



# نظریه یادگیری محاسباتی

امید اعتصامی، محمدهادی فروغمنداعرابی  
بهار ۱۳۹۳

## تیغ اوکام

جلسه‌های هفتم تا نهم

نگارنده: مریم غرقانی

## ۱ تیغ اوکام

تیغ اوکام<sup>۱</sup> یک اصل حل مسأله است که توسط ویلیام اوکام<sup>۲</sup>، فیلسوف و متخصص الهیات انگلیسی، بیان شده است. طبق این اصل، از بین فرضیات مختلفی که به یک اندازه پدیده‌های طبیعی را توجیه می‌کنند، بهترین فرضیه، کوتاهترین فرضیه است. در این بخش، از منظر یادگیری، به این اصل توجه می‌کنیم و یادگیری اوکام را معرفی می‌کنیم.

### ۱.۱ یادگیری اوکام

فرض کنید  $X = \bigcup_{n \geq 1} X_n$  فضای موردها،  $C = \bigcup_{n \geq 1} C_n$  کلاس مفهوم و  $\mathcal{H} = \bigcup_{n \geq 1} \mathcal{H}_n$  کلاس نمایش فرضیه است. در این بخش فرض می‌کنیم نمایش فرضیه‌ها به صورت دودویی، با الفبای  $\{0, 1\}$ ، است و بنابراین  $\text{size}(h)$  طول رشته دودویی

<sup>۱</sup>Occam's razor

<sup>۲</sup>William of Ockham

$h$  است.

فرض کنید  $c \in C_n$ ، مفهوم هدف باشد. نمونه برجسب‌دار  $S$  با اندازه  $m$ ، مجموعه‌ای از زوج‌ها به صورت زیر است:

$$S = \{ \langle x_1, c(x_1) \rangle, \dots, \langle x_m, c(x_m) \rangle \}$$

**تعریف ۱.** فرض کنید  $0 < \beta < 1$  و  $\alpha \geq 0$  دو عدد ثابت هستند. گوییم  $L$  یک الگوریتم  $(\alpha, \beta)$ -اوکام برای کلاس مفهوم  $C$  با استفاده از  $\mathcal{H}$  است، اگر  $L$  با ورودی نمونه برجسب‌دار  $S$  با اندازه  $m$ ، خروجی  $h \in \mathcal{H}$  را بدهد که دارای شرایط زیر است:

۱.  $h$  با نمونه  $S$  سازگار است.

$$\text{size}(h) \leq (n \cdot \text{size}(c))^{\alpha} m^{\beta} \quad ۲.$$

گوییم  $L$  یک الگوریتم  $(\alpha, \beta)$ -اوکام کارا است، اگر زمان اجرای  $L$  بر حسب  $m, n$  و  $\text{size}(c)$  باشد.

همان‌طور که از تعریف برمی‌آید، رشد  $\text{size}(h)$  نسبت به  $m$  کندتر از رشد خطی است ( $\beta < 1$ ). اگر اجازه می‌دادیم  $\beta \geq 1$ ، الگوریتم  $L$  می‌توانست یک حفظ‌کننده باشد؛ یعنی  $h$ ی را تولید می‌کرد که صرفاً به هر مثال، برجسب آن را نسبت می‌دهد.

**قضیه ۲. تیغ اوکام.** فرض کنید  $L$  یک الگوریتم  $(\alpha, \beta)$ -اوکام کارا برای کلاس مفهوم  $C$  با استفاده از  $\mathcal{H}$  باشد. همچنین فرض کنید  $D$ ، یک توزیع احتمال روی فضای مورد‌های  $x$  است و  $c \in C_n$ ، مفهوم هدف است و  $0 < \epsilon, \delta \leq 1$ . در این صورت یک ثابت  $a > 0$  وجود دارد به‌طوری‌که اگر  $L$ ، نمونه ورودی  $S$  شامل  $m$  مثال تصادفی تولید شده از او را کلاً  $Ex(c, D)$  را به عنوان ورودی دریافت کند، که  $m$  دارای خاصیت زیر باشد:

$$m \geq a \left( \frac{1}{\epsilon} \log \frac{1}{\delta} + \left( \frac{(n \cdot \text{size}(c))^{\alpha}}{\epsilon} \right)^{\frac{1}{1-\beta}} \right)$$

آنگاه با احتمال حداقل  $1 - \delta$ ، خروجی  $h$  از الگوریتم  $L$  دارای خطای حداکثر  $\epsilon$  است.

طبق قضیه، اگر  $\beta$  به ۱ نزدیک شود، تعداد مثال‌های لازم  $m$  به سمت بی‌نهایت میل می‌کند. این موضوع بیانگر آن است که اگر اندازه خروجی  $h$  به تعداد مثال‌ها نزدیک شود، الگوریتم  $L$  تقریباً یک حفظ‌کننده خواهد بود که در این صورت برای اینکه خطای خروجی‌اش را پایین بیاورد، باید نمونه‌ای با تعداد بسیار زیادی مثال دریافت کند.

