاری بهینه می‌باشد. در نماد گذاری ما،

$$L_{C_{D|H}} = -\log_2 P(D|h)$$

• بنابراین می‌توان معادله (۴,۱) را به گونه‌ای بازنویسی کرد که نشان دهد h_{MAP} فرضی را که مجموع طول توصیف فرض و طول توصیف داده به شرط داشتن فرض را مینیمم می‌کند بیان می‌نماید،

$$h_{MAP} = \arg \min_{h \in H} L_{C_H}(h) + L_{C_{D|H}}(D|h)$$

که در آن C_H و $C_{D|H}$ متناظراً کد گذاری‌های بهینه برای H و D به شرط H می‌باشند.

بنابراین اصل MDL روشی را برای متعادل کردن پیچیدگی فرض و تعداد خطاهایی که مرتکب می‌شود فراهم می‌نماید. یعنی ممکن است فرض کوتاه‌تری را که کمی خطا می‌کند به فرض بلندی که هیچ خطایی مرتکب نمی‌شود ترجیح دهد. به این ترتیب روشی را برای برخورد با

برازش بیش از حد مدل با داده ارائه می‌نماید.



Quinlan, Rivest (۱۹۸۹) آزمایش‌هایی برای تعیین اندازه درخت تصمیم‌های ساخته شده به روش MDL توصیف می‌نمایند. نتایج آنها حاکی از قابل مقایسه بودن دقت طبقه بندی درخت تصمیم‌های ساخته شده به روش MDL با درخت‌های ساخته شده با روش‌های استاندارد هرس کردن می‌باشد. این نتایج در کارهای Mehta (۱۹۹۵) نیز دیده می‌شود.

آیا از مباحث انجام شده در اصل MDL می‌توان نتیجه گرفت که همواره فرض کوتاه‌تری که توسط این اصل ترجیح داده می‌شود بهترین فرض مورد نظر است؟ خیر. تنها مطلبی که در اینجا نشان دادیم این است که اگر ارائه فرض به گونه‌ای باشد که اندازه آن بر اساس $-\log_2 P(h)$ و ارائه داده‌ها به گونه‌ای باشد که طول داده‌های کدگذاری شده D به شرط دانستن فرض h برابر $-\log_2 P(D|h)$ باشد آنگاه فرضی که بر اساس اصل MDL باشد فرض MAP را تولید می‌نماید. اما، برای داشتن چنین ارائه‌ای باید احتمالات پیشین $P(h)$ و همچنین درست‌نمایی $P(D|h)$ را در اختیار داشته باشیم. در غیر این صورت هیچ دلیلی وجود ندارد که فرض MDL را با هر کدگذاری دلخواه C_1 و C_2 بر سایر فرض‌ها ترجیح دهیم.

۵- طبقه بندی بهینه بیزی (Bayes Optimal Classifier)

تا به حال ما این سوال را مورد توجه قرار دادیم که «چه فرضی محتمل‌ترین فرض با داشتن داده‌های آموزشی D می‌باشد؟».

در حقیقت سوال بسیار مهم دیگری که خیلی نزدیک به این سوال است این است که «محتمل‌ترین دسته داده جدید با در اختیار داشتن داده‌های آموزشی D چیست؟». ممکن است به نظر رسد که استفاده از فرض MAP برای دسته بندی داده جدید کفایت می‌کند، اما از آن بهتر نیز می‌توان عمل کرد.



برای درک این موضوع فرض کنید فضای مفروضات ما دارای سه فرض h_1, h_2, h_3 می‌باشد. همچنین احتمال پسین این فرضیات با توجه به داده‌های آموزشی به ترتیب برابر $0/3$ و $0/4$ و $0/3$ می‌باشد. بنابراین h_1 فرض MAP است. فرض کنید نمونه جدید x مشاهده می‌شود که به وسیله فرض h_1 مثبت و توسط سایر فرض‌ها منفی دسته بندی می‌شود. با در نظر گرفتن کل مفروضات احتمال اینکه x مثبت باشد $0/4$ و احتمال اینکه منفی باشد $0/6$ می‌باشد. محتمل‌ترین دسته‌بندی در این مسأله با دسته‌ای که محتمل‌ترین فرض (فرض MAP) را تولید می‌کند متفاوت می‌باشد.

