

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

نظریه یادگیری ماشین

حمید بیگی

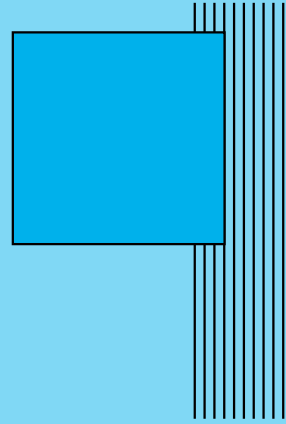
دانشگاه صنعتی شریف

بهار ۱۴۰۱

حق انتشار

درس نامه پیش رو از روی یادداشت‌های درس نظریه یادگیری ماشین که در دانشگاه صنعتی شریف تدریس می‌شود نوشته شده است. با توجه به اینکه ممکن است ایرادات نگارشی در متن این درس نامه وجود داشته باشد خواهشمند است در صورت وجود هر گونه اشتباهی اعم از نگارشی و یا منطقی خواهشمند است از طریق رایانامه beigy@sharif.edu به آگاهی بنده برسانید. همچنین از انتشار این درس نامه به دلیل امکان وجود ایرادهای نگارشی به افراد دیگر بدون آگاهی بنده خوداری نمایید.

حق چاپ برای ناشر محفوظ است.



فهرست مطالب

ت	فهرست نمادها
۱	۱ دیباچه
۱	۱.۱ یادگیری ماشین چیست؟
۲	۲.۱ اهداف پژوهش های یادگیری ماشین
۲	۳.۱ نظریه یادگیری ماشین
۵	۲ مدل های رسمی یادگیری
۵	۱.۲ مولفه های مدل رسمی یادگیری
۸	۲.۲ مدل سازگاری
۱۱	۳.۲ مدل یادگیری احتمالا تقریبا درست
۱۷	۱.۳.۲ پیچیدگی نمونه ای برای فضای فرضیه متناهی
۲۰	۴.۲ مدل یادگیری احتمالا تقریبا درست بدون پیش فرض
۲۷	۳ بعد VC و پیچیدگی رادمیچرا
۲۸	۱.۳ تابع رشد
۳۰	۲.۳ بعد VC

۳۷	پیچیدگی رادامیچر	۳.۳
۴۵	یادگیرهای عمومی	۴.۳

۴ یادگیری نایکنواخت ۴۷

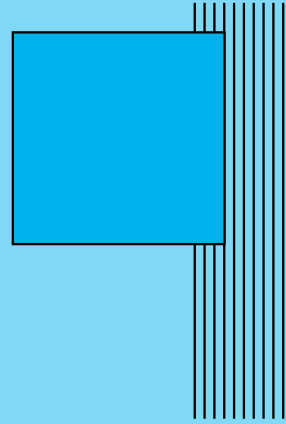
۴۷	قابلیت یادگیری نایکنواخت	۱.۴
۴۹	کمیته سازی هزینه ساختاری	۲.۴
۵۵	کمیته طول توصیف	۳.۴
۵۸	سازگاری	۴.۴
۵۸	جمع بندی نمادهای مختلف یادگیری	۵.۴
۶۰	پرسش‌ها	

۵ پیچیدگی محاسباتی الگوریتم‌های یادگیری ۶۱

۶۲	پیچیدگی محاسباتی یادگیری	۱.۵
۶۵	پیاده سازی قاعده کمیته سازی خطای تجربی	۲.۵
۶۵	فضای فرضیه متناهی	۱.۲.۵
۶۶	مستطیل‌هایی با اضلاع موازی محورهای مختصات	۲.۲.۵
۶۷	ترکیب عطفی تعدادی ویژگی دودویی	۳.۲.۵
۶۷	یادگیری کلاس مفاهیم 3-term DNF	۴.۲.۵
۶۸	یادگیری مستقل از نمایش، سخت نیست	۳.۵
۶۸	سختی یادگیری	۴.۵

۷۱ واژه‌نامه فارسی به انگلیسی

۷۳ نمایه



فهرست نمادها

۶	فضای نمونه	\mathcal{X}
۶	مجموعه برچسب ها	\mathcal{Y}
۶	مجموعه آموزشی	S
۶	فرضیه	h
۶	توزیع احتمال نمونه برداری	\mathcal{D}
۶	مفهوم	c
۱۲	خطای واقعی فرضیه h	$\mathbf{R}(h)$
۱۳	خطای تجربی فرضیه h	$\hat{\mathbf{R}}(h)$
۱۳	پیچیدگی نمونه‌ای	m_H



دیباچه

۱.۱ یادگیری ماشین چیست؟

یادگیری ماشین مطالعه روش های خودکاری است که پیش بینی های دقیق یا کنش های مناسب را براساس تجربیات گذشته یا مشاهده ها انجام می دهد. الگوریتم های یادگیری ماشین باید در زمان مناسب و با حافظه مناسب قابل اجرا باشند و به حجم داده زیادی نیاز نداشته باشند. به این الگوریتم ها، الگوریتم های کارا^۱ می گویند. در چه زمانی به یادگیری ماشین نیاز مندیم؟

۱. وظایف بسیار پیچیده هستند تا بتوان آنها را برنامه نویسی (پیاده سازی) نمود.

(آ) وظایفی که انجام آنها برای انسان ها یا حیوان ها بسیار ساده هستند اما برای رایانه ها بسیار سخت هستند.

(ب) وظایفی که انجام آنها بیشتر ز توان انسان ها است. برای نمونه تحلیل داده های ستاره شناسی، پیش بینی هوا، تحلیل داده های زیستی، تجارت الکترونیک و جویشرهای وب نمونه ای از وظایف هستند.

۲. بسیاری از وظایف با زمان تغییر می کنند و برنامه های رایانه ای می بایست قابلیت تطبیق داشته باشند تا توانایی کارکردن برای این وظایف را داشته باشند. برای نمونه تشخیص هرزنامه نمونه ای از وظایفی است که نیاز به تطبیق دارد.

یادگیری ماشین دارای چه کاربردهایی هستند؟ برخی از کاربردهای یادگیری ماشین عبارتند از

۱. دسته بندی متون مانند دسته بندی نامه ها

۲. پردازش زبان طبیعی مانند بازشناسی موجودیت های نامدار

¹Efficient

۳. بازشناسی گفتار و تبدیل گفتار به متن

۴. بازشناسی نوری حروف

۵. کاربردهای زیست شناسی

۶. بینایی ماشین

۷. تشخیص تقلب

۸. بازی های رایانه ای

۹. تشخیص های پزشکی

۱۰. سامانه های تحلیل نظر و پیشنهاد دهنده ها

دسته بندی های مختلف روش های یادگیری ماشین

۱. روش های با نظارت در مقابل روش های بی نظارت، نیمه نظارتی، و یادگیری تقویتی

۲. روش های یادگیری فعال در مقابل روش های یادگیری غیر فعال

۳. روش های یادگیری برخط در مقابل روش های یادگیری دسته ای

مسائل یادگیری ماشین را می توان به دسته های مختلفی تقسیم نمود. مهم ترین این دسته بندی ها عبارتند از

۱. دسته بندی

۲. رگرسیون

۳. رتبه بندی

۴. خوشه بندی

۵. کاهش بعد

۶. یادگیری متغیرهای پنهان

۲.۱ اهداف پژوهش های یادگیری ماشین

۱. کارا باشند

۲. پاسخ های تولید شده تفسیر پذیر باشند

۳. پاسخ تولید شده تا حد ممکن دقیق باشد

۳.۱ نظریه یادگیری ماشین

هدف نظریه یادگیری ماشین توسعه و تحلیل مدل های رسمی است که به فهم مسایل مهم در یادگیری ماشین کمک کند. این نظریه کمک می کند تا بفهمیم کدام مفاهیم را می توان به صورت کارا یادگرفت و چه اندازه داده برای یادگیری این مفاهیم مورد نیاز است. همچنین این نظریه تلاش می کند تا پاسخ پرسشی که چرا یک الگوریتم در همه شرایط خوب کار نمی کند را بیابد. برخی از پرسش های مهم که در نظریه یادگیری ماشین مطرح هستند عبارتند از

۱. تعداد نمونه مورد نیاز برای یک یادگیری موفق
۲. پیچیدگی محاسباتی یک مساله و/یا الگوریتم یادگیری
۳. پیدا نمودن پیچیدگی محاسباتی یک الگوریتم یادگیری براساس تعداد نمونه های آموزشی
۴. دسته بندی مسایل یادگیری به مسایل ساده و سخت.
۵. چگونگی به کار بردن اطلاعات زمینه و پیشین در مورد مساله یادگیری برای یادگیری موثر
۶. یافتن پاسخ برای بهتر بودن فرضیه های ساده تر
۷. چگونگی تغییر کارایی یادگیر بر اساس چگونگی تحویل داده های آموزشی
۸. اگر یادگیر نمونه آموزشی درخواست نماید پیچیدگی نمونه ای یادگیر چگونه تغییر می یابد؟
۹. آیا داده های برجسب نخورده می توانند کمکی به آموزش / افزایش کارایی الگوریتم یادگیری نمایند؟
۱۰. اگر یادگیری به صورت توزیع شده صورت پذیرد هزینه ارتباطی چه اندازه خواهد بود؟
۱۱. چگونه می توان از یادگیر ضعیف استفاده نمود؟

مدل های رسمی یادگیری

برای تحلیل و مطالعه الگوریتم های یادگیری ماشین نیاز است تا مساله یادگیری به صورت رسمی و دقیق تعریف شود. این تعریف دقیق، مدل یادگیری نام دارد. مدل یادگیری باید به اندازه کافی دقیق باشد و توانایی مدل سازی وجه های مهم مسایل یادگیری واقعی را داشته باشد. دقت مدل به این دلیل مورد نیاز است که بتوان مساله یادگیری را به صورت ریاضی مطالعه نمود. همچنین یک مدل یادگیری خوب میبایست در مقابل تغییرات جزئی در تعریف آن مقاوم باشد. برای هر مدل رسمی فرض های ساده کننده غیرقابل اجتناب است. مدل های یادگیری میبایست توانایی پاسخ به پرسش های گوناگون همانند پرسش های زیر را داشته باشد.

۱. چه چیزی باید یادگرفته شود؟

۲. داده ها چگونه تولید میشوند؟

۳. داده ها چگونه به یادگیر داده میشوند؟ برای نمونه داده ها یکی یکی یا به صورت دسته ای به یادگیر داده میشوند؟

۴. هدف از یادگیری در این مدل چیست؟

در ادامه این بخش، نخست به بررسی مولفه های مدل های رسمی یادگیری می پردازیم و سپس مدل های مختلف یادگیری بررسی و نقاط ضعف و قوت آنها تحلیل می گردند.

۱.۲ مولفه های مدل رسمی یادگیری

مدل های رسمی یادگیری دارای پنج مولفه اصلی هستند: ورودی یادگیر، خروجی یادگیر، مدل تولید داده، معیار کارایی و دانش پیشین. در ادامه هر کدام از این مولفه ها شرح داده می شوند.

۱. یکی از مولفه های هر مدل یادگیری، ورودی یادگیر است. در مدل رسمی یادگیری، یادگیر به داده های زیر دسترسی دارد.
- فضای نمونه که با \mathcal{X} نشان داده می شود عبارت است از مجموعه نمونه هایی که ممکن است در آینده برچسب آنها تعیین شوند. به هریک از اعضا درون \mathcal{X} یک نمونه می گویند. یک نمونه به وسیله ی مجموعه ای از خصوصیات یا ویژگی ها توصیف شده و معمولاً با یک بردار n بعدی بازنمایی می شود.
 - مجموعه برچسب ها که با \mathcal{Y} نشان داده می شود عبارت است از برچسب یا گروهی که می خواهیم یادگیر آن را پیش بینی نماید. برای نمونه در تشخیص اعداد، برچسب نمونه ها از مجموعه $\mathcal{Y} = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$ انتخاب می شوند و همچنین $\mathcal{Y} = \{0, 1\}$ و $\mathcal{Y} = \{-1, +1\}$ دو نمونه از مجموعه برچسب های ممکن برای دسته بندی دودویی هستند.
 - در هنگام آموزش یادگیر، مجموعه ای از نمونه های برچسب دار که مجموعه آموزشی نامیده میشود به عنوان ورودی به الگوریتم یادگیری داده می شود و در هنگام آزمایش، تنها نمونه های آزمایشی که برچسب ندارند به یادگیر داده می شود و یادگیر می بایست برچسب آنها را تعیین نماید. مجموعه آموزشی به صورت $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$ نشان داده می شود که دنباله ای از زوج های مرتب $(x, y) \in \mathcal{X} \times \mathcal{Y}$ است. به این مجموعه نقاط برچسب خورده دامنه نیز گفته می شود.
۲. مولفه دیگر مدل رسمی یادگیر، خروجی آن است. یک یادگیر می بایست یک خروجی به شکل $\mathcal{Y} : \mathcal{X} \rightarrow \mathcal{Y}$ تولید نماید. به این خروجی فرضیه یا در مسایل دسته بندی، دسته بند می گویند. این فرضیه بعداً برای پیش بینی برچسب نمونه های آزمایشی (نمونه های برچسب نخورده) که از فضای \mathcal{X} با توزیع D نمونه برداری شده اند به کار می رود. در اینجا نماد $A(S)$ را برای نشان دادن فرضیه ای که الگوریتم یادگیری A با دریافت مجموعه آموزشی S تولید می نماید به کار می بریم. در واقع یادگیری به نوعی ایجاد یک الگوریتم A است که اگر به آن مجموعه آموزشی S را بدهیم به ما h را به عنوان خروجی تحویل خواهد داد.
۳. شیوه نمونه برداری نمونه ها یا به عبارتی دیگر مدل تولید داده ها بخش دیگری از مدل رسمی یادگیر است. در این جا فرض می کنیم داده های آموزشی و آزمایشی از یک توزیع احتمال D از فضای نمونه ها، نمونه برداری شده باشند و یادگیر هیچ گونه اطلاعی از این توزیع ندارد. برای ایجاد مجموعه آموزشی، پس از نمونه برداری هر نمونه x ، برچسب آن توسط تابع $\mathcal{Y} : \mathcal{X} \rightarrow \mathcal{Y} : c$ تعیین میشود. به تابع c مفهوم می گویند و آموزگار آن را می داند. یادگیر شیوه برچسب زنی آموزگار یا همان تابع c را نمی داند. اغلب فرض می کنیم که نمونه ها با یک مفهوم ناشناخته که متعلق به یک کلاس مفاهیم شناخته شده است برچسب زده می شود.
۴. مولفه دیگر مدل رسمی یادگیری، معیار ارزیابی کارایی فرضیه تولید شده است. کارایی یک فرضیه را می توان بر اساس معیارهای متفاوتی ارزیابی نمود. یکی از این معیارها خطای واقعی فرضیه است که به صورت زیر تعریف میشود. خطای واقعی فرضیه h برابر است با احتمال اینکه نمونه تصادفی x که با توزیع D نمونه برداری شده است به فرضیه h داده شود و برچسب فرضیه h با برچسب واقعی آن یکسان نباشد. به عبارتی دیگر خطای فرضیه h به صورت زیر تعریف میشود.

$$\mathbf{R}(h) = \mathbb{P}_{x \sim D} [h(x) \neq c(x)] \quad (1.2)$$

- این خطا با نام های دیگری مانند خطای تعمیم، ریسک و خطای واقعی نام برده می شود و معمولاً به جای همدیگر به کار می روند.
۵. علاوه بر مولفه های یاد شده، یک یادگیر به دانش پیشین در فرایند یادگیری نیازمند است. یادگیر اطلاعاتی در باره توزیع D یا تابع c ندارد و تنها به مجموعه آموزشی دسترسی دارد. ممکن است برای یادگیری فرضیاتی در باره مساله داشته باشیم که به آنها Inductive Bias می گویند. بر اساس این فرضیات تعیین شده، الگوریتم یادگیری طراحی می شود.
- در ادامه این بخش برخی از بخش ها که پیش از این به آنها اشاره شد به صورت رسمی تعریف می گردد.