**قضیه ۳. تیغ اوکام (نسخه اندازه‌ای).** فرض کنید  $C$  کلاس مفهوم و  $\mathcal{H}$  کلاس نمایش فرضیه باشد. فرض کنید الگوریتمی باشد که به ازای هر عدد طبیعی  $n$  و هر  $c \in C$ ، اگر  $L$  نمونه ورودی  $S$  را با  $m$  مثال برجسب‌دار از  $c$  دریافت کند، در زمان چندجمله‌ای بر حسب  $m, n$  و  $\text{size}(c)$  اجرا می‌شود و خروجی  $h \in \mathcal{H}_{n,m}$  را تولید می‌کند که با نمونه  $S$  سازگار است. آنگاه ثابت  $b > 0$  وجود دارد به‌طوری‌که به ازای هر  $n$ ، هر توزیع  $D$  روی  $X_n$  و هر مفهوم هدف  $c \in C_n$ ، اگر  $L$  یک نمونه  $S$  شامل  $m$  مثال تصادفی تولید شده از  $Ex(c, D)$  را به عنوان ورودی بگیرد، که در آن

$$\log |\mathcal{H}_{n,m}| \leq b \epsilon m - \log \frac{1}{\delta}$$

آنگاه  $L$  خروجی  $h \in \mathcal{H}$  را می‌دهد که با احتمال حداقل  $1 - \delta$ ، خطای  $h$  حداکثر  $\epsilon$  است.

**اثبات.** فرض کنید  $\text{error}(h) > \epsilon$ . در این صورت با آمدن هر مثال  $h$  با احتمال حداقل  $\epsilon$  رد می‌شود. پس احتمال رد نشدن  $h$  با یک مثال حداکثر  $1 - \epsilon$  است. بنابراین احتمال رد نشدن  $h$  با  $m$  نمونه مستقل حداکثر  $(1 - \epsilon)^m$  است. در نتیجه داریم:

$$\Pr\{\exists h \in \mathcal{H}_{n,m} : \text{error}(h) > \epsilon \text{ و } h \text{ پس از } m \text{ نمونه باقی بماند}\} \leq |\mathcal{H}_{n,m}| (1 - \epsilon)^m$$

قرار می‌دهیم:

$$|\mathcal{H}_{n,m}|(1-\epsilon)^m \leq \delta \quad (۱)$$

$$\Leftrightarrow |\mathcal{H}_{n,m}| \leq \frac{\delta}{(1-\epsilon)^m} \quad (۲)$$

$$\Leftrightarrow \log |\mathcal{H}_{n,m}| \leq m \log \frac{1}{1-\epsilon} - \log \frac{1}{\delta} \leq m\epsilon - \log \frac{1}{\delta} \quad (۳)$$

□

بنابراین احتمال اینکه خطای خروجی الگوریتم بیشتر از  $\epsilon$  باشد، حداکثر  $\delta$  است.

حال قضیه ۲ را ثابت می‌کنیم:

اثبات. فرض کنید  $L$  یک الگوریتم  $(\alpha, \beta)$ -اوکام است و  $\mathcal{H}_{n,m}$  فضای تمام خروجی‌های تولید شده الگوریتم  $L$  باشد در صورتی‌که ورودی  $L$  یک نمونه با اندازه  $m$  برچسب‌دار شده با یک مفهوم  $c \in \mathcal{C}$  باشد. اگر  $h \in \mathcal{H}_{n,m}$  خروجی  $L$  باشد، طبق تعریف،  $\text{size}(h) \leq (n \cdot \text{size}(c))^{\alpha m \beta}$  و در نتیجه داریم،  $|\mathcal{H}_{n,m}| \leq 2^{(n \cdot \text{size}(c))^{\alpha m \beta}}$ . کافی است ثابت کنیم اگر  $m \geq a \left( \frac{1}{\epsilon} \log \frac{1}{\delta} + \left( \frac{(n \cdot \text{size}(c))^{\alpha}}{\epsilon} \right)^{\frac{1}{1-\beta}} \right)$ ، آنگاه ثابت  $b > 0$  وجود دارد به‌طوری‌که  $\log |\mathcal{H}_{n,m}| \leq b\epsilon m - \log \frac{1}{\delta}$ . در این صورت با استفاده از نسخه اندازه‌ای قضیه تیغ اوکام، می‌توانیم نتیجه مطلوب را به‌دست آوریم. در واقع می‌خواهیم شرط زیر برقرار باشد:

$$m \geq \frac{1}{b\epsilon} \left( \log |\mathcal{H}_{n,m}| + \log \frac{1}{\delta} \right) \quad (۴)$$

برای رسیدن به این نامساوی، کافی است دو شرط زیر برقرار باشد:

$$m \geq \frac{1}{b\epsilon} \log |\mathcal{H}_{n,m}| \quad ۱.$$

$$m \geq \frac{1}{b\epsilon} \log \frac{1}{\delta} \quad ۲.$$

شرط ۱ معدل است با:

$$m \geq \frac{1}{b\epsilon} \left( \frac{1}{\epsilon} \log \frac{1}{\delta} + \left( \frac{(n \cdot \text{size}(c))^{\alpha}}{\epsilon} \right)^{\frac{1}{1-\beta}} \right) \quad (۵)$$

$$\Leftrightarrow m \geq \left( \frac{1}{b\epsilon} (n \cdot \text{size}(c))^{\alpha} \right)^{\frac{1}{1-\beta}} \quad (۶)$$

بنابراین کافی است شرط زیر برقرار باشد تا شرایط ۱ و ۲ ارضا شوند:

$$m \geq \log \frac{1}{\delta} + \left( \frac{1}{b\epsilon} (n \cdot \text{size}(c))^{\alpha} \right)^{\frac{1}{1-\beta}} \quad (۷)$$

□

کافی است قرار دهیم  $b = \frac{1}{a} > 0$  تا شرط بالا ارضا شود.