در حالت کلی محتمل‌ترین دسته بندی نمونه جدید به وسیله ترکیب احتمالات تمام مفروضات به وسیله وزن دهی آنها بر اساس احتمال پسین مفروضات انجام می‌گیرد. اگر دسته محتمل v_j از یک مجموعه مانند V انتخاب گردد آنگاه احتمال $P(v_j|D)$ که بیانگر احتمال درستی دسته نمونه جدید به شرط مشاهده D می‌باشد به صورت زیر محاسبه می‌گردد،

$$P(v_j | D) = \sum_{h_i \in H} P(v_j | h_i) P(h_i | D) \quad (5,1)$$

بهترین دسته بندی برای نمونه جدید دسته‌ای خواهد بود که عبارت $5,1$ را ماکزیمم کند، یعنی،

$$\arg \max_{v_j \in V} \sum_{h_i \in H} P(v_j | h_i) P(h_i | D) \quad (5,2)$$

برای اینکه فرمول $(5,2)$ را برای مثال قبل توضیح دهیم مجموعه $V = \{\oplus, \otimes\}$ را برای دسته بندی نمونه جدید در نظر گیرید،

$$\begin{aligned} P(h_1 | D) &= .4, & P(\otimes | h_1) &= 0, & P(\oplus | h_1) &= 1 \\ P(h_2 | D) &= .3, & P(\otimes | h_2) &= 1, & P(\oplus | h_2) &= 0 \\ P(h_3 | D) &= .3, & P(\otimes | h_3) &= 1, & P(\oplus | h_3) &= 0 \end{aligned}$$

بنابر این،



$$\sum_{h_i \in H} P(\oplus | h_i) P(h_i | D) = .4$$

$$\sum_{h_i \in H} P(\otimes | h_i) P(h_i | D) = .6$$

و در نتیجه،

$$\arg \max_{v_j \in \{\oplus, \otimes\}} \sum_{h_i \in H} P(v_j | h_i) P(h_i | D) = \otimes$$

هر سیستمی که نمونه‌های جدید را بر اساس فرمول (۵,۲) دسته‌بندی نماید یک روش طبقه‌بندی بهینه بیزی خواهد بود. هیچ روش دیگر وجود ندارد که با داشتن فضای مفروضات یکسان و دانش پیشین مفروضات یکسان بتواند بهتر از این روش عمل نماید. این روش با توجه به داده‌های موجود، فضای مفروضات و احتمالات پیشین بر روی مفروضات، احتمال درست دسته‌بندی کردن نمونه جدید را ماکزیم می‌کند.

یک نکته جالب در مورد طبقه‌بندی بهینه بیزی است این است که تخمین‌های آن می‌تواند مربوط به فرضی باشد که خارج از فضای مفروضات H است! در نظر بگیرید که با استفاده از فرمول (۵,۲) تمام نمونه‌ها در X را دسته‌بندی کردیم. بر چسب گذاری نمونه‌ها به این صورت نیازی ندارد معادل با برچسب گذاری یکی از فرض‌ها مانند h از فضای مفروضات H باشد. می‌توان اینگونه تصور کرد که طبقه‌بندی بهینه بیزی از فضای مفروضات H که متفاوت و ابرمجموعه مفروضات H است استفاده می‌کند. در حقیقت فضای H شامل هر ترکیب خطی از فضای H می‌باشد.

۶- الگوریتم Gibbs

اگر چه روش طبقه‌بندی بهینه بیزی بهترین کارایی را دارد، اما از نظر محاسباتی بسیار پیچیده است. این هزینه زیاد به علت محاسبه احتمال پسین تمامی مفروضات و وزن دهی بر اساس آن برای دسته‌بندی هر نمونه جدید می‌باشد.