تعریف ۱.۲ مفهوم

یک مفهوم یک تابع دودویی در فضای نمونه است و به صورت زیر تعریف می‌گردد.

$$c: \mathcal{X} \mapsto \{0, 1\} \quad (2.2)$$

تابع c یا همان مفهوم ناشناخته است و آموزگار این تابع را می‌داند و بر اساس این تابع نمونه های دریافت شده از فضای نمونه را برچسب می‌زند. برای نمونه مفهوم $c(x_1, x_2, \dots, x_n) = x_1 \wedge x_2$ یک تابع دودویی است و مقدار آن برای یک نمونه زمانی یک است که مقدار دو ویژگی نخست نمونه برابر یک باشد. در بیشتر مواقع از واژه های دسته بند، فرضیه و قاعده پیش بینی برای معنای یک مفهوم به کار می‌بریم. برای نمونه در فضای نمونه حیوانات، تفکیک حیوانات از هم یک کلاس مفاهیم را تشکیل می‌دهد و تشخیص سگ از بقیه حیوانات یک مفهوم است.

تعریف ۲.۲ کلاس مفاهیم

منظور از کلاس مفاهیم مجموعه‌ای از مفاهیم همراه با بازنمایی متناظر آن است.

بازنمایی یک کلاس از مفاهیم به این علت مهم است که اندازه‌ی مفاهیم به کمک آن به دست می‌آید.

تعریف ۳.۲ اندازه مفهوم

منظور از اندازه مفهوم عبارت است از تعداد بیت‌های لازم برای بازنمایی مفهوم در بازنمایی تعریف شده آن.

برای روشن شدن موضوع، در ادامه یک مثال از یک مفهوم به همراه بازنمایی‌های مختلف آن و اندازه مفهوم در هر کدام از این بازنمایی‌ها بیان می‌شود.

مثال ۱.۲ (ترکیب عطفی بکنوا).

کلاس عطف بکنوا را در نظر بگیرید. برای نمونه $c(x_1, x_2, \dots, x_n) = x_i \wedge x_j$ یکی از اعضای این کلاس است. در این توابع نقیض متغیرها را نداریم. این کلاس را می‌توان به راه‌های متفاوتی بازنمایی داد که در ادامه به چند نمونه از آن‌ها می‌پردازیم.

- یک راه بازنمایی این کلاس مفاهیم استفاده از یک آرایه‌ی n بیتی است به نوعی که اگر x_i در مفهوم وجود داشته باشد، بیت i ام برابر مقدار ۱ و در غیر اینصورت برابر ۰ است. در این حالت اندازه‌ی مفهوم از مرتبه $O(n)$ خواهد بود.
- یک راه دیگر بازنمایی این کلاس از مفاهیم این است که به ازای هر ورودی ممکن برچسب آن (۰ یا ۱ بودن) را در نظر بگیریم. در این حالت با توجه به این‌که 2^n ورودی ممکن خواهیم داشت اندازه‌ی مفهوم نیز از مرتبه $O(2^n)$ خواهد بود.
- راه دیگر این است که زیرنویس و اندیس متغیرهای استفاده شده در ترکیب عطفی را به شکل عددهای $\lg n$ بیتی در کنار هم بنویسیم. در این حالت اگر مفهوم مورد نظر k متغیر عطفی داشته باشد، آنگاه اندازه‌ی مفهوم از مرتبه $O(k \times \lg n)$ خواهد بود.

لذا روشن است که انتخاب شیوه مناسب بازنمایی مفاهیم، اندازه مفهوم را تغییر می‌دهد و در هنگامی که الگوریتم یادگیری تابعی از اندازه مفهوم باشد شیوه بازنمایی مفهوم در کارایی الگوریتم یادگیری تاثیر دارد.

تمرین ۱.۲ (اندازه مفهوم درخت تصمیم) اندازه مفهوم کلاس مفاهیم درخت تصمیم در فضای $\{0, 1\}^n$ چه اندازه است؟

۲.۲ مدل سازگاری

در ادامه این بخش، نخستین مدل رسمی یادگیری بیان میشود. این مدل بر اساس سازگاری فرضیه تولید شده توسط الگوریتم یادگیری بنا شده است. این مدل بسیار ساده است و توانایی بکارگیری در کاربردهای واقعی را نداشته اما برای آغاز بررسی مدل های یادگیری، مدل بسیار مناسبی است.

تعریف ۴.۲ فرضیه سازگار

فرضیه h با داده های آموزشی سازگار است اگر و تنها اگر برای همه ی نمونه های مجموعه آموزشی، داشته باشیم

$$h(x) = c(x) \quad \forall x \in S \quad (۳.۲)$$

یا به عبارتی دیگر برای همه نمونه های آموزشی، برچسب تولید شده توسط فرضیه h با برچسب اصلی آن یکسان باشد.

بر اساس تعریف فرضیه های سازگار، مدل سازگاری به صورت زیر تعریف میشود.

تعریف ۵.۲ یادگیری در مدل سازگاری

الگوریتم یادگیری A کلاس مفاهیم C را در مدل سازگاری یاد می گیرد، اگر برای هر مجموعه ی آموزشی S ، الگوریتم A مفهوم $c \in C$ که سازگار با S باشد را در صورت وجود تولید و در غیر این صورت پاسخ "فرضیه ای وجود ندارد" را تولید می نماید.

تعریف ۶.۲ یادگیری کارا در مدل سازگاری

کلاس مفاهیم C در مدل سازگاری به صورت کارا یادگرفتنی است اگر یک الگوریتم کارای A وجود داشته باشد که کلاس مفاهیم C را در مدل سازگاری یاد بگیرد. کارایی براساس اندازه ی مجموعه S که با m نشان داده می شود و اندازه ی نمونه های ورودی که با n نشان داده می شود تعریف می شود.

الگوریتم کارا معمولاً به الگوریتمی گفته می شود که در زمان چندجمله ای براساس اندازه S و اندازه x مساله را حل می کند.

مثال ۲.۲ (توابع عطفی یکنوا).

این کلاس مفاهیم شامل تمامی مفاهیمی است که می توان از ترکیب عطفی یکنوا حداکثر n متغیر دودویی ساخت. در این ترکیب عطفی، نقیض متغیرها ظاهر نمی شوند. برای نمونه، مجموعه آموزشی زیر که در آن یک تابع عطفی سازگار وجود دارد را در نظر بگیرید.

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	y
1	0	1	1	0	0	1	1	+
1	1	1	1	1	0	1	0	+
0	1	1	1	0	0	1	1	+
0	0	0	1	1	1	1	1	-
1	1	1	1	1	0	0	0	-

الگوریتم زیر را برای حل این مثال در نظر بگیرید :

۱. در نمونه های مثبت، تمام ویژگی هایی که دارای مقدار صفر هستند را حذف کن.

۲. ترکیب عطفی ویژگی‌های باقیمانده را حساب کن.
 ۳. اگر فرضیه تولید شده با نمونه‌های منفی هم سازگار بود این فرضیه را به عنوان خروجی اعلام و در غیر اینصورت پاسخ "فرضیه‌ای وجود ندارد" را تولید کن.
 با مجموعه آموزشی یاد شده در بالا، حاصل اعمال الگوریتم فوق $c(x) = x_۲ \wedge x_۴ \wedge x_۷$ است. با توجه به این‌که این الگوریتم در زمان $O(nm)$ خروجی می‌دهد، بنابراین این الگوریتم کارا است.

در نتیجه این کلاس مفاهیم به صورت کارا یادگرفتنی است و در نتیجه قضیه زیر را خواهیم داشت.

قضیه ۱.۲ (یادگیری ترکیب عطفی یکنوا در مدل سازگاری).
 کلاس مفاهیم ترکیب عطفی یکنوا در مدل سازگاری به صورت کارا یادگرفتنی است.

تمرین ۲.۲ (یادگیری ترکیب عطفی غیر یکنوا) کلاس مفاهیم ترکیب عطفی غیر یکنوا عبارت است از ترکیب عطفی تعدادی متغیر یا نقیض آنها. نشان دهید که این کلاس مفاهیم بصورت کارا یادگرفتنی هستند و برای یادگیری این کلاس مفاهیم یک الگوریتم یادگیری طراحی نموده و کارایی آن را بررسی نمایید.

برای حل این مساله نقیض هر متغیر را با یک متغیر جدید نشان می‌دهیم و مساله را به یادگیری ترکیب عطفی یکنوا تبدیل میکنیم. برای نمونه تابع اگر هدف یادگیری یک ترکیب عطفی غیر یکنوا روی n متغیر است. با تعریف $z_i = \neg x_i$ ، مساله به یادگیری ترکیب عطفی با $2n$ متغیر تبدیل می‌شود و از الگوریتمی که پیش از این ارائه شد می‌توان استفاده نمود.

مثال ۳.۲ (ترکیب فصلی یکنوا).

این کلاس مفاهیم از ترکیب فصلی حداکثر n متغیر (ویژگی) دودویی تولید میشود. در این ترکیب فصلی، نقیض متغیرها ظاهر نمی‌شوند. برای نمونه تابع $c(x_1, x_2, \dots, x_n) = x_۲ \vee x_۶ \vee x_۸$ یک ترکیب فصلی است. به سادگی میتوان نشان داد که با کمک قانون دمورگان، هر ترکیب فصلی را میتوان به ترکیب عطفی تبدیل نمود. برای نمونه اگر $c(x) = x_i \vee x_j$ باشد، می‌توان آن را به صورت ترکیب عطفی زیر نیز درآورد $c(x) = x_i \vee x_j = \overline{x_i} \wedge \overline{x_j}$ که $\overline{x_i}$ یعنی نقیض x_i است. یعنی اگر مقدار x_i برابر ۰ باشد $\overline{x_i}$ برابر ۱ خواهد بود و اگر مقدار x_i برابر ۱ باشد $\overline{x_i}$ برابر ۰ خواهد بود. همچنین علامت تغییر مقدار که بالای مقدار $\overline{x_i} \wedge \overline{x_j}$ قرار گرفته است به معنی تغییر برچسب است.
 فرض کنید می‌خواهیم کلاس ترکیب فصلی را یاد بگیریم. الگوریتمی که ارائه می‌دهیم با کمی تغییر الگوریتم قبلی به دست می‌آید. بدین ترتیب که ابتدا برچسب‌ها و بیت‌ها را وارون می‌کنیم (یعنی اگر مقدار ۰ داشتند مقدارشان را ۱ می‌کنیم و اگر مقدار ۱ داشتند مقدارشان را ۰ در نظر می‌گیریم) و سپس با استفاده از الگوریتم قبلی جواب را به دست می‌آوریم و در نهایت همه چیز را به حالت اولیه برمی‌گردانیم.

مثال ۴.۲ (مساله‌ی فرم نرمال عطفی درجه).

مساله‌ی k فرم نرمال عطفی درجه را در نظر بگیرید. اگر ترکیب فصلی حداکثر k متغیر را یک «عبارت» بنامیم، آنگاه مساله‌ی فرم نرمال عطفی درجه k عبارت است از ترکیب عطفی تعدادی «عبارت». یادگیری در این مساله بدین معنی است که k داده شود و پرسش این باشد که آیا میتوان کلاس مفاهیمی که با فرم نرمال عطفی درجه k نشان داده میشود را یادگرفت یا خیر؟ ما فرض میکنیم که k یک عدد کوچک باشد. ما میتوانیم این مساله را به یادگیری ترکیب عطفی کاهش دهیم. برای این کار متغیر جدیدی برای هر عبارت ممکن در فرم نرمال عطفی درجه k را ایجاد میکنیم.
 در حالت کلی متغیرهای z_{s_i} را از روی ترکیب فصلی s تا از متغیرهای اولیه‌ی x_j ها می‌سازیم. بدین ترتیب تمامی حالات s متغیری را در z_{s_i} ایجاد می‌نماییم. مثلاً $z_{s_i} = \{x_{j_1} \vee x_{j_2} \vee \dots \vee x_{j_s}\}$ برابر i امین متغیر جدید z است که دقیقاً از ترکیب فصلی s متغیر ایجاد شده است. برای نمونه اگر $k = ۲$ باشد برای کلاس مفاهیم فرم نرمال عطفی درجه k متغیرهای زیر را

خواهیم داشت.

$$\begin{aligned} z_1 &= (x_1) \\ z_2 &= (x_2) \\ z_3 &= (x_1 \vee x_2) \\ z_4 &= (x_1 \vee \bar{x}_2) \\ z_5 &= (\bar{x}_1 \vee x_2) \\ z_6 &= (\bar{x}_1 \vee \bar{x}_2) \end{aligned}$$

اگر تعداد متغیرهای جدید ایجاد شده را P بنامیم، خواهیم داشت

$$P = \sum_{i=1}^k \binom{2n}{i} \leq (2n)^k \quad (4.2)$$

بنابراین اگر تعداد نمونه‌ها را m عدد فرض کنیم، الگوریتم ذکر شده از $O(m(2n)^k)$ است که برای k های کوچک، مثلاً ۲ یا ۳، الگوریتم کارا است و در غیراینصورت الگوریتم کارا نخواهد بود.

^ak-CNF

در این مثال، نمونه های ۲ بیتی به نمونه های ۶ بیتی تبدیل میشوند که در آن بیت i ام در صورتی برابر یک است اگر z_i برابر یک باشد و در غیر اینصورت برابر با صفر خواهد بود. اگر نمونه های آموزشی را به الگوریتم یادگیری که برای یادگیری ترکیب عطفی بدهیم مفهوم سازگار با داده های آموزشی را بر اساس z_i تولید خواهد نمود. سپس با جایگزینی z_i به ترکیب فصلی تعریف شده آن، مفهوم مساله‌ی فرم نرمال عطفی درجه k سازگار را پیدا میکند. برای یادگیری این کلاس مفاهیم در حالت کلی نیاز به تعریف $O((2n)^k)$ متغیر z_i نیاز است. این مقدار براساس رابطه زیر بدست می‌آید که تعداد فرم نرمال عطفی درجه k هایی که از n متغیر می‌توان تولید نمود برابر است با

$$(2n)(2n-1)(2n-2)\dots(2n-k) = O((2n)^k). \quad (5.2)$$

زیرا هر مکان در فرم نرمال عطفی درجه k را می‌توان از $2n$ متغیر (خود یا نقیض آن) استفاده نمود. بنابراین اگر k عددی کوچک باشد، این الگوریتم کارا است و اگر k بزرگ باشد زمان اجرای الگوریتم نمایی خواهد بود.