## ۲.۱ بهبود اندازه نمونه برای یادگیری عطف

در جلسات قبل یک الگوریتم برای یادگیری عطف<sup>۳</sup> ارائه شد که فرمول عطف  $h \in \mathcal{H}_{n,m}$  را تولید می‌کرد. اگر نمونه  $S$  حداقل یک مثال داشته باشد ( $m \geq 1$ )، به ازای هر متغیر  $x_i$ ،  $h$  یا شامل  $x_i$  است، یا  $\bar{x}_i$  و یا هیچ‌یک از  $x_i$  و  $\bar{x}_i$  در  $h$  وجود ندارد. اگر

<sup>۳</sup>conjunction

هیچ مثالی تولید نشده باشد ( $m = 0$ )، شامل همه  $2n$  لیترال ممکن است. بنابراین تعداد کل حالاتی که  $h$  می‌تواند داشته باشد، به صورت زیر است:

$$\mathcal{H}_{n,m} \leq 3^n + 1 \leq 4^n \Rightarrow \log \mathcal{H}_{n,m} \leq 2n \quad (8)$$

قرار می‌دهیم:

$$2n \leq b\epsilon m - \log \frac{1}{\delta} \quad (9)$$

$$\Rightarrow m \geq \frac{1}{b\epsilon} (2n + \log \frac{1}{\delta}) \quad (10)$$

$$\Rightarrow m = O\left(\frac{n}{\epsilon} + \frac{1}{\epsilon} \log \frac{1}{\delta}\right) \quad (11)$$

کرانی که برای اندازه نمونه از قضیه اوکام حاصل شد، از شرط  $m \geq \frac{\log \frac{1}{\delta}}{\epsilon}$  که در الگوریتم یادگیری عطف به‌دست آوردیم، بهتر است. الگوریتم قبلی، فرمول عطف با بیشترین طول را که با نمونه سازگار بود، به‌دست می‌آورد. در ادامه الگوریتمی را برای یادگیری عطف ارائه می‌دهیم که در آن سعی می‌کنیم یک فرمول عطف سازگار با نمونه با کوتاهترین طول را به‌دست آوریم.

### ۱.۲.۱ یادگیری عطف با تعداد کمی لیترال

در این قسمت یک الگوریتم کارتر برای یادگیری عطف ارائه می‌دهیم. الگوریتم جدید علاوه بر مثال‌های مثبت از مثال‌های منفی نیز استفاده می‌کند برای اینکه طول خروجی را کم کند. قبل از ارائه الگوریتم ابتدا یک مسأله ترکیبیاتی و یک راه‌حل تقریبی را برای آن معرفی می‌کنیم.

**تعریف ۴. مسأله پوشش مجموعه‌ای<sup>۴</sup>.** فرض کنید یک گردایه  $S$  از زیرمجموعه‌های  $U = \{1, \dots, m\}$  را به عنوان ورودی داریم؛ هدف پیدا کردن زیرگردایه  $T \subseteq S$  است که  $|T|$  کمینه باشد و  $U, T$  را بپوشاند:

$$\cup_{t \in T} t = U \quad (12)$$

مسأله پوشش مجموعه‌ای یک مسأله  $NP - hard$  معروف است. فرض کنید  $c = OPT(S)$ ، اندازه پوشش کمینه برای مورد  $S$  در این مسأله است. در ادامه، یک الگوریتم تقریبی حریصانه را معرفی می‌کنیم که یک پوشش  $R$ ، با اندازه حداکثر  $c \log m$  پیدا می‌کند:

الگوریتم به این صورت است که در ابتدا  $R$  را برابر با یک گردایه خالی قرار می‌دهد. در هر مرحله، مجموعه  $s^* \in S$  با بزرگترین اندازه را به  $R$  اضافه می‌کند؛ سپس  $S$  را به این صورت به‌روز می‌کند که به ازای هر  $s$  در  $S$ ،  $s - s^*$  را جایگزین می‌کند. این رویه را تا زمانی ادامه می‌دهد که  $R$ ، همه  $U$  را بپوشاند.

**قضیه ۵.** الگوریتم حریصانه یک پوشش مجموعه‌ای با اندازه حداکثر  $c \log m$  پیدا می‌کند.

**اثبات.** فرض کنید  $U_i \subseteq U$  مجموعه‌ای از عناصر است که بعد از  $i$  مرحله در الگوریتم، توسط  $R$  پوشیده نشده‌اند. یک زیرگردایه از  $S$  با اندازه حداکثر  $c$  وجود دارد که  $U_i$  را می‌پوشاند. (چون یک زیرگردایه با اندازه  $c$  وجود دارد که  $U$  را می‌پوشاند.) پس در مرحله  $i+1$ ام، یک مجموعه در  $S$  وجود دارد که اندازه آن حداقل  $\frac{|U_i|}{c}$  است. پس داریم:

$$|U_{i+1}| \leq |U_i| - \frac{|U_i|}{c} = |U_i| \left(1 - \frac{1}{c}\right) \quad (13)$$

<sup>۴</sup> the set cover problem

با استقرا روی  $i$  نتیجه می‌گیریم:

$$|U_i| \leq m \left(1 - \frac{1}{c}\right)^i \leq m e^{-\frac{i}{c}} = e^{\log m - \frac{i}{c}} \quad (14)$$

قرار می‌دهیم:

$$e^{\log m - \frac{i}{c}} < 1 \Rightarrow \log m - \frac{i}{c} < 0 \Rightarrow c \log m < i \quad (15)$$

بنابراین الگوریتم حداکثر  $c \log m$  مرحله ادامه می‌یابد؛ در نتیجه با حداکثر  $c \log m$  مجموعه از  $S$ ، کل  $U$  را می‌پوشاند.  $\square$

حال یک الگوریتم برای یادگیری عطف ارائه می‌کنیم و سپس با استفاده از قضیه تیغ اوکام، یک کران پایین برای اندازه نمونه به‌دست می‌آوریم.