یک روش که بهینگی و پیچیدگی کمتری دارد الگوریتم Gibbs می باشد (Oppper و Haussler، ۱۹۹۱) که به صورت زیر است :

۱- فرض h را از فضای H به صورت تصادفی، براساس توزیع احتمال پسین بر روی H انتخاب کن.

۲- از h برای دسته بندی نمونه جدید x استفاده کن .

هنگام دریافت نمونه جدید الگوریتم Gibbs فرض را به تصادف براساس توزیع احتمال پسین کنونی انتخاب می کند . می توان نشان داد که در چنین شرایطی امید دسته بندی کردن اشتباه حداکثر دو برابر خطای طبقه بندی بهینه بیزی می باشد (Haussler، ۱۹۹۴).

۷- طبقه بندی ساده بیزی

یک روش بسیار کاربردی یادگیری بیزی روش یادگیرنده ساده بیزی می باشد که عموماً روش طبقه بندی ساده بیزی نامیده می شود . در برخی زمینه ها نشان داده شده است که کارایی آن قابل قیاس با کارایی روش هایی مانند شبکه عصبی و درخت تصمیم می باشد . این بخش، روش طبقه بندی ساده بیزی را معرفی می کند . بخش بعدی کاربرد این روش را در یک مثال عملی نشان می دهد .

طبقه بندی ساده بیزی برای مسائلی که هر نمونه x در آن توسط مجموعه ای از مقادیر صفات و تابع هدف $f(x)$ از مجموعه ای مانند V انتخاب می گردد کاربرد دارد . مجموعه ای از داده های آموزشی و خروجی تابع هدف و یا طبقه ای که نمونه جدید به آن تعلق دارد مورد نظر است .

روش بیزی برای طبقه بندی نمونه جدید این است که محتمل ترین طبقه یا مقدار هدف v_{MAP} را با داشتن مقادیر صفات $\langle a_1, a_2, \dots, a_n \rangle$ که توصیف کننده نمونه جدید است شناسایی کند،



$$v_{MAP} = \arg \max_{v_j \in V} P(v_j | a_1, a_2, \dots, a_n) \quad (7,1)$$

با استفاده از قضیه بیز می توان عبارت (7,1) را به صورت زیر بازنویسی کرد،

$$\begin{aligned} v_{MAP} &= \arg \max_{v_j \in V} \frac{P(a_1, a_2, \dots, a_n | v_j) P(v_j)}{P(a_1, a_2, \dots, a_n)} \\ &= \arg \max_{v_j \in V} P(a_1, a_2, \dots, a_n | v_j) P(v_j) \end{aligned} \quad (7,2)$$

حال با استفاده از داده های آموزشی سعی می کنیم دو جمله معادله (7,2) را تخمین بزنیم . محاسبه از روی داده های آموزشی به این صورت که میزان تکرار v_j در داده ها چقدر است، آسان می باشد . اما محاسبه جملات مختلف $P(a_1, a_2, \dots, a_n | v_j)$ به این صورت قابل قبول نخواهد بود مگر اینکه حجم بسیار زیاد از داده های آموزشی در اختیار داشته باشیم . مشکل اینجاست که تعداد این جملات برابر تعداد نمونه های ممکن ضرب در تعداد مقادیر تابع هدف می باشد . بنابراین باید هر نمونه را چندین بار مشاهده کنیم تا تخمین مناسبی از آن بدست آید .