مثال ۵.۲ (کلاس مفاهیم DNF (Or of and)).

فرض کنید می‌خواهیم کلاس مفاهیم DNF را یاد بگیریم. این کلاس به صورت ترکیب فصلی (\vee) تعدادی عبارت که هرکدام به صورت ترکیب عطفی (\wedge) تعدادی متغیر است. برای نمونه

$$c(x_1, x_2, \dots, x_n) = (x_2 \wedge x_7 \wedge x_{10}) \vee (x_1 \wedge x_5) \vee (x_3 \wedge x_4 \wedge x_6) \quad (6.2)$$

یک عبارت DNF است. برای یادگیری این کلاس مفاهیم میتوانیم از الگوریتم زیر استفاده کنیم.

۱. ترکیب فصلی نمونه های مثبت را محاسبه کن.
۲. اگر این ترکیب فصلی با نمونه های منفی سازگار بود آن را بعنوان خروجی و در غیر اینصورت پاسخ "فرضیه‌ای وجود ندارد" را تولید کن.

این الگوریتم کارا است و در زمان $O(nmk)$ قابل پیاده‌سازی است که k حداکثر تعداد عبارت‌های هر نمونه است. بنابراین این کلاس

نیز در مدل سازگاری بصورت کارا یادگرفتنی است. اما فرضیه تولید شده کاربرد ندارد زیرا

۱. این فرضیه داده های مثبت را ذخیره نموده و بیش برآزش رخ میدهد و هر نمونه ای که در مجموعه آموزشی وجود نداشته باشد را به

- عنوان نمونه منفی برجسب میزند. این ایراد اختلاف بین قدرت حافظه و قدرت تعمیم فرضیه را نشان می‌دهد.
۲. اندازه فرضیه تولید شده برابر است با مجموع اندازه نمونه های مثبت مجموعه آموزشی.
- مدل سازگاری، مدلی بسیار ساده است اما در عمل کاربرد ندارد زیرا این مدل دارای ایرادهایی است. برخی از ایرادهای کلی و عمومی مدل سازگاری در زیر آمده است.
۱. در این مدل هیچ سخنی در باره قدرت تعمیم فرضیه تولید شده به میان نیامده است. یک مدل یادگیری خوب می‌بایست در مورد توان (قابلیت) فرضیه تولید شده روی داده های جدید سخن بگوید.
 ۲. مشکل دوم این مدل این است که اندازهی مفاهیم به میزان اندازهی ورودی‌ها (و طبعاً بزرگ‌تر از حالت ایده‌آل) هستند.
 ۳. ایراد دیگر این مدل این است که این مدل به نوبه مقاوم نیست. با توجه به اینکه در کاربردهای واقعی داده ها نویزی هستند و به علت وجود داده‌های نویزی، خطای آموزشی صفر نمی‌شود لذا این مدل در عمل قابل استفاده نیست.
 ۴. و آخرین مشکلی که در این جا بیان می‌کنیم این است که اگر برای کلاس مفاهیم C الگوریتم کارا وجود داشته باشد، برای زیرمجموعه‌های C' مانند C' که $C' \subseteq C$ لزوماً نمی‌توان الگوریتم کارا پیدا نمود. برای بیان بهتر این ایراد به مثال زیر توجه نمایید.

مثال ۶.۲ (یادگیری کلاس مفاهیم فرم نرمال فصلی دو لفظی).

فرض کنید که کلاس مفاهیم C شامل همه فرم نرمال فصلی دو لفظی (یا دو ترکیب عطفی با طول دلخواه) باشد. برای نمونه تابع $c(x_1, \dots, x_n) = (x_2 \wedge x_3 \wedge x_4) \vee (x_1 \wedge x_5)$ یک فرم نرمال فصلی دو لفظی^a است. حال چگونه میتوانیم یک فرم نرمال فصلی دو لفظی از روی داده های آموزشی یاد بگیریم؟ اگر بتوانیم یادگیری مساله فرم نرمال فصلی دو لفظی را به مساله یادگیری فرم نرمال عطفی درجه k تبدیل کنیم کار ساده تر میشود. اگر از قواعد منطق (با در نظر گرفتن \wedge معادل جمع و \vee معادل ضرب) استفاده کنیم می‌توانیم این تبدیل را انجام بدهیم. برای نمونه

$$c(x_1, \dots, x_n) = (x_1 \wedge x_2) \vee (x_3 \wedge x_4) = (x_1 \vee x_3) \wedge (x_1 \vee x_4) \wedge (x_2 \vee x_3) \wedge (x_2 \vee x_4).$$

این مثال نشان می‌دهد که ما همیشه می‌توانیم فرم نرمال فصلی دو لفظی را به فرم نرمال عطفی درجه k تبدیل نماییم. آیا این به معنای این است که این کلاس مفاهیم یادگرفتنی است؟ اگر بتوانیم برعکس تبدیل بالا را یاد بگیریم مساله ساده می‌شود. با تبدیل فرم نرمال فصلی دو لفظی به فرم نرمال عطفی درجه k می‌توانیم فرم نرمال عطفی درجه k را یاد بگیریم و سپس باید پاسخ فرم نرمال عطفی درجه k تولید شده را بتوانیم به فرم نرمال فصلی دو لفظی تبدیل نماییم. مشکل اصلی این است که ما همیشه می‌توانیم یک فرم نرمال فصلی دو لفظی را به فرم نرمال عطفی درجه k تبدیل نماییم اما الزاماً نمی‌توانیم یک فرم نرمال عطفی درجه k را به یک فرم نرمال فصلی دو لفظی تبدیل نماییم. در نتیجه اگر ما یک فرم نرمال عطفی درجه k سازگار را پیدا کنیم بدان معنی نیست که یک فرم نرمال فصلی دو لفظی سازگار وجود داشته باشد. یا بطور غیر رسمی فضای فرم نرمال عطفی درجه k یک ابر فضا از فضای فرم نرمال فصلی دو لفظی است و یا به عبارتی دیگر

$$(2 - \text{termDNF}) \subseteq (2 - \text{CNF}). \quad (7.2)$$

در زمینه نظریه یادگیری، یادگیری فرم نرمال عطفی درجه k آسان است اما یادگیری کلاس مفاهیم به شکل فرم نرمال فصلی دو لفظی یک مساله ان‌پی سخت است.

^a2-term DNF

۳.۲ مدل یادگیری احتمالاً تقریباً درست

در بخش پیش مدل سازگاری بیان گردید. همچنین درباره دو کاستی بزرگ این مدل سخن به میان آمد. نخست آن که این مدل تعریفی از

قدرت تعمیم فرضیه ها ارایه نمی‌کند و دوم آن‌که این مدل امکان پردازش داده‌های نویزی را ندارد. این دو کاستی سبب غیر کاربردی شدن این مدل می‌شود. در این بخش تلاش می‌کنیم که مدل را توسعه دهیم تا در تعریف مدل قدرت تعمیم فرضیه ها ارایه شود و پردازش داده های نویزی را در آینده بررسی خواهیم نمود. در این بخش مدل احتمالا تقریبا درست بیان می‌شود که در آن قدرت تعمیم فرضیه ها در مدل گنجانده شده‌است. برای بررسی قدرت تعمیم فرضیه ها فرض های ساده کننده زیر را در نظر می‌گیریم که بیشتر این فرض ها در آینده تعدیل خواهند شد.

۱. فرض می‌کنیم که نمونه‌های آزمون و آموزش از توزیع ثابت و یکسان اما ناشناخته \mathcal{D} به صورت مستقل نمونه برداری می‌شوند. این توزیع برای یادگیر ناشناخته است.

۲. نمونه ها با کمک یک تابع ثابت و ناشناخته $c \in \mathcal{C}$ برچسب زده می‌شوند. یادگیر نمی‌داند که این تابع چیست و دارای چه شکلی است.

در ادامه نخست خطای واقعی فرضیه h که نشان دهنده قدرت تعمیم آن است را تعریف می‌کنیم.

تعریف ۷.۲ خطای واقعی

خطای واقعی هر فرضیه $h \in H$ که خطای تعمیم آن نیز نامیده میشود به صورت زیر تعریف می‌شود.

$$\begin{aligned} \mathbf{R}(h) &= \mathbb{P}_{x \sim \mathcal{D}} [c(x) \neq h(x)] \\ &= \sum_{c(x) \neq h(x)} \mathcal{D}(x). \end{aligned} \quad (8.2)$$

هدف آرمانی الگوریتم‌های یادگیری این است که فرضیه‌ای را پیدا نمایند که خطای واقعی آن را برابر صفر باشد یعنی همه نمونه هایی که در هنگام آموزش ندیده‌است را بدون خطا برچسب بزند. اما این کار به دو دلیل زیر ممکن نیست. (۱) مجموعه آموزشی یک زیرمجموعه از فضای نمونه است و ممکن است چندین فرضیه سازگار با داده‌های آموزشی داشته‌باشیم و هیچکدام از آنها با مفهوم اصلی (c) یکسان نباشند. (۲) احتمال اینکه نمونه های آموزشی از همه بخش های فضای نمونه، نمونه برداری نشده باشند بزرگتر از صفر است. لذا هیچگاه امکان تولید فرضیه‌ای با خطای واقعی صفر وجود ندارد و به همین جهت خطای واقعی صفر را تعدیل نموده و به جای آن خطای واقعی کوچکتر از ϵ را در نظر می‌گیریم و از آنجایی که داده‌ها از یک فضا به وسیله توزیع \mathcal{D} نمونه برداری می‌شوند یادگیر می‌بایست برای هر ترتیب از داده‌های آموزشی فرضیه با خطای کم را تولید نماید لذا احتمال شکست δ را در نظر گرفته و تلاش می‌کنیم که احتمال شکست کوچکتر از δ باشند. اگر خطای فرضیه h کوچک باشد فرضیه h را یک فرضیه تقریبا درست^۱ می‌گویند. اما همیشه نمی‌توان تضمین نمود که خطای واقعی کوچک باشد زیرا ممکن است داده های آموزشی نماینده خوبی از مفهوم واقعی نباشند به همین جهت باید تعریف بازنگری شده و به دنبال فرضیه‌ای باشیم که با احتمال زیاد تقریبا درست باشد یا به عبارتی دیگر فرضیه h احتمالا تقریبا درست باشد. درستی تقریبی یک فرضیه با کمک پارامتر دقت (ϵ) بیان می‌گردد و به عبارتی دیگر یک فرضیه h تقریبا درست است اگر شرط $\mathbf{R}(h) \leq \epsilon$ برقرار باشد و احتمال درستی با کمک پارامتر اطمینان ($1 - \delta$) بیان می‌گردد و به عبارتی فرضیه h را احتمالا تقریبا درست می‌گویند اگر شرط $\mathbb{P}[\mathbf{R}(h) \leq \epsilon] \geq 1 - \delta$ برقرار باشد براساس این تعاریف، مدل یادگیری احتمالا تقریبا درست به صورت زیر تعریف می‌شود.

¹Approximately correct

تعریف ۸.۲ قابلیت یادگیری احتمالا تقریبا درست

کلاس مفاهیم C قابلیت یادگیری احتمالا تقریبا درست را دارد اگر یک الگوریتم A و یک تابع چند جمله ای p وجود داشته باشد به گونه ای که برای هر $\epsilon \in (0, \frac{1}{8})$ و هر $\delta \in (0, \frac{1}{8})$ برای هر توزیع D روی X و برای تمامی مفاهیم $c \in C$ و برای همه مجموعه های آموزشی S با انداز $m \geq m_H(\epsilon, \delta) = p(\frac{1}{\epsilon}, \frac{1}{\delta}, |x|, |C|)$ بصورت مستقل از توزیع D نمونه برداری شده اند رابطه زیر برقرار باشد

$$\mathbb{P}_{S \sim D^m} [\mathbf{R}(h) \leq \epsilon] \geq 1 - \delta. \quad (9.2)$$

اگر الگوریتم یادگیری A در زمان $p(\frac{1}{\epsilon}, \frac{1}{\delta}, |x|, |C|)$ اجرا شود آنگاه میگوییم که کلاس مفاهیم C به صورت کارا قابل یادگیری است و اگر الگوریتم یادگیری L وجود داشته باشد به آن الگوریتم یادگیری احتمالا تقریبا درست برای کلاس مفاهیم C میگویند. در این مدل اگر فرض کنیم که یک $h \in H$ وجود دارد به گونه ای که خطای واقعی آن صفر است ($\mathbf{R}(h) = 0$) یعنی خطای خطای تجربی آن که به صورت زیر تعریف میشود نیز صفر است.

$$\hat{\mathbf{R}}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}[h(x_i) \neq c(x_i)]. \quad (10.2)$$

مدل یادگیری احتمالا تقریبا درست دارای ویژگی های بسیار جالبی است که عبارتند از (۱) این چارچوب مستقل از توزیع نمونه برداری D است. (۲) نمونه های آموزشی و آزمایشی که خطا را مشخص می کنند از یک توزیع نمونه برداری می شوند. (۳) این چارچوب برای یادگیری یک کلاس مفاهیم C است نه یک مفهوم $c \in C$. اگر کلاس مفاهیم C قابلیت یادگیری احتمالا تقریبا درست را داشته باشد آنگاه (۱) تعداد نمونه های آموزشی مورد نیاز جهت تولید فرضیه h یک تابع چند جمله ای براساس پارامترهایش است و (۲) اگر الگوریتم یادگیری برای هر نمونه زمان ثابتی را مصرف کند آنگاه زمان پردازش برای هر مجموعه آموزشی یک تابع چند جمله ای است. در این مدل سه پارامتر ϵ, δ, m وجود دارد و دو پارامتر از این سه پارامتر را مشخص می کنیم و در مسایل یادگیری مختلف به دنبال یافتن پارامتر سوم هستیم. اگر ϵ و δ مشخص شده باشند ما به دنبال پیدا کردن اندازه مجموعه آموزشی هستیم که بتواند یک یادگیر احتمالا تقریبا درست با پارامترهای مشخص شده را تولید نماید که در تعریف زیر آمده است.

تعریف ۹.۲ پیچیدگی نمونه ای

تابع $m_H : (0, 1)^2 \rightarrow \mathbb{N}$ را پیچیدگی نمونه ای کلاس مفاهیم H می گویند. به بیانی دیگر، به تعداد نمونه های آموزشی مورد نیاز برای یک یادگیری موفق پیچیدگی نمونه ای می گویند.

در ادامه این بخش چند نمونه از مسایلی که قابلیت یادگیری احتمالا تقریبا درست را دارند بررسی میشوند. نخست مجموعه توابع آستانه یک بعدی^۲ که نسخه انتزاعی بسیاری از دسته بندهای واقعی است بیان میشود.