الگوریتم جدید به این صورت است که وقتی یک نمونه  $S$  با  $m$  مثال را به عنوان ورودی دریافت می‌کند، الگوریتم یادگیری عطف قدیمی را روی  $S$  اجرا می‌کند تا خروجی  $h$  را تولید کند. سپس از مثال‌های منفی استفاده می‌کند تا لیتراهای اضافی را حذف کند. توجه کنید که حذف لیتراها، تاثیری روی سازگار بودن  $h$  با مثال‌های مثبت ندارد. فقط باید لیتراها را طوری حذف کنیم که فرضیه  $h$  همچنان با مثال‌های منفی سازگار باشد. برای این کار، این مسأله را به یک مسأله پوشش مجموعه‌ای تبدیل می‌کنیم و الگوریتم تقریبی حریصانه را روی آن اعمال می‌کنیم.

برای هر لیترال  $z$  در  $h$ ،  $N_z$  را مجموعه تمام مثال‌های منفی  $\langle a, \circ \rangle$  در  $S$  تعریف می‌کنیم که  $z$  در  $a$ ،  $\circ$  است. اگر گردایه‌ای از  $N_z$ ها داشته باشیم که کل مثال‌های منفی  $S$  را می‌پوشانند، اگر  $h'$  را عطف لیتراهای این گردایه تعریف کنیم، در این صورت  $h'$  با تمام مثال‌های منفی  $S$  سازگار است. می‌خواهیم کوچکترین گردایه از  $N_z$ ها را پیدا کنیم که تمام مثال‌های منفی  $S$  را بپوشانند. این یک مسأله پوشش مجموعه‌ای است و اگر از الگوریتم حریصانه برای حل آن استفاده کنیم، یک پوشش با اندازه حداکثر  $\text{size}(c) \log m \log n$  به‌دست می‌آوریم. بنابراین الگوریتم جدید، یک فرضیه با اندازه حداکثر  $\text{size}(c) \log m \log n$  تولید می‌کند. (فرضیه، حداکثر  $\text{size}(c) \log m$  لیترال دارد و اندازه هر لیترال  $\log n$  بیت است.) در اینجا،  $\mathcal{H}_{n,m}$ ، مجموعه تمام عطف‌های با حداکثر  $\text{size}(c) \log m$  است که برای آن داریم،  $|\mathcal{H}_{n,m}| \leq 2^{\text{size}(c) \log m \log n}$ . طبق نسخه اندازه‌ای قضیه تیغ اوکام، اگر شرط زیر برقرار باشد، الگوریتم ارائه شده، یک الگوریتم یادگیری با پارامترهای  $\epsilon$  و  $\delta$  است.

$$\frac{c(\log |\mathcal{H}_{n,m}| + \log \frac{1}{\delta})}{\epsilon} \leq m \quad (16)$$

در صورتی که شرط زیر برقرار باشد، این شرط قضیه تیغ اوکام ارضا می‌شود:

$$m \geq c_1 \left( \frac{1}{\epsilon} \log \frac{1}{\delta} + \frac{\text{size}(c) \log n (\log \text{size}(c) + \log \log n)}{\epsilon} \right) \quad (17)$$

با کمی تغییر در این الگوریتم، می‌توانیم کران بهتری برای اندازه نمونه به‌دست آوریم. به‌طور کلی همه الگوریتم‌هایی که تاکنون معرفی کردیم، فرضیه  $h$ ای را پیدا می‌کردند که با همه مثال‌های نمونه سازگار باشد. اگر این قید را برداریم و اجازه دهیم کمی خطا در نمونه داشته باشیم، ممکن است بتوانیم الگوریتم یادگیری بهتری داشته باشیم. همانند الگوریتمی که در این قسمت ارائه دادیم، عمل می‌کنیم؛ فقط الگوریتم تقریبی پوشش مجموعه‌ای را تا جایی ادامه می‌دهیم که به  $|S_i| < \frac{\epsilon}{4} m$  برسیم:

$$|S_i| < m e^{-\frac{i}{c}} < \frac{\epsilon}{4} m \quad (18)$$

$$\Rightarrow e^{-\frac{i}{c}} < \frac{\epsilon}{4} \quad (19)$$

$$\Rightarrow i > c \log \frac{4}{\epsilon} \quad (20)$$

بنابراین اگر  $h \in \mathcal{H}_{n,m}$  خروجی این الگوریتم باشد، تعداد لیترال‌های  $h$  حداکثر  $\text{size}(c) \log \frac{1}{\epsilon}$  است و چون هر لیترال را با  $\log n$  بیت نمایش می‌دهیم،  $\text{size}(h)$  حداکثر  $\text{size}(c) \log \frac{1}{\epsilon} \log n$  است. در نتیجه داریم:

$$|\mathcal{H}_{n,m}| \leq \text{size}(c) \log \frac{1}{\epsilon} \log n \quad (21)$$

**قضیه ۶.** اگر در الگوریتم یادگیری عطف از الگوریتم حریصانه تقریبی با خطای  $\frac{\epsilon}{4}$  استفاده کنیم، خروجی  $h$  تولید می‌شود که خطایش با زیاد شدن نمونه به صورت نمایی کم می‌شود.