فرض روش طبقه بندی ساده بیزی بر اساس این ساده سازی است که مقادیر صفات با داشتن مقادیر تابع هدف از یکدیگر مستقل شرطی می باشند . به عبارت دیگر ، این فرض بیانگر این است که به شرط مشاهده خروجی تابع هدف احتمال مشاهده صفات a_1, a_2, \dots, a_n برابر ضرب احتمالات هر صفت به طور جداگانه می باشد . اگر این را جایگزین معادله (7,2) کنیم روش طبقه بندی ساده بیزی را نتیجه می دهد،

$$v_{NB} = \arg \max_{v_j \in V} P(v_j) \prod_i P(a_i | v_j) \quad (7,3)$$

که v_{NB} خروجی طبقه بندی ساده بیزی برای تابع هدف می باشد . توجه کنید که تعداد جملات $P(a_i | v_j)$ که در این روش باید محاسبه شوند برابر تعداد صفات ضرب در تعداد دسته های خروجی برای تابع هدف می باشد که این مقدار از تعداد جملات $P(a_1, a_2, \dots, a_n | v_j)$ بسیار کمتر است .



نتیجه اینکه یادگیری ساده بیزی سعی در تخمین مقادیر مختلف $P(v_j)$ و $P(a_i | v_j)$ با استفاده از میزان تکرار آنها در داده‌های آموزشی دارد. این مجموعه تخمین‌ها متناظر با فرض یاد گرفته شده است. سپس از این فرض برای طبقه بندی نمونه‌های جدید استفاده می‌شود که این کار با استفاده از فرمول (۷,۳) صورت می‌گیرد. هر گاه فرض مستقل شرطی بودن روش طبقه بندی ساده بیزی بر آورده شود طبقه ساده بیزی معادل طبقه MAP خواهد بود.

۷,۱- یک مثال در توضیح طبقه بندی ساده بیزی

برای درک بهتر روش طبقه بندی ساده بیزی مثالی در رابطه با یادگیری مفهومی به این صورت بیان می‌کنیم:

«طبقه بندی روزها بر اساس اینکه فردی تنیس بازی می‌کند یا خیر!»

جدول (۷,۱) نشان دهنده ۱۴ نمونه آموزشی از مفهوم تنیس بازی کردن می‌باشد که هر روز با صفات وضعیت، دما، رطوبت و باد توصیف می‌شود.

جدول (۷,۱) داده‌های آموزشی برای مفهوم تنیس بازی کردن

تینیس بازی کردن	باد	رطوبت	دما	وضعیت	روز
خیر	کند	زیاد	گرم	آفتابی	۱
خیر	تند	زیاد	گرم	آفتابی	۲
بله	کند	زیاد	گرم	ابری	۳
بله	کند	زیاد	متوسط	بارانی	۴
بله	کند	نرمال	خنک	بارانی	۵



۶	بارانی	خنک	نرمال	تند	خیر
۷	ابری	خنک	نرمال	تند	بله
۸	آفتابی	متوسط	زیاد	کند	خیر
۹	آفتابی	خنک	نرمال	کند	بله
۱۰	بارانی	متوسط	نرمال	کند	بله
۱۱	آفتابی	متوسط	نرمال	تند	بله
۱۲	ابری	متوسط	نرمال	تند	بله
۱۳	ابری	گرم	زیاد	کند	بله
۱۴	بارانی	متوسط	نرمال	کند	خیر

در اینجا ما می‌خواهیم با استفاده از طبقه بندی ساده بیزی و داده‌های آموزشی نمونه جدید زیر را طبقه بندی کنیم .

وضعیت=آفتابی، دما=خنک، رطوبت=زیاد، باد=تند <

وظیفه ما تعیین مقدار هدف (بله یا خیر) مفهوم تنیس بازی کردن برای این نمونه جدید می‌باشد . با استفاده از فرمول (۷,۳) می‌توان v_{NB} را اینگونه محاسبه کرد،

$$\begin{aligned}
 v_{NB} &= \arg \max_{v_j \in \{yes, no\}} P(v_j) \prod_i P(a_i | v_j) \\
 &= \arg \max_{v_j \in \{yes, no\}} P(v_j) P(outlook = sunny | v_j) P(temperature = cool | v_j) \\
 &\quad P(humidity = high | v_j) P(wind = strong | v_j)
 \end{aligned}
 \tag{۷,۱,۱}$$