مثال ۷.۲ (تابع آستانه یک بعدی).

در مجموعه توابع آستانه یک بعدی، فرض کنید $X = \mathbb{R}$ و کلاس مفاهیم C مجموعه ای مفاهیم c بصورت زیر باشد

$$c(x) = \begin{cases} 1 & x \geq \theta \\ 0 & x < \theta \end{cases} \quad (11.2)$$

²One dimensional threshold function

برای هر نقطه x_0 ، تابع آستانه ناحیه‌ای از \mathbb{R} را که در بازه $[x_0, \infty)$ باشد را برچسب مثبت زده و نقاطی که در بازه $(-\infty, x_0)$ باشند را برچسب منفی می‌زند. یا به عبارتی دیگر فرضیه یادگرفته شده به صورت زیر نمونه‌ها را دسته بندی می‌کند.

$$h(x) = \begin{cases} 1 & x \geq x_0 \\ 0 & x < x_0 \end{cases} \quad (12.2)$$

نقطه x_0 هر نقطه‌ای می‌تواند باشد که نقاط مثبت و منفی را از هم جدا نماید. فرضیه h را به گونه‌ای انتخاب می‌کنیم که با داده‌های آموزشی سازگار باشد. برای این منظور از الگوریتم یادگیری ۱ استفاده می‌کنیم که فرضیه h را به صورت زیر انتخاب می‌کند. در این فرضیه x_0 برابر سمت چپ ترین نمونه مثبت است..

الگوریتم ۱ الگوریتم یادگیری تابع آستانه یک بعدی

```

1: procedure OTF(S)
2:    $x_0 \leftarrow \infty$ 
3:   for  $i \leftarrow 1$  to  $m$  do
4:     if  $(c(x_i) = 1) \ \& \ (x_i < x_0)$  then
5:        $x_0 \leftarrow x_i$ 
6:     end if
7:   end for
8:   return  $h = x_0$ 
9: end procedure

```

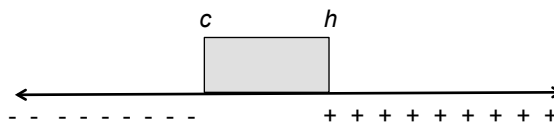
حال نشان می‌دهیم که این کلاس مفاهیم قابلیت یادگیری احتمالا تقریبا درست را دارد.

قضیه ۲.۲ (قابلیت یادگیری احتمالا تقریبا درست توابع آستانه یک بعدی).

الگوریتم ۱، قابلیت یادگیری احتمالا تقریبا درست را دارد و پیچیدگی نمونه‌ای آن برابر است با

$$m \geq \frac{1}{\epsilon} \ln \frac{1}{\delta}. \quad (13.2)$$

برهان برای اثبات، نخست حالت کلی تر را در نظر می‌گیریم و سپس برای الگوریتم یاد شده قضیه را اثبات می‌نماییم. برای اثبات شکل زیر را در نظر بگیرید.



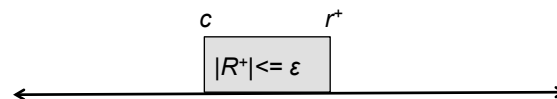
شکل ۱.۲ ناحیه خطا برای تابع آستانه یک بعدی

در این حالت یک ناحیه مانند شکل ۱.۲ مانند شکل بالا بین c و h باقی می‌ماند که فرضیه h نقاط درون این ناحیه را به اشتباه برچسب می‌زند. هدف این است که اندازه این ناحیه کوچکتر یا مساوی ϵ باشد. اندازه به معنای طول بین c و h نیست بلکه جرم تابع توزیع در این بازه است. ما می‌خواهیم تضمین نماییم که برای هر توزیع D و همه نقاط ممکن، احتمال قرار گرفتن نقاط در این بازه بیش از ϵ نباشد که دو حالت زیر پیش می‌آید.

۱. فرضیه h در فاصله‌ای بیشتر از ϵ در سمت راست c قرار گیرد به این حالت B^+ می‌گوییم.

۲. فرضیه h در فاصله‌ای بیشتر از ϵ در سمت چپ c قرار گیرد به این حالت B^- می‌گوییم.

نخست حالت B^+ را در نظر بگیرید. بر اساس شکل زیر فرض کنید R^+ مشخص کننده کوچکترین ناحیه‌ای باشد که c در سمت چپ آن و جرم احتمال در آن دست کم برابر ϵ باشد یعنی $R^+ = [c, r^+]$ که $r^+ = \sup\{r \geq c \mid \mathbb{P}[c, r] < \epsilon\}$. با توجه به الگوریتم ارایه شده، اگر همه نمونه‌ها در خارج از بازه R^+ قرار داشته باشند h در سمت راست ناحیه r^+ قرار می‌گیرد و هیچ نمونه آموزشی در بازه R^+ قرار نمی‌گیرد.



شکل ۲.۲ تعیین کران برای ناحیه خطا برای تابع آستانه یک بعدی

در حالتی که $\mathbb{P}[x \in R^+] < \epsilon$ باشد، مساله حل است و اگر بازه R^+ دارای اندازه دست کم ϵ باشد داریم $\mathbb{P}[x \in R^+] \geq \epsilon$ و به طور معادل رابطه $\mathbb{P}[x \notin R^+] \leq 1 - \epsilon$ نتیجه می‌شود. برای m نمونه آموزشی می‌توان نوشت

$$\mathbb{P}[B^+] = \mathbb{P}[(x_1 \notin \mathbb{R}^+) \wedge \dots \wedge (x_m \notin \mathbb{R}^+)] \quad (14.2)$$

$$= \mathbb{P}[x_1 \notin \mathbb{R}^+] \dots \mathbb{P}[x_m \notin \mathbb{R}^+] \quad (15.2)$$

$$\leq (1 - \epsilon)^m. \quad (16.2)$$

خط دوم رابطه بالا بر این اساس نوشته شده است که نمونه‌ها مستقل از هم و با توزیع یکسان نمونه برداری شده‌اند. تحلیل بالا را می‌توان در باره ناحیه B^- نیز نوشت. بنابراین احتمال اینکه خطایی بیشتر از ϵ داشته باشیم را می‌توانیم به صورت زیر بدست آوریم.

$$\mathbb{P}[\mathbf{R}(h) > \epsilon] \leq \mathbb{P}[B^+ \cup B^-] \quad (17.2)$$

$$\leq \mathbb{P}[B^+] + \mathbb{P}[B^-] \quad (18.2)$$

$$\leq 2(1 - \epsilon)^m \quad (19.2)$$

$$\leq 2e^{-\epsilon m}. \quad (20.2)$$

خط دوم رابطه بالا بر اساس کران اجتماعات و خط سوم بر اساس تقارن نوشته شده‌اند و خط آخر بر اساس نابرابری $1 + x \leq e^x$ نوشته شده است. هدف این است که این احتمال از δ کوچکتر باشد و یا به عبارتی دیگر داشته باشیم

$$\mathbb{P}[\mathbf{R}(h) > \epsilon] \leq 2e^{-\epsilon m} \leq \delta. \quad (21.2)$$

حال اگر نابرابری بالا را بر اساس m حل نماییم خواهیم داشت

$$m \geq \frac{1}{\epsilon} \ln \frac{2}{\delta}. \quad (22.2)$$

با توجه به الگوریتم یاد شده، هیچ گاه حالت B^- رخ نمی‌دهد و بنابراین ضریب ۲ در رابطه بالا حذف شده و قضیه اثبات می‌گردد. ■

در عمل مجموعه آموزشی که تعداد نمونه ها را مشخص می نماید در اختیار الگوریتم قرار داده می شود و هدف این است که کران خطای واقعی را تخمین بزنیم. با استفاده از قضیه بالا و با محاسبات ساده به سادگی می توان نشان داد که با احتمال دست کم $1 - \delta$ نابرابری زیر برقرار است.

$$\mathbf{R}(h) \leq \frac{1}{m} \ln \frac{2}{\delta}. \quad (23.2)$$

برای بدست آوردن معادله بالا، کافی است معادله (22.2) را بر اساس ϵ بازنویسی نماییم که به نتیجه زیر می رسیم.

$$\epsilon \geq \frac{1}{m} \ln \frac{2}{\delta}. \quad (24.2)$$

و از تعریف یادگیری احتمالا تقریبا درست، نابرابری $\mathbf{R}(h) \leq \epsilon$ را داریم. لذا کوچکترین مقداری که برای ϵ می توان متصور شد برابر است با $\frac{1}{m} \ln \frac{2}{\delta}$ و در نتیجه خواهیم داشت

$$\mathbf{R}(h) \leq \frac{1}{m} \ln \frac{2}{\delta}. \quad (25.2)$$

حال کلاس فرضیه را از توابع آستانه یک بعدی به بازه گسترش می دهیم.

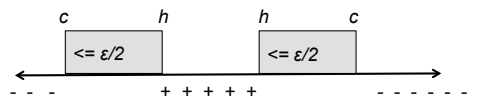
مثال ۸.۲ (یادگیری بازه).

در مجموعه توابع بازه، فرض کنید $X = \mathbb{R}$ و کلاس مفاهیم C مجموعه ای مفاهیم c بصورت زیر باشد

$$c(x) = \begin{cases} 1 & x \in [a, b] \\ 0 & x \notin [a, b] \end{cases} \quad (26.2)$$

این مساله را با کمک توابع آستانه یک بعدی حل می کنیم. در این مساله دو ناحیه مرزی خطا مطابق شکل ۳.۲ داریم و می خواهیم که اندازه هر کدام حداکثر $\frac{\epsilon}{2}$ باشد و بنابراین چهار حالت مختلف داریم. اگر از پاسخ مساله قبل استفاده کنیم خواهیم داشت

$$m \geq \frac{2}{\epsilon} \ln \frac{4}{\delta}. \quad (27.2)$$



شکل ۳.۲ تعیین کران برای ناحیه خطا برای تابع بازه

مثال ۹.۲ (پیچیدگی نمونه ای مستطیل های با اضلاع موازی محورها).

چند نمونه برای یادگیری موفق مساله مستطیل های با اضلاع موازی محورها مورد نیاز است؟ در این مثال، فضای فرضیه ها را مستطیل هایی که اضلاع آنها موازی محورهای مختصات باشد در نظر می گیریم و الگوریتم یادگیری تلاش می کند که کوچکترین مستطیل سازگار با داده های آموزشی را پیدا کند. بنابراین می بایست $\hat{\mathbf{R}}(h) = 0$ فرض کنید c مفهوم مستطیل باشد همانگونه که در شکل ۴.۲ نشان داده شده است فرضیه تولید شده که با h نشان داده می شود درون مستطیل c قرار دارد. نوارهای بین مستطیل های c و h خطای واقعی فرضیه تولید شده h را نشان می دهند. هدف پیدا کردن تعداد نمونه های لازم است به گونه ای که احتمال قرار گرفتن یک نمونه مثبت در این چهار نوار کوچکتر یا مساوی ϵ باشد. اگر برای هر کدام از این

نوارها بتوانیم تضمین نماییم که کران بالای این احتمال $\frac{\epsilon}{4}$ است لذا مجموع این احتمال ها حداکثر برابر خواهد بود با $\epsilon = 4(\frac{\epsilon}{4})$. دقت شود که چهار گوشه نوارها را دو بار محاسبه نموده ایم لذا خطا از $\epsilon = 4(\frac{\epsilon}{4})$ کمتر خواهد بود. احتمال اینکه یک نمونه تصادفی در یک نوار نباشد برابر است با $1 - \frac{\epsilon}{4}$ و احتمال اینکه m نمونه در این نوار نباشند برابر است با $(1 - \frac{\epsilon}{4})^m$ و احتمال اینکه این m نمونه در هیچکدام از این چهار نوار نباشند حداکثر برابر است با $4(1 - \frac{\epsilon}{4})^m$ که میبایست حداکثر برابر باشد با δ . بنابراین داریم.

$$4 \left(1 - \frac{\epsilon}{4}\right)^m \leq \delta$$

$$\left(1 - \frac{\epsilon}{4}\right)^m \leq \frac{\delta}{4}$$

با استفاده از نابرابری $1 - x \leq e^{-x}$ خواهیم داشت.

$$\left(1 - \frac{\epsilon}{4}\right)^m \leq \exp\left[-\frac{\epsilon m}{4}\right] \leq \frac{\delta}{4}$$

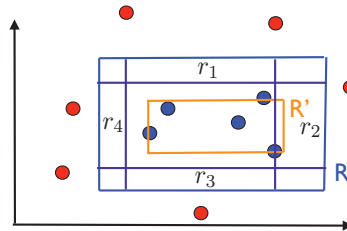
اگر از طرفین نابرابری بالا لگاریتم گرفته شود خواهیم داشت

$$-\frac{\epsilon m}{4} \leq \log \frac{\delta}{4}$$

بنابراین کمترین تعداد نمونه های مورد نیاز برای یادگیری موفق این مساله برابر است با

$$m \geq \frac{4}{\epsilon} \log \frac{4}{\delta}$$

در نتیجه مساله یادگیری مستطیل های با اضلاع موازی محورها قابلیت یادگیری احتمالا تقریباً درست را دارا است.



شکل ۴.۲ خطای بین فرضیه تولید شده و مفهوم

پیچیدگی نمونه ای برای فضای فرضیه متناهی

۱.۳.۲

در مثال های یاد شده فرضیه ای که توسط الگوریتم تولید می شود با داده های آموزشی سازگار بودند. در این بخش کران پیچیدگی نمونه ای و به صورت معادل کران خطای واقعی را برای فرضیه های سازگار و حالتی که $|H|$ متناهی باشد محاسبه می کنیم. چون هدف پیدا کردن فرضیه سازگار با داده های آموزشی است لذا فرض می کنیم که فضای فرضیه H حاوی مفهوم c باشد یا به عبارتی دیگر شرط $c \in H$ برقرار باشد. قضیه زیر این کران را مستقل از الگوریتم یادگیری محاسبه می نماید.

قضیه ۳.۲ (پیچیدگی نمونه‌ای فضای فرضیه متناهی).

اگر H یک مجموعه متناهی از فرضیه ها و مجموعه $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$ مجموعه آموزشی باشد که از توزیع D و بطور مستقل از هم نمونه برداری شده باشند و توسط مفهوم $c \in \mathcal{C}$ برچسب گذاری شده باشند. آنگاه برای هر $0 < \delta < \frac{1}{4}$ و $0 < \epsilon < \frac{1}{4}$ ، اگر الگوریتم یادگیری A فرضیه سازگار $h \in H$ را تولید نماید نابرابری‌های زیر با احتمال دست کم $1 - \delta$ برقرار هستند..

$$m \geq \frac{1}{\epsilon} \left[\log |H| + \log \frac{1}{\delta} \right] \quad (28.2)$$

$$\mathbf{R}(h) \leq \frac{1}{m} \left[\log |H| + \log \frac{1}{\delta} \right] \quad (29.2)$$

در رابطه‌های بالا، $\log |H|$ را می‌توان اندازه بازنمایی فرضیه های فضای H یعنی تعداد بیت های لازم برای بازنمایی فرضیه های این فضا دانست. در واقع این معیار نشان دهنده غنای فضای فرضیه است به گونه‌ای که هرچه این فضا بزرگتر باشد احتمالاً انعطاف پذیرتر است و می‌تواند فرضیه نزدیک به مفهوم را در خود داشته باشد. در ادامه نشان می‌دهیم که این معیار الزماً به خوبی غنای فضای فرضیه را نشان نمی‌دهد.