قبل از اثبات این قضیه، ابتدا کران چرنوف را معرفی می‌کنیم:

**لم ۷. کران چرنوف.** فرض کنید  $X_1, \dots, X_m$ ، متغیرهای تصادفی برنولی مستقل با احتمال موفقیت  $p$  هستند. فرض کنید  $S = X_1 + \dots + X_m$  تعداد کل موفقیت‌ها در  $X_1, \dots, X_m$  است؛ بنابراین  $\mathbb{E}[S] = \sum_{i=1}^m \mathbb{E}[X_i] = pm$  است. در این صورت اگر  $0 \leq \gamma \leq 1$ ، آنگاه داریم:

$$\mathbb{P}[S > (1 + \gamma)pm] \leq e^{-\frac{m p \gamma^2}{4}} \quad (22)$$

$$\mathbb{P}[S < (1 - \gamma)pm] \leq e^{-\frac{m p \gamma^2}{4}} \quad (23)$$

حال به اثبات قضیه می‌پردازیم:

**اثبات.** فرض کنید  $\text{error}(h) > \epsilon$ . در نتیجه احتمال لو رفتن  $h$  با یک مثال حداقل  $\epsilon$  است. متغیر نشان‌گر  $X_i$  ( $1 \leq i \leq m$ ) را به این صورت تعریف می‌کنیم که اگر  $h$  با مثال  $i$  سازگار نباشد، مقدار آن ۱ و در غیر این صورت، مقدار آن ۰ است.

$$\mathbb{E}[X_i] = \Pr\{X_i = 1\} > \epsilon \Rightarrow \mathbb{E}\left[\sum_{i=1}^m X_i\right] > \epsilon m \quad (24)$$

با استفاده از کران چرنوف داریم:

$$\Pr\left\{\sum_{i=1}^m X_i < \frac{\epsilon}{4}m\right\} = \Pr\left\{\sum_{i=1}^m X_i < (1 - \frac{3}{4})\epsilon m\right\} \leq e^{-m \frac{\epsilon^2}{16}} \quad (25)$$

پس احتمال اینکه خطای یک  $h$  بد روی نمونه  $S$  کمتر از  $\frac{\epsilon}{4}m$  باشد، حداکثر  $e^{-m \frac{\epsilon^2}{16}}$  است. بنابراین احتمال اینکه یک  $h$  بد وجود داشته باشد که خطایش روی نمونه  $S$  کمتر از  $\frac{\epsilon}{4}m$  باشد، حداکثر  $|\mathcal{H}_{n,m}|e^{-m \frac{\epsilon^2}{16}}$  است. کافی است این مقدار را کمتر از  $\delta$  قرار دهیم تا مطمئن باشیم احتمال اینکه خطای خروجی الگوریتم بیشتر از  $\epsilon$  باشد، کمتر از  $\delta$  خواهد بود.

$$\delta > |\mathcal{H}_{n,m}|e^{-m \frac{\epsilon^2}{16}} \Rightarrow \log \delta > \log |\mathcal{H}_{n,m}| - m \frac{\epsilon^2}{16} \quad (26)$$

بنابراین کران زیر برای اندازه نمونه حاصل می‌شود: (به ازای یک ثابت  $c_1$ )

$$m > c_1 \left( \frac{1}{\epsilon} \log |\mathcal{H}_{n,m}| + \frac{1}{\epsilon} \log \frac{1}{\delta} \right) \quad (27)$$

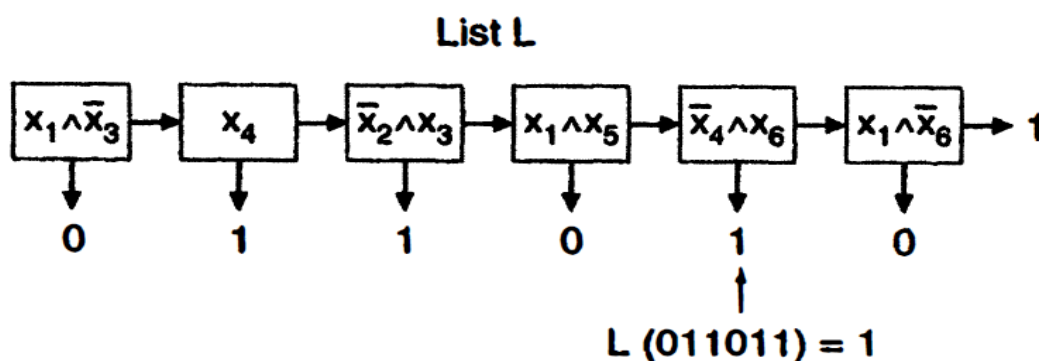
$$\Rightarrow m > c_1 \left( \frac{1}{\epsilon} \log \frac{1}{\delta} + \frac{\text{size}(c) \log \frac{1}{\epsilon} \log n}{\epsilon} \right) \quad (28)$$

□

### ۳.۱ یادگیری لیست تصمیم

تعریف ۸. یک  $k$ -لیست تصمیم<sup>۵</sup>  $(k - DL)$  روی متغیرهای بولی  $x_1, \dots, x_n$ ، یک دنباله مرتب  $L = (c_1, b_1), \dots, (c_l, b_l)$  با یک بیت  $b$  است؛ به این صورت که هر  $c_i$ ، عطف حداکثر  $k$  لیترال است. برای هر ورودی  $a \in \{0, 1\}^n$ ، اگر  $j$  کوچکترین اندیسی باشد که  $c_j(a) = 1$ ، آنگاه  $L(a) = b_j$ ؛ اگر چنین اندیسی وجود نداشته باشد، آنگاه  $L(a) = b$ .