برای محاسبه v_{NB} باید مقادیر ۱۰ احتمال تعیین شوند که با استفاده از داده‌های آموزشی تخمین زده می‌شوند ابتدا احتمال مقادیر هدف را به راحتی با استفاده از فرکانس تکرار آن در ۱۴ نمونه آموزشی تخمین می‌زنیم،

$$P(\text{playTennis} = \text{yes}) = \frac{9}{14} = .64$$

$$P(\text{playTennis} = \text{no}) = \frac{5}{14} = .36$$

به طور مشابه می‌توانیم احتمالات شرطی را تخمین بزنیم. برای مثال آنهایی که مربوط به «باد=تند» می‌شوند عبارتند از،

$$P(\text{wind} = \text{strong} | \text{playTennis} = \text{yes}) = \frac{3}{9} = .33$$

$$P(\text{wind} = \text{strong} | \text{playTennis} = \text{no}) = \frac{3}{5} = .60$$

با استفاده از این احتمالات و تخمین‌های مشابهی برای سایر صفات، از فرمول (۷,۱,۱) برای محاسبه v_{NB} استفاده می‌کنیم (در اینجا نام صفات برای سادگی حذف شده‌اند)

$$P(\text{yes})P(\text{sunny} | \text{yes})P(\text{cool} | \text{yes})P(\text{high} | \text{yes})P(\text{strong} | \text{yes}) = .0053$$

$$P(\text{no})P(\text{sunny} | \text{no})P(\text{cool} | \text{no})P(\text{high} | \text{no})P(\text{strong} | \text{no}) = .0206$$

بنابراین طبقه‌بندی ساده بیزی مقدار «خیر = تنیس بازی کردن» را بر اساس مقادیر یاد گرفته شده از داده‌های آموزشی به نمونه جدید اختصاص می‌دهد. علاوه بر آن با نرمالیزه کردن مقادیر بالا بگونه‌ای که مجموع آن یک شود می‌توان احتمال شرطی آنکه مقدار هدف «خیر» شود را محاسبه نمود.

$$\frac{.0206}{.0206 + .0053}$$

برای مثال در این نمونه احتمال مورد نظر برابر است با

۷,۲- تخمین زدن احتمالات

تا به حال، احتمالات را با استفاده از کسری از دفعات که مشاهده می‌شد نسبت به تعداد کل مشاهدات تخمین می‌زدیم. به عنوان مثال ما تخمین (خیر = تنیس بازی کردن | تند = باد) را با



کسر $\frac{n_c}{n}$ محاسبه کردیم که $n=5$ تعداد کل داده‌های آموزشی که «خیر=تنیس بازی کردن» است

و $n_c=3$ تعداد داده‌هایی از آنهاست که برای آنها «تند=باد» می‌باشد.

با وجود اینکه این کسر در بسیاری از موارد تخمین خوبی از احتمال مورد نظر است، هنگامی که n_c

بسیار کوچک باشد تخمین خوبی بدست نمی‌دهد. برای دیدن این مشکل فرض کنید مقدار واقعی

(خیر=تنیس بازی کردن) P برابر 0.8 باشد و در داده‌های آموزشی، فقط 5 نمونه وجود دارد که

«خیر = تنیس بازی کردن» باشد. محتمل‌ترین مقدار n_c برابر صفر خواهد بود. این اثر، دو مسأله

را به وجود می‌آورد. اول اینکه $\frac{n_c}{n}$ یک احتمال بایاس شده تولید می‌کند. دوم اینکه وقتی تخمین

این احتمال صفر باشد، این احتمال بر نمونه‌هایی که در آینده مقدار «تند=باد» داشته باشند در طبقه

بندی بیز احتمال غالب خواهد بود. این مشکل به علت ضرب این احتمال صفر در فرمول $(7,3)$ با