برهان نخست نابرابری را (۲۸.۲) بدست می‌آوریم و سپس نابرابری (۲۹.۲) را. احتمال اینکه فرضیه سازگاری وجود داشته باشد که خطای واقعی آن بیشتر از ϵ برابر است با

$$\begin{aligned} \mathbb{P} \left[\exists h \in H \mid \hat{\mathbf{R}}(h) = \circ \wedge \mathbf{R}(h) > \epsilon \right] &= \mathbb{P} \left[(\hat{\mathbf{R}}(h_1) = \circ \wedge \mathbf{R}(h_1) > \epsilon) \vee (\hat{\mathbf{R}}(h_2) = \circ \wedge \mathbf{R}(h_2) > \epsilon) \vee \dots \right] \\ &\leq \sum_{h_i \in H} \mathbb{P} \left[\hat{\mathbf{R}}(h_i) = \circ \wedge \mathbf{R}(h_i) > \epsilon \right] \end{aligned} \quad (30.2)$$

$$\leq \sum_{h_i \in H} \mathbb{P} \left[\hat{\mathbf{R}}(h_i) = \circ \mid \mathbf{R}(h_i) > \epsilon \right]. \quad (31.2)$$

نابرابری (۳۰.۲) با استفاده از نابرابری کران اجتماعات بدست آمده است

$$\mathbb{P} \left[\bigcup_{i=1}^K A_i \right] \leq \sum_{i=1}^K \mathbb{P} [A_i] \quad (32.2)$$

و نابرابری (۳۱.۲) نیز با استفاده از قاعده احتمال شرطی بدست آمده است. حال هر $h \in H$ را در نظر بگیرید که $\mathbf{R}(h) > \epsilon$ باشد آنگاه احتمال اینکه h با داده‌های آموزشی سازگار باشد را می‌توان به صورت زیر محاسبه نمود.

$$\mathbb{P} \left[\hat{\mathbf{R}}(h) = \circ \mid \mathbf{R}(h) > \epsilon \right] \leq (1 - \epsilon)^m \quad (33.2)$$

نابرابری بالا از روش زیر بدست آمده است. فرض کنید که $\mathbf{R}(h) > \epsilon$ باشد. احتمال اینکه h یک نمونه را اشتباه دسته بندی نماید بزرگتر از ϵ است و در نتیجه احتمال اینکه h یک نمونه را درست دسته بندی نماید کوچکتر یا مساوی $1 - \epsilon$ خواهد بود و احتمال اینکه h همه m نمونه را درست دسته بندی نماید کوچکتر یا مساوی $(1 - \epsilon)^m$ است. باقرار دادن رابطه بالا در نابرابری (۳۱.۲) خواهیم داشت.

$$\sum_{h_i \in H} \mathbb{P} \left[\hat{\mathbf{R}}(h_i) = \circ \mid \mathbf{R}(h_i) > \epsilon \right] \leq |H| (1 - \epsilon)^m \leq \delta \quad (34.2)$$

با استفاده از نابرابری $1 - x \leq e^{-x}$ و نابرابری بالا خواهیم داشت.

$$|H| (1 - \epsilon)^m \leq |H| \exp(-\epsilon m) \leq \delta \quad (35.2)$$

با ساده سازی این نابرابری، نابرابری زیر بدست می‌آید.

$$m \geq \frac{1}{\epsilon} \left[\log|H| + \log \frac{1}{\delta} \right]. \quad (۳۶.۲)$$

برای بدست آوردن نابرابری (۲۹.۲) کافی است نابرابری (۲۸.۲) را براساس ϵ حل نماییم، نابرابری زیر بدست می‌آید.

$$\epsilon \geq \frac{1}{m} \left[\log|H| + \log \frac{1}{\delta} \right]. \quad (۳۷.۲)$$

با توجه به اینکه کران بالای خطای واقعی ϵ است لذا همیشه خطای واقعی از کمترین مقدار ϵ کمتر خواهد بود و در نتیجه نابرابری (۲۹.۲) بدست می‌آید و قضیه اثبات می‌گردد.

قضیه بالا بیان می‌دارد که هنگامی که H متناهی است هر الگوریتم سازگار A یک الگوریتم یادگیری احتمالا تقریبا درست است. هزینه ای که برای الگوریتم‌های سازگار می‌پردازیم این است که از مجموعه H های بزرگتری فرضیه را انتخاب می‌کنیم که مفهوم نهایی c در آن است یعنی $c \in H$. یکی از مشکلات نابرابری (۳۵.۲) این است که اگر $|H|$ خیلی بزرگ باشد ممکن است سمت چپ این نابرابری بیشتر از یک باشد. یکی از اشکالات این کران این است که این کران برای $|H|$ های بزرگ سفت نیست. یعنی اگر $|H| = \infty$ باشد این کران اصلا مناسب نیست. حال با توجه به قضیه بالا، تعداد نمونه های لازم برای چند مساله یادگیر را بدست می‌آوریم.

مثال ۱۰.۲ (یادگیری ترکیب عطفی از n متغیر).

فرض کنید H فضای فرضیه‌ای باشد که از ترکیب عطفی چند متغیر (ویژگی) ساخته شود. برای یادگیری این مساله به چند نمونه نیاز است؟ در این مساله به جای هر متغیر می‌توان آن متغیر، نقیض آن و یا ؟ را قرار داد در نتیجه $|H| = 3^n$ خواهد بود که n تعداد ویژگی‌های مساله است. بنابراین با استفاده از نابرابری (۲۸.۲) داریم

$$\begin{aligned} m &\geq \frac{1}{\epsilon} \left[\log|H| + \log \frac{1}{\delta} \right] \\ &= \frac{1}{\epsilon} \left[\log 3^n + \log \frac{1}{\delta} \right] \\ &= \frac{1}{\epsilon} \left[n \log 3 + \log \frac{1}{\delta} \right]. \end{aligned}$$

چون پیچیدگی نمونه‌ای برای یادگیری این مساله یک تابع چندجمله‌ای از پارامترهایش است پس این مساله قابلیت یادگیری احتمالا تقریبا درست را دارد.

در ادامه مساله یادگیری تابع دودویی را بررسی می‌کنیم.

مثال ۱۱.۲ (یادگیری از فضای همه توابع دودویی با n متغیر).

فرض کنید H فضای فرضیه همه توابع دودویی باشد که از n متغیر ساخته شود. برای یادگیری این مساله به چند نمونه نیاز است؟ در این مساله فضای فرضیه دارای $|H| = 2^{2^n}$ فرضیه خواهد بود. بنابراین با استفاده از نابرابری (۲۸.۲) داریم

$$\begin{aligned} m &\geq \frac{1}{\epsilon} \left[\log|H| + \log \frac{1}{\delta} \right] \\ &= \frac{1}{\epsilon} \left[\log 2^{2^n} + \log \frac{1}{\delta} \right] \\ &= \frac{1}{\epsilon} \left[2^n \log 2 + \log \frac{1}{\delta} \right]. \end{aligned}$$

با توجه به اینکه پیچیدگی نمونه ای چند جمله ای بر اساس n نیست لذا قابلیت یادگیری احتمالا تقریبا درست را ندارد.

اگر n را یک عدد ثابت و کوچک در نظر بگیریم میتوانیم ادعا کنیم که این کلاس مفاهیم قابلیت یادگیری احتمالا تقریبا درست را دارد اما همچنان پیچیدگی نمونه ای با n به صورت نمایی رشد می کند که با تعریف یادگیری احتمالا تقریبا درست در تناقض است. برای برخی از توابع دودویی، تعداد عبارت های داخل فرضیه 2^n است که سبب نمایی شدن پیچیدگی نمونه ای می شود. اما برای بسیاری از مسایل ساده تر، بازنمایی بسیار ساده تر است و نیاز به پیچیدگی نمونه ای نمایی ندارد. از آنجایی که میبایست تعداد نمونه ها برای تمام مفاهیم $c \in C$ یک تابع چندجمله ای براساس پارامترهایش باشد کران قضیه برای مثال بالا نتیجه نمی دهد که این مساله قابلیت یادگیری احتمالا تقریبا درست را دارد اما این قضیه بیان نمی کند که این مساله قابلیت یادگیری احتمالا تقریبا درست را ندارد. زیرا کران مشخص شده در قضیه، یک کران سفت نیست و ما نمی توانیم از این قضیه برای نشان دادن قابلیت یادگیری احتمالا تقریبا درست این مساله استفاده کنیم. پرسش اینکه آیا این مساله قابلیت یادگیری احتمالا تقریبا درست را دارد یا نه برای بیش از دو دهه، یک مساله باز است.

قضیه ۴.۲ (قابلیت یادگیری احتمالا تقریبا درست فضای فرضیه متناهی).
هر فضای فرضیه متناهی بالقوه قابلیت یادگیری احتمالا تقریبا درست را دارد.

۴.۲ مدل یادگیری احتمالا تقریبا درست بدون پیش فرض

پیش از این فرض کردیم که $C \subset H$ است و الگوریتم یادگیری یک فرضیه سازگار را تولید میکند. براین اساس نابرابری (۳۵.۲) مهم است زیرا تعداد نمونه های کافی برای یک یادگیری موفق در حالتی که خطای آموزشی صفر است را مشخص مینماید. در این نابرابری فرض شده است که $\hat{\mathbf{R}}(h) = 0$ باشد. اما در بیشتر کاربردها به دلایل زیر ممکن است فرضیه تولید شده سازگار نباشد و در نتیجه خطای آموزشی بیش از صفر داریم. الف) فضای فرضیه H به اندازه کافی غنی نیست که توانایی نشان دادن مفهوم c را داشته باشد یا به عبارتی دیگر $c \notin H$. ب) از نظر محاسباتی پیدا کردن فرضیه h سازگار امکان پذیر نباشد و ج) ممکن است c قطعی وجود نداشته باشد و یک تابع تصادفی باشد. در صورتیکه H شامل c نباشد ممکن است خطای آموزش صفر نباشد. در نتیجه در این مسایل، خطای آموزشی بیشتر از صفر داریم. البته فرضیه بدست آمده در صورتی مناسب است که دارای خطای آموزش و واقعی کم باشد. علاوه بر مشکل یاد شده، ممکن است داده ها نیز نویزی باشند و نتوان یک فرضیه سازگار با داده های آموزشی پیدا کرد. لذا در این حالت نیز خطای آموزشی بزرگتر از صفر داریم. با توجه به دلایلی که گفته شد ما نیاز داریم تا تعاریف و فرضیات مدل های پیشین را بازنگری کنیم. پیش از این فرض کردیم که نمونه ها از توزیع D نمونه برداری می شوند و سپس با تابع معین c که نقش آموزگار را دارد برچسب زده می شوند. بنابراین برچسب ها به صورت احتمالاتی مشخص نمی شوند. در این جا فرض می کنیم که یک تابع توزیع احتمال D روی زوج نمونه های (x, y) وجود دارد که نمونه ها از این توزیع نمونه برداری می شوند. بنابراین خطای واقعی به صورت زیر اصلاح می شود.

$$\begin{aligned} \mathbf{R}(h) &= \mathbb{P}_{(x,y) \sim D} [h(x) \neq y] \\ &= \sum_{h(x) \neq y} D(x, y). \end{aligned} \quad (38.2)$$

هدف این است که فرضیه h پیدا کنیم که با احتمال بالا دارای خطای واقعی پایینی باشد. اما در شرایط بیان شده ممکن است چنین فرضیه ای وجود نداشته باشد و ما می توانیم فرضیه ای را پیدا کنیم که نزدیک به بهترین فرضیه درون H باشد. در نتیجه هدف پیدا کردن فرضیه h

است که خطای واقعی آن نزدیک به $\min_h \mathbf{R}(h)$ باشد. همچنین مدل های پیشین فرض می‌کردند که الگوریتمی وجود دارد که توانایی پیدا کردن فرضیه سازگار با داده‌های آموزشی را دارد. این فرض در این حالت دیگر وجود ندارد. بنابراین الگوریتم‌هایی را بررسی میکنیم که فرضیه $h \in H$ را پیدا می‌کنند که بیشترین سازگاری را با داده‌ها داشته باشد. این سازگاری بر اساس خطای تجربی روی داده‌های آموزشی تعریف می‌شود. یک روش این است که در این مسایل ما به دنبال فرضیه‌ای باشیم که خطای آموزش را کمینه نماید. بیشتر الگوریتم‌های یادگیری ماشین که هم اکنون به کار می‌روند از این گونه هستند. به این الگوریتم‌ها، الگوریتم‌های کمینه سازی خطای تجربی^۳ می‌گویند. دلیل به کارگیری این الگوریتم‌ها این است که کمینه سازی خطای تجربی تقریب خوبی به کمینه سازی خطای واقعی است. دلایل دیگر برای بکارگیری این الگوریتم‌ها در یادگیری ماشین را در بخش‌های بعدی بیان خواهیم نمود. در واقع الگوریتم‌های کمینه‌سازی خطای تجربی فرضیه h_{erm} را برای مجموعه آموزشی S به صورت زیر انتخاب میکنند.

$$h_{erm} = \underset{h \in H}{\operatorname{argmin}} \hat{\mathbf{R}}(h). \quad (۳۹.۲)$$

در الگوریتم‌های کمینه سازی خطای تجربی، دو نکته وجود دارد که باید به درستی انتخاب شوند. نخست آنکه فضای فرضیه H به درستی انتخاب شود و دوم آنکه شیوه کمینه سازی خطای تجربی که بر اساس تابع هزینه (تابع خطا) تعیین می‌گردد به درستی انتخاب گردد. حال برای روشن شدن موضوع یک مثال بیان می‌شود.

مثال ۱۲.۲ (رگرسیون خطی).

مساله رگرسیون خطی را در نظر بگیرید که مجموعه آموزشی به صورت $S = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$ باشد که $y_i \in \mathbb{R}$ و $\mathbf{x}_i \in \mathbb{R}^n$ و نمونه‌های مجموعه آموزشی بر اساس توزیع \mathcal{D} و به صورت مستقل نمونه برداری شده باشند. همچنین فرض کنید که \mathbf{A} ماتریسی باشد که سطرهایش مقادیر \mathbf{x}_i ها و مجموعه فرضیه‌ها برابر $H = \{h_w : \mathcal{X} \mapsto \mathcal{Y} | h_w(x) = \langle \mathbf{x}, \mathbf{w} \rangle\}$ باشد. اگر تابع خطای تجربی (تابع هزینه) را به صورت زیر تعریف کنیم

$$\hat{\mathbf{R}}(h) = \frac{1}{m} \sum_{i=1}^m [h(\mathbf{x}_i) - y_i]^2, \quad (۴۰.۲)$$

یک الگوریتم کمینه سازی خطای تجربی، فرضیه زیر را تولید میکند.

$$\mathbf{w}^* = \underset{\mathbf{w} \in \mathbb{R}^n}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m [h(\mathbf{x}_i) - y_i]^2 \quad (۴۱.۲)$$

$$= \underset{\mathbf{w} \in \mathbb{R}^n}{\operatorname{argmin}} \|\mathbf{A}\mathbf{w} - \mathbf{y}\|^2 \quad (۴۲.۲)$$

که \mathbf{y} برداری است با مقادیر y_i ها پر شده است. با حل این معادله فرضیه زیر بدست می‌آید.

$$\mathbf{w}^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}. \quad (۴۳.۲)$$

در واقع در این مساله، الگوریتم کمینه سازی خطای تجربی یک پاسخ بسته برای ما تولید مینماید اما همیشه چنین نیست.