شکل زیر نشان‌دهنده یک  $2 - DL$  به همراه مقداردهی آن روی یک ورودی خاص است:



شکل ۱: نمایش یک  $2-DL$

قضیه ۹. به ازای هر عدد ثابت  $1 \leq k$ ،  $k - DL$  قابل یادگیری کارای PAC است.

اثبات. یک الگوریتم یادگیری برای  $1 - DL$  ارائه می‌دهیم. برای  $k$  کلی، می‌توانیم به راحتی مسأله یادگیری  $k - DL$  را به این مسأله تحویل کنیم. این کار همانند روشی است که برای یادگیری  $k - CNF$  ارائه دادیم؛ در واقع کافی است برای هر عطف  $k$  متغیری، یک متغیر جدید تعریف کنیم و یک  $1 - DL$  با متغیرهای جدید را یاد بگیریم.

الگوریتم نمونه  $S$  را به عنوان ورودی می‌گیرد و ابتدا با یک لیست تصمیم خالی شروع می‌کند. در هر مرحله یک لیترال  $z$  را پیدا می‌کند که  $S_z \subseteq S$  (مجموعه تمام مثال‌هایی که  $z$  در آن‌ها ۱ است)، ناتهی باشد و فقط شامل مثال‌های مثبت یا فقط شامل مثال‌های منفی باشد. در این صورت به  $z$  یک لیترال مفید می‌گوییم. سپس الگوریتم جعبه  $z$  را به انتهای لیست تصمیم اضافه می‌کند. اگر  $S_z$  فقط شامل مثال‌های مثبت باشد، بیت مربوط به  $z$  را ۱ قرار می‌دهد؛ در غیر این صورت، بیت مربوط به  $z$  را ۰ قرار می‌دهد. سپس به جای  $S$ ،  $S - S_z$  را جاگذاری می‌کند و همین رویه را تکرار می‌کند تا به  $S = \emptyset$  برسد. در این صورت همه مثال‌ها با لیست تصمیمی که الگوریتم تولید می‌کند، سازگار هستند.

باید ثابت کنیم در هر مرحله که یک مجموعه  $S$  از مثال‌ها باقی مانده، حتماً یک لیترال مفید وجود دارد. باید توجه کنیم که مفهوم هدف، یک  $1 - DL$  است که نمونه از آن تولید شده است. در  $1 - DL$  اصلی اولین لیترال  $z$  را در نظر بگیرید که  $S_z - S$  ناتهی است؛ در واقع  $z$  اولین لیترالی است که در  $1 - DL$  اصلی، حداقل یکی از اعضای  $S$  (مثال‌های باقی‌مانده) در آن حذف می‌شود. در این صورت  $z$  یک لیترال مفید است.

هر  $1 - DL$  با  $n$  متغیر، با  $n \log n$  بیت کد می‌شود. زمان اجرای الگوریتم ارائه شده، بر حسب  $m$  چند جمله‌ای است. طبق نسخه اندازه‌ای قضیه تیغ اوکام، اگر  $m \geq \alpha((\frac{1}{\epsilon})(\log \frac{1}{\delta} + n \log n))$ ، یک الگوریتم یادگیری کارای PAC با پارامترهای  $\epsilon$  و  $\delta$  داریم.

□

<sup>۵</sup>k- decision list

## ۲ یادگیری PAC انکاری

تا کنون فرض می‌کردیم مثال‌ها از یک توزیع  $D$  روی فضای موردها تولید شده و با یک مفهوم هدف  $c$  برچسب‌گذاری شده‌اند. در حالت کلی می‌توانیم فرض کنیم که لزوماً یک مفهوم هدف وجود ندارد و مثال‌ها از یک توزیع  $D$  روی  $\mathcal{X} \times \mathcal{Y}$  تولید می‌شوند که  $\mathcal{X}$ ، فضای موردها و  $\mathcal{Y}$ ، مجموعه برچسب‌هاست. در این حالت باید یک مدل کلی‌تر برای یادگیری PAC ارائه دهیم.

خطای  $h$  نسبت به توزیع  $D$  به صورت زیر تعریف می‌شود:

$$L_D(h) = \Pr_{(x,y) \sim D}[h(x) \neq y] \quad (29)$$

**تعریف ۱۰.** گوییم  $\mathcal{H}$  قابل یادگیری PAC انکاری<sup>۶</sup> است، اگر یک الگوریتم یادگیری  $L$  وجود داشته باشد که به ازای هر  $\epsilon, \delta < 1$  و هر توزیع  $D$  روی  $\mathcal{X} \times \mathcal{Y}$ ، یک عدد طبیعی  $m$  وجود داشته باشد که اگر  $m$  مثال مستقل با توزیع  $D$  دریافت کند، خروجی  $h$  را تولید می‌کند که با احتمال حداقل  $1 - \delta$  خاصیت زیر را دارد:

$$L_D(h) \leq \min_{h' \in \mathcal{H}} L_D(h') + \epsilon \quad (30)$$

## ۳ کمینه‌سازی ریسک تجربی (ERM)

هدف کلی در یادگیری این است که فرضیه‌ای را تولید کنیم که خطای آن نسبت به مفهوم هدف و توزیع  $D$  روی تمام فضای موردها کم باشد. اما یادگیرنده مفهوم هدف و توزیع  $D$  را نمی‌داند؛ یک روش برای تخمین خطای فرضیه این است که خطای تجربی (خطا روی نمونه) فرضیه  $h$  را به دست آوریم. خطای  $h$  روی نمونه  $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$  به صورت زیر تعریف می‌شود:

$$L_S(h) = \frac{1}{m} \sum_{i=1}^m I(h(x_i) \neq y_i) \quad (31)$$

**تعریف ۱۱.** اگر الگوریتم یادگیری  $L$  به ازای هر نمونه  $S$ ، خروجی  $h$  با  $L_S(h)$  کمینه را تولید کند، به آن یک الگوریتم  $ERM$ <sup>۷</sup> می‌گوییم.