دیگر به وجود می‌آید. برای جلوگیری از این امر می‌توان از یک راه کار بیزی به نام

$m_estimate$ که به صورت زیر تعریف می‌گردد استفاده شود،

$$\frac{n_c + mp}{n + m} \quad (7,2,1)$$

در اینجا n و n_c مانند قبل می‌باشند و p تخمین اولیه از احتمالی است که باید محاسبه شود. m

ضریبی است که به «اندازه نمونه متناظر» معروف است و کار آن وزن دهی به میزان اعتقاد ما از

احتمال پیشینی که در نظر گرفتیم می‌باشد. یعنی اگر صفتی دارای k حالت ممکن می‌باشد $p = \frac{1}{k}$

برای مثال، در تخمین (خیر=تنیس بازی کردن|تند=باد) صفت باد دارای 2 مقدار ممکن می‌باشد

، بنابراین $p=0.5$. دقت کنید که اگر m در اینجا برابر صفر اختیار شود، فرمول $(7,2,1)$ به

تخمین معمولی تبدیل می‌گردد. اگر هر دوی m و n غیر صفر باشند آنگاه ضریب $\frac{n_c}{n}$ و p



متناسب با وزن m با یکدیگر ترکیب می‌شوند. دلیل نامگذاری m به «اندازه نمونه متناظر» این است که فرمول $(\gamma, \mu, 1)$ می‌تواند به این صورت تعبیر شود که n مشاهده واقعی را با m نمونه مجازی که با احتمال p توزیع شده‌اند ترکیب می‌کنیم.

۸- الگوریتم EM:

به طور کلی الگوریتم EM در مسائلی به کار می‌رود که می‌خواهیم مجموعه‌ای از پارامترهای θ را تخمین بزنیم که مبتنی بر یک توزیع احتمالی اند و تنها بخشی از داده‌هایی که توسط این توزیع احتمالی تولید شده‌اند را داریم نه کل آنها را. مجموعه داده‌های مشاهده شده را به صورت $X = \{x_1, x_2, \dots, x_m\}$ که مستقل از همند و مجموعه داده‌های نامرئی را به صورت $Z = \{z_1, z_2, \dots, z_m\}$ و کل داده‌ها را به صورت $Y = X \cup Z$ در نظر می‌گیریم. الگوریتم EM فرض با بیشینه نزدیکی (ML) h' می‌گردد که $P(Y|h')$ را بیشینه کند. این امید ریاضی روی توزیع احتمال Y که توسط پارامترهای معلوم θ اندازه‌گیری می‌شود. الگوریتم EM برای تخمین توزیع احتمالی روی Y به جای پارامترهای واقعی θ از فرض فعلی h استفاده می‌کند. تابع $Q(h|h')$ را به گونه‌ای تعریف می‌کنیم که $E(\ln P(Y|h'))$ را به صورت تابعی از h' تعریف کند. با فرض $\theta = h$ و اینکه داده‌های X از کل داده‌های Y مشاهده شده‌اند، داریم:

$$Q(h|h') = E[\ln p(Y|h') | h, X] \quad (8,1)$$

مرحله ۱: مرحله تخمین (E). با استفاده از مقدار فرض فعلی h و داده‌های مشاهده شده X مقدار $Q(h|h')$ را محاسبه می‌کند تا توزیع احتمالی روی Y را تخمین زند.

مرحله ۲: مرحله بیشینه سازی (M). فرض h را با فرض h' ای جایگزین می‌کند که تابع Q را بیشینه کند.