در الگوریتم‌های کمینه سازی خطای تجربی یک پرسش مهم وجود دارد که بهترین فرضیه یعنی فرضیه‌ای که کمترین خطای واقعی را داشته باشد کدام است و مقدار خطای واقعی این فرضیه چه اندازه است؟ بهترین فرضیه را به صورت h_{opt} و خطای واقعی آن را به صورت

³Empirical risk minimization algorithms (ERM)

\mathbf{R}^* نشان می‌دهیم و برابر است با

$$\mathbf{R}^* = \inf_{h \in H} \mathbf{R}(h). \quad (44.2)$$

بهترین مقدار خطای واقعی که به خطای بیز معروف است با دسته بند زیر که به دسته بند بهینه بیز معروف است قابل حصول است

$$h_{opt} = \begin{cases} 1 & \text{if } \mathbb{P}[y = 1 | x] \geq \frac{1}{2} \\ 0 & \text{if } \mathbb{P}[y = 1 | x] < \frac{1}{2} \end{cases} \quad (45.2)$$

با توجه به اینکه یادگیر به توزیع D دسترسی ندارد و آن را نمی‌شناسد برای دسته بند بهینه بیز نمی‌توان مقدار خطای بیز را محاسبه نمود و بر اساس قضیه زیر می‌توان نشان داد که هیچ دسته بندی خطای واقعی کمتری نسبت به دسته بند بهینه بیز تولید نمی‌کند.

قضیه ۵.۲ (خطای دسته بند بهینه بیز).

خطای دسته بند بهینه بیز برابر \mathbf{R}^* است. به عبارتی دیگر هیچ دسته بندی وجود ندارد که خطایی کمتر از خطای دسته بند بهینه بیز را تولید نماید.

برهان اثبات این قضیه را به عنوان تمرین انجام دهید.

هدف همه روش های یادگیری ماشین این است که $\min_{h \in H} \mathbf{R}(h)$ را پیدا نمایند. برای این منظور برای مجموعه آموزشی S داده شده، خطای تجربی را کمینه می‌کنیم. برای اینکه بتوان از این رویکرد استفاده نمود باید بتوانیم نشان دهیم که کمینه‌سازی خطای تجربی سبب کمینه‌سازی خطای واقعی می‌شود یا دست کم می‌توان کران خطای واقعی را بر اساس خطای تجربی پیش‌بینی نمود. به بیانی دیگر باید بگوییم که برای همه فرضیه‌های $h \in H$ با احتمال $1 - \delta$ نابرابری زیر برقرار است.

$$|\mathbf{R}(h) - \hat{\mathbf{R}}(h)| \leq \epsilon \quad (46.2)$$

یا به بیانی دیگر کران پایین m را باید به گونه‌ای پیدا کنیم که برای همه $\epsilon, \delta \in (0, 1)$ نابرابری زیر برقرار باشد.

$$\mathbb{P} \left[\left| \mathbf{R}(h_{erm}) - \min_{h \in H} \mathbf{R}(h) \right| > \epsilon \right] < \delta \quad (47.2)$$

که h_{erm} فرضیه‌ای است که الگوریتم کمینه سازی خطای تجربی آن را تولید می‌نماید. نابرابری (46.2) این نتیجه را می‌دهد که اگر h_{erm} خطای تجربی را کمینه نماید آنگاه برای همه $h \in H$ داریم

$$\mathbf{R}(h_{erm}) \leq \hat{\mathbf{R}}(h_{erm}) + \epsilon \quad (48.2)$$

$$\leq \hat{\mathbf{R}}(h) + \epsilon \quad (49.2)$$

$$\leq (\mathbf{R}(h) + \epsilon) + \epsilon \quad (50.2)$$

$$= \mathbf{R}(h) + 2\epsilon. \quad (51.2)$$

بدلیل اینکه h_{erm} خطای تجربی را کمینه می‌کند رابطه (49.2) بدست آمده است. رابطه (50.2) به صورت زیر بدست آمده است. چون

h_{erm} خطای تجربی را کمینه می‌نماید از نابرابری (46.2) داریم

$$-\epsilon \leq \mathbf{R}(h) - \hat{\mathbf{R}}(h) \leq \epsilon \quad (52.2)$$

با ساده سازی نابرابری بالا به نابرابری زیر می‌رسیم.

$$\hat{\mathbf{R}}(h) \leq \mathbf{R}(h) + \epsilon. \quad (۵۳.۲)$$

نابرابری بالا بیان می‌دارد که خطای واقعی فرضیه‌ای که از کمینه سازی خطای تجربی بدست می‌آید از کمینه خطای همه فرضیه‌ها بعلاوه 2ϵ بیشتر نخواهد بود. در نتیجه این فرضیه خطای واقعی نزدیک به کران پایین خطا برای همه فرضیه‌های مجموعه H را تولید می‌نماید برای ثابت نمودن این نابرابری، می‌بایست همگرایی یکنواخت را برای همه $h \in H$ ثابت نماییم یعنی نابرابری (۴۶.۲) را اثبات نماییم. این نابرابری بیان می‌کند که الگوریتم کمینه‌سازی خطا یک یادگیر احتمالا تقریبا درست بدون پیش فرض است. پیش از اینکه همگرایی یکنواخت را برای همه $h \in H$ ثابت نماییم، نخست مدل توسعه یافته مدل یادگیری احتمالا تقریبا درست را برای حالتی که فرضیه بدست آمده سازگار نباشد بیان می‌کنیم. این مدل یادگیری را مدل یادگیری بدون پیش فرض^۴ می‌گویند و به صورت زیر تعریف می‌شود زیرا فرضیه‌ای که پیش از این داشتیم را دیگر ندارد. این یادگیری که فرض نمی‌کند که $c \notin H$ و تنها فرضیه‌ای را پیدا می‌کند که خطای آموزش را کمینه نماید یادگیر بدون پیش فرض نامیده می‌شود زیرا الزامی در باره شرط $C \subset H$ ندارد.

تعریف ۱۰.۲ یادگیری احتمالا تقریبا درست بدون پیش فرض فرض کنید H یک فضای فرضیه باشد. الگوریتم یادگیری A

یک الگوریتم یادگیری احتمالا تقریبا درست بدون پیش فرض است اگر یک چند جمله‌ای p وجود داشته باشد به گونه‌ای که برای هر $\epsilon \in (0, 1)$ و $\delta \in (0, 1)$ و برای هر توزیع D روی $\mathcal{X} \times \mathcal{Y}$ برای هر مجموعه آموزشی با اندازه

$$m \geq p\left(\frac{1}{\epsilon}, \frac{1}{\delta}, |\mathbf{x}|, |C|\right) \quad (۵۴.۲)$$

نابرابری زیر برای فرضیه h_s که توسط الگوریتم یادگیری A با دیدن مجموعه آموزشی S تولید می‌شود برقرار است.

$$\mathbb{P}\left[\left|\mathbf{R}(h_s) - \min_{\hat{h} \in H} \mathbf{R}(\hat{h})\right| \leq \epsilon\right] \geq 1 - \delta \quad (۵۵.۲)$$

براساس تعریف بالا، می‌توان الگوریتم کارا را در این مدل به صورت زیر تعریف نمود.

تعریف ۱۱.۲ الگوریتم یادگیری احتمالا تقریبا درست بدون پیش فرض کارا اگر الگوریتم A که در تعریف ۱۰.۲ آمده است در

زمان چند جمله‌ای $p(\frac{1}{\epsilon}, \frac{1}{\delta}, |x|, |c|)$ اجرا شود انگاه به این الگوریتم، یک الگوریتم یادگیری احتمالا تقریبا درست بدون پیش فرض کارا می‌گویند.

همان‌گونه که بیان شد الگوریتم کمینه‌سازی خطا یک الگوریتم یادگیر احتمالا تقریبا درست بدون پیش فرض است. برای اینکه این ویژگی را نشان بدهیم باید اثبات کنیم که برای همه فرضیه‌های $h \in H$ ، با احتمال دست‌کم $1 - \delta$ روی یک مجموعه‌ای که به صورت تصادفی از فضای نمونه انتخاب شده باشد نابرابری (۵۵.۲) برقرار است. این نیازمندی‌ها توسط ویژگی همگرایی یکنواخت بیان می‌شود که در ادامه بیان می‌شود.

⁴Agnostic learning model

تعریف ۱۲.۲ همگرایی یکنواخت فضای ویژگی H با توجه به دامنه $\mathcal{X} \times \mathcal{Y}$ و تابع خطا (هزینه یا زیان) دارای ویژگی همگرایی یکنواخت است اگر یک تابع $\mathbb{N} \rightarrow (0, 1)^2 : m_H^{UC} \mapsto (\epsilon, \delta)$ وجود داشته باشد به گونه ای که برای هر $\epsilon \in (0, 1)$ و $\delta \in (0, 1)$ و برای هر توزیع \mathcal{D} روی $\mathcal{X} \times \mathcal{Y}$ ، اگر مجموعه آموزشی S با اندازه $m \geq m_H^{UC}(\epsilon, \delta)$ به صورت مستقل و با توزیع \mathcal{D} نمونه برداری شده باشند آنگاه نابرابری زیر برقرار است.

$$\mathbb{P} \left[\forall h \in H \left| \mathbf{R}(h) - \hat{\mathbf{R}}(h) \right| \leq \epsilon \right] \geq 1 - \delta$$

واژه یکنواخت در تعریف بالا به این معنی است که داشتن یک مجموعه آموزشی با اندازه مشخص که برای همه فرضیه های $h \in H$ و همه توزیع های ممکن روی دامنه کار می کند. مشابه پیچیدگی نمونه ای که پیش از این داشتیم، تابع m_H^{UC} به کمترین اندازه مجموعه آموزشی اشاره می کند که برای آن ویژگی همگرایی یکنواخت برقرار است. در فصل های بعد نشان خواهیم داد که ویژگی یکنواخت شرط لازم و کافی برای قابلیت یادگیری است.

پیش از اینکه رابطه (۴۶.۲) را اثبات کنیم نخست حالتی ساده که تنها یک فرضیه داشته باشیم در نظر می گیریم و آن را بررسی می کنیم.

قضیه ۶.۲ (کران خطا برای یادگیر بدون پیش فرض).

اگر S یک مجموعه آموزشی با m نمونه که بصورت مستقل و با توزیع یکسان از \mathcal{D} نمونه برداری شده باشد آنگاه برای هر ثابت $\delta > 0$ و هر فرضیه ثابت $h : \mathcal{X} \mapsto \{0, 1\}$ نابرابری زیر برقرار است

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \sqrt{\frac{\ln \frac{2}{\delta}}{2m}} \quad (۵۶.۲)$$

پیش از اثبات این قضیه، نخست نابرابری هافدینگ را که در اثبات قضیه به کار می رود را شرح می دهیم. فرض کنید که X_1, \dots, X_m متغیرهای تصادفی مستقل باشند که برای همه i ها شرط $X_i \in [a_i, b_i]$ برقرار باشد. آنگاه برای هر $\epsilon > 0$ ، نابرابری زیر برای $S_m = \sum_{i=1}^m X_i$ برقرار است.

$$\mathbb{P} \left[|S_m - \mathbb{E}[S_m]| \geq \epsilon \right] \leq 2 \exp \left[-2 \frac{\epsilon^2}{\sum_{i=1}^m (b_i - a_i)^2} \right] \quad (۵۷.۲)$$

برهان اثبات این قضیه با کمک نابرابری (۵۷.۲) به صورت زیر انجام می شود. برای هر $i \in \{1, 2, \dots, m\}$ ، متغیر z_i را به صورت زیر تعریف می کنیم.

$$z_i = \begin{cases} 1 & \text{if } h(x_i) \neq y_i \\ 0 & \text{if } h(x_i) = y_i \end{cases} \quad (۵۸.۲)$$

با توجه به $\hat{\mathbf{R}}(h) = \frac{1}{m} \sum_{i=1}^m z_i$ و $\mathbf{R}(h) = \mathbb{E}[\hat{\mathbf{R}}(h)]$ و نابرابری هافدینگ داریم

$$\mathbb{P} \left[\left| \hat{\mathbf{R}}(h) - \mathbf{R}(h) \right| \geq \epsilon \right] \leq 2 \exp(-2m\epsilon^2) \leq \delta. \quad (۵۹.۲)$$

با توجه به اینکه سمت راست نابرابری بالا باید کوچکتر از δ باشد و از طرفین این نابرابری \log گرفته شود خواهیم داشت.

$$\epsilon \geq \sqrt{\frac{\ln \frac{2}{\delta}}{2m}} \quad (۶۰.۲)$$

و با استفاده از تعریف مدل یادگیری احتمالا تقریبا درست بدون پیش فرض برای هر $\delta > 0$ ، نابرابری زیر با احتمال دست کم $(1 - \delta)$ برقرار است.

$$|\mathbf{R}(h) - \hat{\mathbf{R}}(h)| \leq \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (۶۱.۲)$$

نتیجه ۱.۲.

فرضیه $\{0, 1\} : \mathcal{X} \rightarrow h$ را در نظر بگیرید آنگاه برای هر $\delta > 0$ ، نابرابری زیر با احتمال دست کم $(1 - \delta)$ برقرار است.

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (۶۲.۲)$$

برهان اگر سمت راست معادله (۵۶.۲) را برابر δ قرار دهیم

$$2e^{(-2m\epsilon^2)} = \delta \quad (۶۳.۲)$$

و سپس برای ϵ حل نماییم خواهیم داشت

$$\epsilon^2 = \frac{\ln \frac{1}{\delta}}{2m} \quad (۶۴.۲)$$

اما ϵ کران پایین $|\hat{\mathbf{R}}(h) - \mathbf{R}(h)|$ در قضیه ۶.۲ است. لذا داریم

$$|\hat{\mathbf{R}}(h) - \mathbf{R}(h)|^2 \geq \epsilon^2 = \frac{\ln \frac{1}{\delta}}{2m} \quad (۶۵.۲)$$

$$|\hat{\mathbf{R}}(h) - \mathbf{R}(h)| \geq \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (۶۶.۲)$$

$$\hat{\mathbf{R}}(h) - \mathbf{R}(h) \geq \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (۶۷.۲)$$

و لذا خواهیم داشت

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (۶۸.۲)$$

آیا کران بدست آمده در نتیجه‌ی ۱.۲ را می‌توان برای فرضیه h که توسط الگوریتم یادگیری براساس داده‌های آموزشی S تولید می‌شود را تعمیم داد؟ پاسخ منفی است زیرا (الف) h یک فرضیه ثابت است و در اینجا توسط الگوریتم یادگیری مشخص می‌شود و (ب) h یک متغیر تصادفی است که به مجموعه آموزشی S وابسته است. بنابراین در نتیجه‌ی ۱.۲ دیگر رابطه $\mathbb{E}[\hat{\mathbf{R}}(h)] = \mathbf{R}(h)$ برقرار نیست زیرا $\mathbf{R}(h)$ یک متغیر تصادفی است و $\hat{\mathbf{R}}(h)$ یک عدد ثابت. در قضیه‌ای که در ادامه بیان می‌شود نتیجه‌ی ۱.۲ را برای این منظور تعمیم می‌دهیم.