## ۴ همگرایی یکنواخت

**تعریف ۱۲.** گوییم نمونه  $S$  یک  $\epsilon$ -نماینده از  $\mathcal{X}, D$  و  $\mathcal{H}$  است، اگر

$$\forall h \in \mathcal{H} : |L_D(h) - L_S(h)| \leq \epsilon \quad (32)$$

**قضیه ۱۳.** اگر  $S$  یک  $\frac{\epsilon}{4}$ -نماینده از  $\mathcal{X}, D$  و  $\mathcal{H}$  باشد و  $h$  خروجی یک الگوریتم  $ERM$  با ورودی  $S$  باشد، داریم:

$$L_D(h) \leq \min_{h' \in \mathcal{H}} L_D(h') + \epsilon \quad (33)$$

اثبات. به ازای هر  $h' \in \mathcal{H}$  داریم:

$$L_D(h) \leq L_S(h) + \frac{\epsilon}{4} \leq L_S(h') + \frac{\epsilon}{4} \leq L_D(h') + \epsilon \quad (34)$$

□

چون این نامساوی به ازای هر  $h' \in \mathcal{H}$  برقرار است، حکم قضیه ثابت می‌شود.

<sup>۶</sup> Agnostic PAC Learnable

<sup>۷</sup> Empirical Risk Minimization



**تعریف ۱۴.** گوئیم کلاس فرضیه  $\mathcal{H}$  همگرایی یکنواخت<sup>۸</sup> دارد، اگر به ازای هر  $\epsilon$  و  $\delta$  و هر توزیع  $D$ ، یک عدد طبیعی  $m$  وجود داشته باشد به طوری که اگر نمونه  $S$  شامل  $m$  مثال مستقل تولید شده از توزیع  $D$  باشد، با احتمال حداقل  $1 - \delta$ ، نمونه  $S$  یک  $\epsilon$ -نماینده از  $\mathcal{X}$ ،  $D$  و  $\mathcal{H}$  باشد.

**لم ۱۵.** اگر  $\mathcal{H}$  خاصیت همگرایی یکنواخت داشته باشد،  $\mathcal{H}$  قابل یادگیری PAC انکاری است.

**اثبات.** پارامترهای  $\epsilon$  و  $\delta$  و توزیع دلخواه  $D$  در نظر بگیرید. به ازای  $\epsilon$  و  $\delta$ ، یک  $m$  وجود دارد که اگر نمونه  $S$  شامل  $m$  مثال مستقل تولید شده از توزیع  $D$  باشد، به احتمال حداقل  $1 - \delta$ ،  $S$  یک  $\epsilon$ -نماینده از  $\mathcal{X}$ ،  $D$  و  $\mathcal{H}$  است. بنابراین اگر نمونه  $S$  را به یک الگوریتم  $ERM$  بدهیم، خروجی  $h$  را تولید می‌کند که دارای خاصیت زیر است:

$$L_D(h) \leq \min_{h' \in \mathcal{H}} L_D(h') + \epsilon \quad (35)$$

□

**مثال ۱۶.** کلاس همه مستطیل‌های با اضلاع موازی محورهای مختصات خاصیت همگرایی یکنواخت دارد.

**مثال ۱۷.** کلاس همه توابع  $f: \{0, 1\}^m \rightarrow \{0, 1\}$  قابل یادگیری PAC انکاری نیست و در نتیجه همگرایی یکنواخت ندارد.

**قضیه ۱۸.** اگر  $\mathcal{H}$  متناهی باشد، قابل یادگیری PAC انکاری است.

**اثبات.**

$$Pr_S[\exists h \in \mathcal{H} : |L_D(h) - L_S(h)| > \epsilon] \leq \sum_{h \in \mathcal{H}} Pr_S[|L_D(h) - L_S(h)| > \epsilon] \leq \sum_{h \in \mathcal{H}} \sum_{\epsilon \in \mathcal{H}} \leq 2|\mathcal{H}|e^{-2m\epsilon^2} \quad (36)$$

توجه کنید که  $L_S(h)$  یک متغیر تصادفی با میانگین  $L_D(h)$  است و نامساوی بالا با استفاده از کران چرنوف به دست آمده است. قرار می‌دهیم:

$$2|\mathcal{H}|e^{-2m\epsilon^2} < \delta \quad (37)$$

$$m > \frac{\log \frac{2|\mathcal{H}|}{\delta}}{2\epsilon^2} \quad (38)$$

□

در این صورت،  $\mathcal{H}$  همگرایی یکنواخت دارد و در نتیجه قابل یادگیری PAC انکاری است.