$h \leftarrow \arg \max_{h'} Q(h' h)$	(۸,۲)
---	-------

۸,۱ تولید الگوریتم K-Mean از الگوریتم EM :

مسئله K-Mean عبارتست از تخمین پارامترهای $\theta = \langle \mu_1, \mu_2, \dots, \mu_k \rangle$. این پارامترها میانگین K توزیع نرمال را تعریف می کنند. مجموعه داده های مشاهده شده $X = \{x_i\}$ و داده های نامرئی $Z = \{z_{i,1}, z_{i,2}, z_{i,3}, \dots, z_{i,k}\}$ را در نظر بگیرید. داده های نامرئی مشخص می کنند که کدام یک از K توزیع نرمال برای تولید x_i به کار رفته اند. برای استفاده از الگوریتم EM برای عبارت $Q(h|h')$ باید مقداری مناسب با مسئله K-Mean بیابیم. احتمال یک نمونه $y_i = \langle x_i, z_{i,1}, z_{i,1}, \dots, z_{i,k} \rangle$ از کل داده ها می تواند به صورت زیر نوشته شود :

$p(y_i h') = p(x_i, z_{i,1}, z_{i,1}, \dots, z_{i,k} h') = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} \sum_{j=1}^k z_{ij} (x_i - \mu'_j)^2}$	(۸,۱,۱)
---	---------

تنها یکی از z_{ij} می تواند ۱ باشد و بقیه باید ۰ باشند. لذا این عبارت توزیع احتمال تولید x_i توسط توزیع انتخاب شده را نشان می دهد. لگاریتم احتمال $P(Y|h')$ برای همه m نمونه به صورت زیر است :

$P(Y h') = \ln \prod_{i=1}^m p(y_i h')$ $= \sum_{i=1}^m \ln p(y_i h')$ $= \sum_{i=1}^m \left(\ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} \sum_{j=1}^k z_{ij} (x_i - \mu'_j)^2 \right)$	(۸,۱,۲)
---	---------



اکنون باید امید ریاضی $P(Y|h')$ را روی توزیع احتمال Y یا به طور معادل روی قسمت‌های نامرئی

Y یعنی مقادیر z_{ij} حساب کنیم. عبارت بالا برای $P(Y|h')$ تابع خطی از مقادیر z_{ij} است. لذا

خواهیم داشت:

$$E[P(Y|h')] = E\left[\sum_{i=1}^m \left(\ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} \sum_{j=1}^k z_{ij} (x_i - \mu'_j)^2\right)\right] \quad (۸,۱,۳)$$

$$= \sum_{i=1}^m \left(\ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} \sum_{j=1}^k E[z_{ij}] (x_i - \mu'_j)^2\right)$$

لذا تابع $Q(h|h')$ به صورت زیر در خواهد آمد:

$$Q(h'|h) = \sum_{i=1}^m \left(\ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} \sum_{j=1}^k E[z_{ij}] (x_i - \mu'_j)^2\right) \quad (۸,۱,۴)$$

که $h' = \langle \mu'_1, \mu'_2, \dots, \mu'_k \rangle$ فرض جدید (بهبود یافته) است. $E[z_{ij}]$ بر اساس مقادیر فعلی h و مقادیر

مشاهده شده X محاسبه می شود.

$$E[z_{ij}] = \frac{e^{-\frac{1}{2\sigma^2}(x_i - \mu'_j)^2}}{\sum_{n=1}^k e^{-\frac{1}{2\sigma^2}(x_i - \mu'_n)^2}} \quad (۸,۱,۵)$$

مرحله دوم الگوریتم EM مقادیر $\langle \mu'_1, \mu'_2, \dots, \mu'_k \rangle$ را به گونه ای تعیین می کند که مقدار تابع Q

بیشینه شود:



$$\begin{aligned} \arg \max_{h'} Q(h' | h) &= \arg \max_{h'} \sum_{i=1}^m \left(\ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} \sum_{j=1}^k E[z_{ij}] (x_i - \mu'_j)^2 \right) \\ &= \arg \min_{h'} \sum_{i=1}^m \sum_{j=1}^k E[z_{ij}] (x_i - \mu'_j)^2 \end{aligned} \quad (۸,۱,۶)$$

لذا فرض ML (بیشینه کردن نزدیکی) به صورت کمینه کردن مجموع مربعات خطاها در می آید .