قضیه ۷.۲ (کران خطا برای یادگیر بدون پیش فرض). فرض کنید که H مجموعه ای متناهی از فرضیه ها باشد آنگاه برای هر $\delta > 0$ با احتمال دست کم $(1 - \delta)$ ، نابرابری زیر برای همه فرضیه های $h \in H$ برقرار است.

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \sqrt{\frac{\log|H| + \log \frac{2}{\delta}}{2m}} \quad (۶۹.۲)$$

برهان فرض کنید که فرضیه های $\{h_1, h_2, \dots, h_{|H|}\}$ اعضا مجموعه H باشند، آنگاه بر اساس کران اجتماعات و سپس اعمال نتیجه‌ی ۱.۲ خواهیم داشت

$$\begin{aligned} \mathbb{P} \left[\exists h \in H : \left| \hat{\mathbf{R}}(h) - \mathbf{R}(h) \right| \geq \epsilon \right] &= \mathbb{P} \left[\left(\left| \hat{\mathbf{R}}(h_1) - \mathbf{R}(h_1) \right| \geq \epsilon \right) \vee \dots \vee \left(\left| \hat{\mathbf{R}}(h_{|H|}) - \mathbf{R}(h_{|H|}) \right| \geq \epsilon \right) \right] \\ &\leq \sum_{h_i \in H} \mathbb{P} \left[\left| \hat{\mathbf{R}}(h_i) - \mathbf{R}(h_i) \right| \geq \epsilon \right] \end{aligned} \quad (۷۰.۲)$$

$$\leq 2|H| \exp(-2m\epsilon^2) \leq \delta. \quad (۷۱.۲)$$

و با حل معادله بالا براساس ϵ اثبات قضیه کامل می شود.

اگر نابرابری (۷۱.۲) را براساس m حل کنیم پیچیدگی نمونه‌ای یادگیری فضا بصورت زیر خواهد شد.

$$m \geq \frac{1}{2\epsilon^2} \left[\log|H| + \log \frac{2}{\delta} \right] \quad (۷۲.۲)$$

قضیه بالا نشان می دهد که هرچه m بزرگتر باشد خطای واقعی کمتر می شود و هرچه $|H|$ بزرگتر باشد خطای تجربی کمتر می شود اما خطای واقعی افزایش می یابد. بنابراین باید مصالحه‌ای بین فضای فرضیه و خطای واقعی صورت گیرد. این موضوع تیغ اوکام را نتیجه می دهد.

بعد VC و پیچیدگی رادمیچرا

در بخش پیش نشان داده شد که اگر H فضای فرضیه متناهی باشد آنگاه ممکن است مساله قابل یادگیری باشد. این اثبات براساس متناهی بودن H بنا شده است. پیش از این مثال‌های مختلفی از فضاهای فرضیه همانند آستانه یک بعدی، مستطیل‌هایی با اضلاع موازی محورها، خط و بازه داشتیم که متناهی نبودند. اما در مدل‌های مختلف یادگیری، یادگرفتنی بودند. در نتیجه کران‌های بدست آمده در بخش پیش قابل تعمیم به فضای فرضیه‌های نامتناهی همانند فضای خط‌ها و فضای توابع چند جمله‌ای نیست. به همین دلیل اندازه فضای H مهم نیست بلکه قدرت توصیف یا غنی بودن این فضا مهم است. برای اینکه بتوانیم از این فضای فرضیه‌ها استفاده نماییم می‌بایست جمله $\ln |H|$ در کران خطای واقعی یا پیچیدگی نمونه ای را با معیارهای دیگری که پیچیدگی یا غنا یا قدرت توصیف فضای فرضیه را نشان دهند جایگزین نماییم. برای اینکه شهود بهتری از پیچیدگی فرضیه‌ها را داشته باشیم فرض کنید که در فضای H ، هر فرضیه با تابع k پارامتر حقیقی نمایش داده شود. برای نمونه در فضای فرضیه‌های مکعب مستطیل‌های n بعدی که اضلاع آن مواز محورها است $k = 2n$ پارامتر و در یک ابر صفحه در فضای n بعدی $k = n + 1$ پارامتر داریم. در این حالت ما می‌خواهیم این فرضیه‌ها را در حافظه محدود کامپیوتر نشان دهیم. اگر هر پارامتر را با b بیت نمایش دهیم نیاز به kb بیت برای نمایش تابع نیاز خواهیم داشت و در نتیجه فضای فرضیه H توان نمایش 2^{kb} فرضیه متفاوت را خواهد داشت. برای سادگی فرض کنید که فرضیه پیدا شده با داده‌های آموزشی سازگار باشد در این صورت با توجه به اینکه فضای فرضیه دارای $|H| = 2^{kb}$ فرضیه است و با جایگذاری این مقدار در کران پیچیدگی نمونه ای، پیچیدگی دارای کران پایین زیر خواهد بود.

$$m = O\left(\frac{1}{\epsilon} \left[k + \ln \frac{1}{\delta}\right]\right) \quad (1.3)$$

در این حالت تعداد نمونه‌های مورد نیاز یک تابع خطی از تعداد پارامترها است. هر چند کران بالا یک رابطه خطی را برای پیچیدگی نمونه ای نشان می‌دهد اما گویای قدرت توصیف یا دسته بندی فضای فرضیه H نیست. در ادامه این بخش به بیان تعدادی معیار که توانایی

بیشتری در بیان پیچیدگی و قدرت توصیف فضای فرضیه دارند پرداخته می شود

۱.۳ تابع رشد

فرض کنید که H فضای فرضیه و $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$ مجموعه آموزشی باشد. تابع رشد فضای فرضیه H برای مجموعه آموزشی S که با $\Pi_H(S)$ نشان داده می شود به صورت زیر تعریف می شود.

تعریف ۱.۳ تابع رشد فضای فرضیه H برای مجموعه آموزشی S برابر است با تعداد دسته بندی های ممکن S بوسیله H .

به عبارتی دیگر $\Pi_H(S)$ برابر است با تعداد بردارهای مجزای $(h(x_1), h(x_2), \dots, h(x_m))$. لذا برای هر فضای فرضیه $h \in H$ و برای هر مجموعه آموزشی S می توان نشان داد که نتیجه زیر برقرار است.

نتیجه ۱.۳. هر فضای فرضیه H و برای مجموعه آموزشی S نامساوی $\Pi(S) \leq 2^m$ برقرار است.

یک فاکتور مهم در نمایش قدرت توصیف فضای فرضیه ها بیشترین تعداد دسته بندی های ممکن است زیرا بسیاری از توابع در فضای H رفتارهای مشابهی دارند. در نتیجه تابع رشد برابر است با

$$\Pi_H(m) = \max\{\Pi_H(s) \mid s \in \mathcal{X}^m\}. \quad (2.3)$$

به تابع $\Pi_H(S)$ ، تابع رشد می گویند. در ادامه چند مثال از توابع رشد برای فضاهای فرضیه مختلف بررسی می شود.

مثال ۱.۳ (توابع آستانه یک بعدی).

فرض کنید که H مجموعه ای از توابع آستانه یک بعدی باشد. اگر مجموعه آموزشی دارای یک نقطه باشد به دو حالت می توان این نمونه را برچسب زد. لذا $\Pi_H(S) = 2$. اگر مجموعه آموزشی دارای دو نقطه باشد به سه حالت می توان این نمونه را برچسب زد. لذا $\Pi_H(S) = 3$. و اگر مجموعه آموزشی دارای m نقطه باشد تابع رشد برابر است با $\Pi_H(S) = m + 1$. همان گونه که تابع رشد نشان می دهد مقدار این تابع بسیار کوچکتر از 2^m است. فرض کنید که نمونه های آموزشی به صورت $x_1 < x_2 < \dots < x_m$ مرتب شده باشند. حال برای یک مقدار $\theta \in \mathbb{R}$ ، $h_\theta(x) = 1$ است اگر $x \geq \theta$. در نتیجه برای هر $k \in [1, m-1]$ اگر $h_\theta(x_k) = 1$ باشد نتیجه می گیریم که $h_\theta(x_{k+1}) = 1$ است. در نتیجه بردار $(h(x_1), h(x_2), \dots, h(x_m))$ دارای $m + 1$ مقدار زیر خواهد بود.

$$(\circ, \circ, \dots, \circ), (\circ, \circ, \dots, 1), \dots, (\circ, 1, \dots, 1), (1, 1, \dots, 1)$$

در نتیجه مقدار تابع رشد برای این فضا برابر است با

$$\Pi_H(S) = m + 1.$$

همان گونه که روشن است اندازه این فضا بی نهایت است یعنی بی نهایت مقدار برای θ وجود دارد در حالیکه بسیاری از این فرضیه ها برای مجموعه آموزشی همانند هم هستند و قابل تفکیک نیستند. بنابراین تابع رشد یک پارامتر بهتری نسبت به تعداد فرضیه ها است.

مثال ۲.۳ (فضای فرضیه های نیم خط).

این کلاس فرضیه بسیار شبیه توابع آستانه یک بعدی است با این تفاوت که نمونه های با مقدار بزرگتر از θ می توانند برچسب مثبت یا منفی داشته باشند. در حالیکه در توابع آستانه یک بعدی نمونه های با مقدار بزرگتر از θ تنها می توانستند برچسب مثبت داشته باشند. بنابراین تعداد حالت های مجزا برای دسته بندی دو برابر تعداد حالت های مجزا برای دسته بندی در توابع آستانه یک بعدی است. اگر تعداد حالت هایی که دو بار شمرده شده اند را کم نماییم خواهیم داشت.

$$\Pi_H(m) = 2(m+1) - 2 = 2m. \quad (۳.۳)$$

مثال ۳.۳ (فضای فرضیه های بازه).

فرض کنید که H مجموعه ای از توابع بازه باشد. یا به عبارتی $H = \{[a, b] \mid a < b \in \mathbb{R}\}$. یک نمونه به صورت زیر دسته بندی می شود.

$$h(x) = \begin{cases} 1 & \text{if } x \in [a, b] \\ 0 & \text{if } x \notin [a, b] \end{cases} \quad (۴.۳)$$

اگر m نمونه آموزشی داشته باشیم $m+1$ حالت برای قرارداد هر کدام از مرزها خواهیم داشت. در نتیجه $\binom{m+1}{2} + 1 = \Pi_H(m)$. حالت $+1$ زمانی رخ می دهد که بازه بین دو نمونه متوالی قرار بگیرند و در نتیجه خواهیم داشت

$$\Pi_H(m) = \binom{m+1}{2} + 1 = O(m^2) \quad (۵.۳)$$

که باز هم این مقدار بسیار کمتر از 2^m است.

در بدترین حالت به یک فضای فرضیه H نیاز خواهیم داشت که $\Pi_H(S) = 2^m$. از آنجایی که شرط $\Pi_H(S) < |H|$ برقرار است می توانیم با بکار بردن $\Pi_H(S)$ به جای H کران سفت تری را برای پیچیدگی نمونه ای پیدا کنیم و لذا قضیه زیر را می توان اثبات نمود.

قضیه ۱.۳ (کران خطا برای فرضیه های سازگار).

فرض کنید H خانواده ای از توابع باشد که مقدار $\{-1, +1\}$ را می پذیرند آنگاه برای همه $h \in H$ های سازگار، برای هر $\delta > 0$ ، با احتمال دست کم $1 - \delta$ رابطه زیر برقرار است.

$$\mathbf{R}(h) = O\left(\frac{\ln \Pi_H(2m) + \ln \frac{1}{\delta}}{m}\right). \quad (۶.۳)$$

برهان اثبات این قضیه را به عنوان تمرین انجام دهید.

قضیه ۲.۳ (کران خطا برای فرضیه های ناسازگار).

فرض کنید H خانواده ای از توابع باشد که مقدار $\{-1, +1\}$ را می پذیرند آنگاه برای همه $h \in H$ ها، برای هر $\delta > 0$ ، با احتمال دست کم $1 - \delta$ رابطه زیر برقرار است.

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \sqrt{\frac{2 \ln \Pi_H(m)}{m}} + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}}. \quad (۷.۳)$$

برهان اثبات این قضیه را به عنوان تمرین انجام دهید.

۲.۳ بعد VC

در بخش پیش تابع رشد بیان گردید و اشاره شد که تابع رشد معیار مناسب‌تری نسبت به اندازه فضای فرضیه بود و کران سفت‌تری را نسبت به آن نتیجه می‌دهد یعنی تعداد دسته بندهای مختلفی که یک کلاس فرضیه می‌تواند ایجاد نماید یک معیار خوب است. همچنین تعداد دسته بندی‌های ممکن بوسیله H برای هر مجموعه آموزشی با اندازه m و داده‌های دو به دو مجزا حداکثر 2^m است. اما مثال‌های بیان شده نشان دادند که معمولاً $\Pi_H(m) \leq 2^m$ است و محاسبه تابع رشد برای بسیاری از فضاها سخت و پیچیده می‌باشد. حال اگر ما حالتی را فرض کنیم که H امکان ایجاد همه 2^m دسته بندی ممکن را داشته باشد در اینصورت اندازه مجموعه آموزشی تغییر می‌کند. اما نشان داده شد که اگر اندازه تابع رشد از 2^m کمتر باشد کلاس فرضیه مورد نظر نمی‌تواند برخی از برچسب گذاری‌های ممکن را پوشش دهد. بنابراین این یک معیار خوب می‌تواند این باشد که یک کلاس فرضیه توانایی دسته بندی تمام برچسب گذاری‌های ممکن یک مجموعه آموزشی با چه اندازه را خواهد داشت. این نوع نگرش معیاری دیگر را نتیجه می‌دهد که نام آن بعد VC^۱ است و در این بخش به معرفی آن می‌پردازیم. این معیار همانند تابع رشد نیز یک معیار پیچیدگی کاملاً ترکیب‌یاتی است و محاسبه آن نیز از تابع رشد راحت‌تر است. برای تعریف بعد VC ابتدا باید مفاهیم دورستگی^۲ و تفکیک^۳ را مشخص کرد. مجموعه فرضیه‌های H را در نظر بگیرید، یک دورستگی مجموعه S یکی از راه‌های برچسب‌گذاری داده‌های S بر اساس یکی از فرضیه‌های H است. همچنین گوئیم مجموعه فرضیه‌های H مجموعه S با $m \geq 1$ نمونه را تفکیک کرده‌است، اگر H همه دورستگی‌های ممکن از S را بپوشاند، یعنی $\Pi_H(m) = 2^m$. در ادامه این بخش نخست چند تعریف و مثال آورده می‌شود.