**قضیه ۱۹.** نهار مجانی در کار نیست!<sup>۹</sup> فرض کنید  $A$  یک الگوریتم یادگیری برای فضای موردی  $\mathcal{X}$  است. اگر  $m < \frac{|\mathcal{X}|}{4}$ ، آنگاه یک توزیع  $D$  روی  $\mathcal{X} \times \{0, 1\}$  وجود دارد به طوری که:

$$1. \text{ یک تابع } f: \mathcal{X} \rightarrow \{0, 1\} \text{ وجود دارد که } L_D(f) = 0$$

۲. اگر یک نمونه  $S$  شامل  $m$  مثال مستقل تولید شده از توزیع  $D$  داشته باشیم، با احتمال حداقل  $\frac{1}{4}$ ، داریم،  $L_D(A(S)) \geq \frac{1}{8}$

یک نتیجه مهم که از این قضیه می‌گیریم این است که اگر  $\mathcal{H}$  مجموعه همه توابع  $f: \mathcal{X} \rightarrow \{0, 1\}$  باشد، تقریباً نمی‌توانیم  $\mathcal{H}$  را یاد بگیریم، مگر اینکه اندازه نمونه خیلی بزرگ باشد و نمونه را حفظ کنیم. در واقع باید کلاس فرضیه را محدود کنیم تا بتوانیم با قواعد ساده (نه چندان پیچیده) آن را یاد بگیریم.

حال به اثبات قضیه می‌پردازیم:

<sup>۸</sup>Uniform Convergence

<sup>۹</sup>No free lunch

اثبات. فرض کنید  $m > |\mathcal{X}|$ . توزیع  $\mathcal{D}$  را یک توزیع یکنواخت روی  $\mathcal{X}$  تعریف می‌کنیم. فرض کنید  $f$  یک تابع کاملاً تصادفی از  $\mathcal{X}$  به  $\{0, 1\}$  است. یعنی به ازای هر  $x \in \mathcal{X}$ ،  $f(x)$  با احتمال  $\frac{1}{2}$ ، ۰ و با احتمال  $\frac{1}{2}$ ، ۱ است. فرض کنید نمونه  $S$  شامل  $m$  مثال مستقل با توزیع  $\mathcal{D}$  از  $\mathcal{X}$  تولید شده و با تابع  $f$  برچسب‌گذاری شده است. الگوریتم  $A$ ، نمونه  $S$  را دریافت می‌کند و یک فرضیه  $h = A(S)$  را تولید می‌کند. می‌دانیم که حداقل  $m$  مورد  $x \in \mathcal{X}$  وجود دارد که در نمونه  $S$  نیامده‌اند. چون مفهوم هدف، یک تابع کاملاً تصادفی است، بهترین تصمیمی که الگوریتم  $A$  می‌تواند برای موارد  $x \in \mathcal{X} - S$  بگیرد، این است که برچسب آن‌ها را حدس بزند. یعنی با احتمال  $\frac{1}{2}$ ، هر یک از برچسب‌های ۰ و ۱ را انتخاب کند. در این صورت، به هر یک از این موارد با احتمال  $\frac{1}{2}$ ، برچسب اشتباه نسبت می‌دهد. بنابراین داریم:

$$E_f[E_S[L_{\mathcal{D}}(A(S))]] = E_f\left[\frac{1}{\sqrt{m}} \sum_{x \in \mathcal{X}} \Pr\{h(x) \neq f(x)\}\right] \quad (39)$$

$$\geq E_f\left[\frac{1}{\sqrt{m}} \sum_{x \in \mathcal{X} - S} \Pr\{h(x) \neq f(x)\}\right] \quad (40)$$

$$= \frac{1}{\sqrt{m}} \frac{1}{2} |\mathcal{X} - S| \geq \frac{1}{2} \quad (41)$$

$$\geq \frac{1}{\sqrt{m}} \left(\frac{1}{2} m\right) = \frac{1}{2} \quad (42)$$

بنابراین یک  $f$  وجود دارد که به ازای آن داریم:

$$E_S[L_{\mathcal{D}}(A(S))] \geq \frac{1}{2}$$

که در آن  $\mathcal{D}$  توزیعی روی  $\mathcal{X} \times \{0, 1\}$  به صورت زیر است:

$$D(x, y) = \begin{cases} \frac{1}{|\mathcal{X}|} & \text{اگر } y = f(x) \\ 0 & \text{در غیر این صورت} \end{cases}$$

واضح است که  $L_{\mathcal{D}}(f) = 0$ . در ادامه ثابت می‌کنیم با احتمال حداقل  $\frac{1}{2}$  داریم،  $L_{\mathcal{D}}(A(S)) \geq \frac{1}{2}$ :

$L_{\mathcal{D}}(A(S))$  یک متغیر تصادفی است که مقدار آن در بازه  $[0, 1]$  است. این متغیر تصادفی را با  $Y$  نشان می‌دهیم. همچنین قرار

می‌دهیم،  $p = \Pr\{Y \geq \frac{1}{2}\}$ . داریم:

$$\mathbb{E}[Y] = \mathbb{E}[Y|Y \geq \frac{1}{2}] \Pr\{Y \geq \frac{1}{2}\} + \mathbb{E}[Y|Y < \frac{1}{2}] \Pr\{Y < \frac{1}{2}\} \quad (43)$$

$$= p \mathbb{E}[Y|Y \geq \frac{1}{2}] + (1-p) \mathbb{E}[Y|Y < \frac{1}{2}] \quad (44)$$

$$\leq p + (1-p) \frac{1}{2} = \frac{1}{2} p + \frac{1}{2} \quad (45)$$

$$\Rightarrow \frac{1}{2} p + \frac{1}{2} \geq \mathbb{E}[Y] \geq \frac{1}{2} \quad (46)$$

$$\Rightarrow p \geq \frac{1}{2} \quad (47)$$

□