مقدار فرمول وقتی کمینه می شود که هر μ_j به صورت زیر مقدار دهی شود :

$$\mu_j = \frac{\sum_{i=1}^m E[z_{ij}] x_i}{\sum_{i=1}^m E[z_{ij}]} \quad (۸,۱,۷)$$

مراجع

- [۱] Buntine W. L. (۱۹۹۴). Operations for learning with graphical models. *Journal of Artificial Intelligence Research*, ۲, ۱۵۹-۲۲۵.
- [۲] R. Agrawal, J. Gehrke, D. Gunopulos, and Prabhakar Raghavan. Automatic subspace clustering of high dimensional data for data mining applications. In *ACM SIGMOD Conference*, ۱۹۹۸.
- [۳] R. Neal and G. Hinton. A view of the EM algorithm that justifies incremental, sparse and other variants. Technical report, Dept. of Statistics, University of Toronto, ۱۹۹۳.
- [۴] PELLEG, D. and MOORE, A. ۲۰۰۰. X-means: Extending K-means with Efficient Estimation of the Number of Clusters. In *Proceedings 1۷th ICML*, Stanford University.
- [۵] LEE, C-Y. and ANTONSSON, E.K. ۲۰۰۰. Dynamic partitional clustering using evolution strategies. In *Proceedings of the ۳rd Asia-Pacific Conference on Simulated Evolution and Learning*, Nagoya, Japan.
- [۶] Casella, G., & Berger, R. L. (۱۹۹۰). *Statistical inference*. Pacific Grove, CA: Wadsworth & Brooks/Cole.



- [V] [L. Xu and M. Jordan. On convergence properties of the EM algorithm for Gaussian mixtures. *Neural Computation*, 7, 1995.](#)
- [A] [Cédric Archambeau, John A. Lee, Michel Verleysen. On Convergence Problems of the EM Algorithm for Finite Gaussian Mixtures. ESANN'2003 proceedings – European Symposium on Artificial Neural Networks Bruges \(Belgium\), 23-25 April 2003, d-side publi., ISBN 2-930307-03-X, pp. 99-106](#)
- [9] [P. Langley, W. Iba, and K. Thompson. An Analysis of Bayesian Classifiers. Proc. 10th Nat. Conf. on Artificial Intelligence, 223-228, AAAI Press and MIT Press, USA 1992](#)
- [10] [P. Langley and S. Sage. Induction of Selective Bayesian Classifiers. Proc. 10th Conf. on Artificial Intelligence, 1994](#)
- [11] [B.W. Silverman: Density Estimation for Statistics and Data Analysis. Chapman and Hall, London \(1986\).](#)
- [12] [E. Parzen: On Estimation of a Probability Density Function and Mode. *Annals of Math. Statistics*, 33, 1065-1076 \(1962\).](#)
- [13] [J. Bilmes: A Gentle Tutorial of the EM Algorithm and its Application to Parameter Estimation for Gaussian Mixture and Hidden Markov Models. Technical Report of the International Computer Science Institute, Berkeley, CA \(1998\).](#)
- [14] [Adwait Ratnaparkhi. 1998. *Maximum Entropy Models for Natural Language Ambiguity Resolution*. Ph.D. thesis, the University of Pennsylvania.](#)
- [15] [Robert E. Schapire and Yoram Singer. 2000. Boostexter: A boosting-based system for text categorization. *Machine Learning*, 39\(2/3\):135-168.](#)
- [16] [Andrew McCallum and Kamal Nigam. 1998. A comparison of event models for naive bayes text classification. In *AAAI-98 Workshop on Learning for Text Categorization*.](#)
- [17] [Thorsten Joachims. 1999. Transductive inference for text classification using support vector machines. In *Proc. 16th International Conf. on Machine Learning*, pages 200-209. Morgan Kaufmann, San Francisco, CA.](#)