تعریف ۲.۳ دورستگی دورستگی مجموعه S با اندازه $m > 1$ بوسیله فضای فرضیه H عبارت است از دسته بندی یکی از برچسب گذاری‌های ممکن S به دو دسته.

تعریف ۳.۳ تفکیک داده همه دورستگی‌های مجموعه S با اندازه $m > 1$ بوسیله فضای فرضیه H ایجاد می‌شود اگر H همه برچسب گذاری‌های ممکن یا افزای‌های ممکن S را امکان پذیر نماید. یعنی داشته باشیم $\Pi_H(m) = 2^m$.

براساس دو تعریف دورستگی و تفکیک داده‌ها بعد VC فضای فرضیه H به صورت زیر تعریف می‌شود.

تعریف ۴.۳ بعد VC بعد VC فضای فرضیه H عبارت است از اندازه بزرگترین مجموعه‌ای که توسط H به طور کامل تفکیک شود و به عبارتی

$$VC(H) = \max \{m \mid \Pi_H(m) = 2^m\} \quad (۸.۳)$$

بنابراین اگر شرط $VC(H) = d$ برقرار باشد آنگاه مجموعه‌ای از اندازه d وجود دارد که همه دورستگی‌های آن توسط H شکسته می‌شود. اما به این معنی نیست که دورستگی‌های همه مجموعه‌های از اندازه d و یا کمتر توسط H شکسته می‌شوند. برای محاسبه بعد VC

¹Vapnik-Chervonenkis Dimension

²dichotomy

³shattering

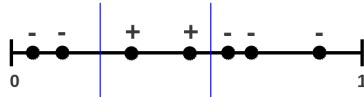
معمولا یک کران پایین برای آن اثبات می‌شود و برای این کار باید نشان داد که همه دویخشی‌های مجموعه‌ای مانند S از اندازه d توسط H شکسته می‌شود. برای پیدا کردن یک کران بالا برای بعد VC باید نشان داد که هیچ مجموعه‌ای از اندازه $d+1$ توسط H شکسته نمی‌شود.

مثال ۴.۳ (بعد VC آستانه یک بعدی).

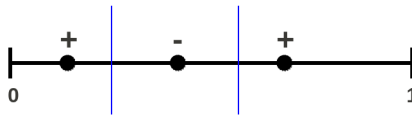
فضای فرضیه آستانه یک بعدی را در نظر بگیرید. این فضا تنها می‌تواند یک نمونه را تفکیک کند و لذا $VC(H) = 1$ خواهد بود.

مثال ۵.۳ (بازه‌های روی یک خط).

کران پایین بعد VC برای این مجموعه فرضیه ۲ است. برای مشاهده این نکته می‌توان همه دسته بندی‌های دو دسته ای دو نقطه را در نظر گرفت که در شکل ۱.۳ نشان داده شده‌است که توسط یک بازه از هم تفکیک شده‌اند. همچنین در شکل ۲.۳ با یک مثال نشان داده شده‌است که هیچ مجموعه سه تایی از نقاط را نمی‌توان توسط این مجموعه فرضیه از هم تفکیک کرد. بنابراین $VC(H) = 2$.



شکل ۱.۳ نمونه ای از نقاط که می‌توان با یک بازه آنها را از هم جدا کرد.



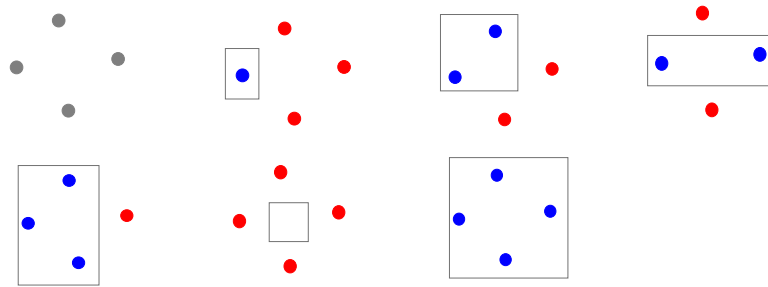
شکل ۲.۳ سه نقطه ای را که نمی‌توان با یک بازه آنها را از هم جدا کرد.

مثال ۶.۳ (مستطیل‌های با اضلاع موازی محورهای مختصات).

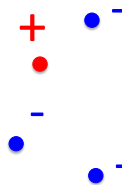
برای این مجموعه از فرضیه‌ها می‌توان نشان داد که کران پایین بعد VC دست کم ۴ است. برای این منظور ۴ نقطه روی یک لوزی در نظر بگیرید که اضلاع آن موازی محورها نیستند. همه حالت‌های دسته بندی دوتایی از این نقاط را می‌توان توسط یک مستطیل تشخیص داد. شکل ۳.۳ یک نمونه دیگر ای ۴ نقطه را نشان می‌دهد که می‌توانیم آنها را با یک مستطیل تفکیک کنیم. اما برای هر ۵ نقطه در صفحه اگر کوچکترین مستطیل دربردارنده همه نقاط را در نظر بگیرید، همواره یک نقطه درون این مستطیل قرار خواهد گرفت و با دادن برچسب - به آن و دادن برچسب + به مابقی نقاط مطابق شکل ۴.۳ هرگز نمی‌توان مستطیلی با اضلاع موازی محورها یافت که این کلاس‌ها را از هم تفکیک کند. بنابراین بعد VC برای این مجموعه فرضیه‌ها برابر با ۴ است.

مثال ۷.۳ (بعد VC خطوط در فضای دو بعدی).

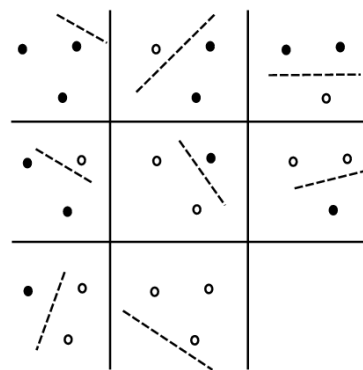
فضای فرضیه خطوط در فضای دو بعدی را در نظر بگیرید. مطابق شکل ۵.۳ دست کم یک سه نقطه‌ای وجود دارد که می‌توانیم ۳ نقطه در فضای دویبعدی را با یک خط تفکیک کنیم. بنابراین حداقل بعد VC برابر با ۳ است اما مطابق شکل نمی‌توانیم یک ۴ نقطه‌ای را در فضای دویبعدی پیدا کنیم که با یک خط از هم تفکیک نماییم. بنابراین بعد VC خط در فضای دویبعدی کمتر از ۵ است و در نتیجه بعد VC خط در فضای دویبعدی مساوی ۴ است.



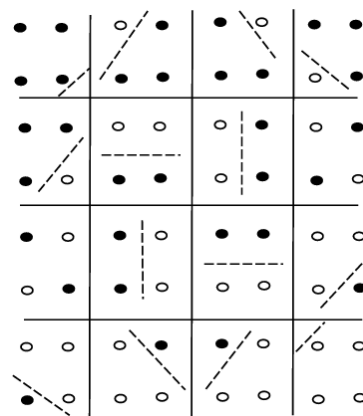
شکل ۳.۳ نمونه ای از ۴ نقطه که می توان با یک مستطیل از هم جدا کرد.



شکل ۴.۳ نمونه ای از ۵ نقطه که نمی توان با یک مستطیل از هم جدا کرد.



شکل ۵.۳ نمونه ای از سه نقطه که می توان با یک خط آنها را از هم جدا کرد.



شکل ۶.۳ چهار نقطه ای را که نمی توان با یک خط آنها را از هم جدا کرد.

قضیه ۳.۳ (کران بعد VC برای فضای فرضیه های متناهی).
اگر H یک فضای فرضیه متناهی باشد، آنگاه خواهیم داشت $VC(H) \leq \log |H|$.

برهان فرض کنید که $VC(H) = d$ باشد. بنابراین رابطه زیر برقرار است.

$$\Pi_H(d) = 2^d. \quad (9.3)$$

اما برای هر مجموعه آموزشی با اندازه $m > 1$ داریم $\Pi_H(m) \leq |H|$ و با توجه به رابطه بالا داریم $2^d = \Pi_H(d) \leq |H|$. اگر از دو طرف رابطه بالا لگاریتم بگیریم قضیه اثبات می شود.

مثال ۸.۳ (بعد VC ترکیب عطفی حداکثر n متغیر).
برای فضای فرضیه H که از ترکیب عطفی حداکثر n ویژگی یا نقیض آن ساخته می شود رابطه زیر برقرار است.

$$n \leq VC(H) \leq n \log 3. \quad (10.3)$$

درحالت کلی برای هر فضای بردارای می توان نشان داد که لم زیر برقرار است.

لم ۱.۳ (بعد VC فضای برداری).
فرض کنید \mathcal{G} مجموعه ای از توابع در فضای برداری \mathbb{R}^d برای $d < \infty$ باشد. حال اگر مجموعه \mathcal{F} را به صورت

$$\mathcal{F} = \{x \mapsto \text{sgn}(g(x)) \mid g \in \mathcal{G}\}$$

تعریف کنیم. اگر بعد \mathcal{G} برابر با k باشد آنگاه بعد VC این فضا برابر است با

$$VC(\mathcal{F}) \leq k.$$

برهان اثبات این لم را به عنوان تمرین انجام دهید.

لم ۲.۳ (لم Sauer).
اگر H فرضیه ای با $VC(H) = d$ ، آنگاه برای هر $m \in \mathbb{N}$ رابطه زیر برقرار است:

$$\Pi_H(m) \leq \sum_{i=0}^d \binom{m}{i} \quad (11.3)$$

برهان اثبات این لم به کمک استقرا روی $m + d$ انجام می گیرد. شما به عنوان تمرین آن را انجام دهید.

نتیجه ۲.۳.

اگر H مجموعه‌ای از فرضیه‌ها باشد به طوری که $VC(H) = d$ ، آنگاه برای هر $m \geq d$ داریم

$$\Pi_H(m) \leq \left(\frac{em}{d}\right)^d = O(m^d) \quad (۱۲.۳)$$

برهان با استفاده از لم ۲.۳ داریم

$$\Pi_H(m) \leq \sum_{i=0}^d \binom{m}{i} \quad (۱۳.۳)$$

$$\leq \sum_{i=0}^d \binom{m}{i} \underbrace{\left(\frac{m}{d}\right)^{d-i}}_{> 1} \quad (۱۴.۳)$$

$$\leq \sum_{i=0}^m \binom{m}{i} \left(\frac{m}{d}\right)^{d-i} \quad (۱۵.۳)$$

$$= \left(\frac{m}{d}\right)^d \sum_{i=0}^m \binom{m}{i} \left(\frac{d}{m}\right)^i \quad \text{بر اساس توزیع دو جمله ای} \quad (۱۶.۳)$$

$$= \left(\frac{m}{d}\right)^d \left(1 + \frac{d}{m}\right)^m \quad \text{بر اساس نامساوی } (1-x) \leq e^{-x} \quad (۱۷.۳)$$

$$\leq \left(\frac{m}{d}\right)^d \left(e^{d/m}\right)^m \quad (۱۸.۳)$$

$$= \left(\frac{m}{d}\right)^d e^d = \left(\frac{me}{d}\right)^d \quad (۱۹.۳)$$



نتیجه ۳.۳.

اگر H مجموعه‌ای از توابع باشد که مقادیر آن از $\{-1, +1\}$ و $VC(H) = d$ باشد. آنگاه برای هر $\delta > 0$ با احتمال دست کم $1 - \delta$ برای هر $h \in H$ داریم:

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \sqrt{\frac{2d \log \frac{em}{d}}{m}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \quad (۲۰.۳)$$

برهان از کران تابع رشد داشتیم

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \sqrt{\frac{2 \ln \Pi_H(m)}{m}} + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}}. \quad (۲۱.۳)$$

همچنین با استفاده از نتیجه لم ۲.۳ خواهیم داشت

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \sqrt{\frac{2 \ln \Pi_H(m)}{m}} + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (۲۲.۳)$$

$$\leq \hat{\mathbf{R}}(h) + \sqrt{\frac{\sqrt{2 \ln \left(\frac{me}{d} \right)^d}}{m}} + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (23.3)$$

$$\leq \hat{\mathbf{R}}(h) + \sqrt{\frac{\sqrt{2d \ln \frac{me}{d}}}{m}} + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}}. \quad (24.3)$$

بنابراین شکل کلی کران به صورت زیر است:

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + O\left(\sqrt{\frac{\log(m/d)}{(m/d)}}\right) \quad (25.3)$$

نتیجه‌گیری ۲.۳ بسیار مهم است زیرا این نتیجه‌گیری می‌گوید که همه کلاس‌های فرضیه‌ها در یکی از دو گروه زیر قرار می‌گیرند.

۱. اگر d نامتناهی باشد آنگاه $\Pi_H(m) = 2^m$ است و در نتیجه کران زیر

$$\mathbf{R}(h) = O\left(\frac{\ln \Pi_H(2^m) + \ln \left(\frac{1}{\delta}\right)}{m}\right). \quad (26.3)$$

بدون معنی می‌باشد.

۲. اگر d متناهی باشد آنگاه رابطه $\Pi_H(m) = O(m^d)$ برقرار است. در این حالت قضیه یاد شده نتایج جالبی دارد. از آنجایی که

$\ln \Pi_H(m) = O(d \ln m)$ ، لذا کران خطا به صورت خطی به d و به صورت وارون به m وابسته است و با افزایش m به سمت

صفر میل می‌کند.

از لم ۲.۳ به سادگی می‌توان قضیه زیر را نتیجه گرفت.

قضیه ۴.۳ (۰)

اگر $VC(H) = d$ آنگاه برای همه فرضیه‌های سازگار درون H ، با احتمال دست کم $1 - \delta$ نابرابری زیر برقرار است.

$$\mathbf{R}(h) = O\left(\frac{d \ln m + \ln \left(\frac{1}{\delta}\right)}{m}\right) \quad (27.3)$$

و به طور مشابه می‌توان نتیجه گرفت که

$$m = O\left(\frac{1}{\epsilon} \log \frac{1}{\delta} + \frac{d}{\epsilon} \log \frac{1}{\epsilon}\right). \quad (28.3)$$

در نتیجه پیچیدگی نمونه‌ای به صورت خطی با d رشد می‌نماید. برای نمونه می‌توانیم با افزودن تعداد ویژگی‌ها، بعد $VC(H)$ را افزایش دهیم. در نتیجه یادگیری هر بعد به تعداد ثابتی نمونه نیاز دارد. برای روشن شدن نتیجه بیان شده، چند مثال بیان می‌کنیم.

مثال ۹.۳ (تابع آستانه یک بعدی).

پیش از این نشان دادیم که برای فضای فرضیه آستانه یک بعدی $VC(H) = 1$ است. بنابراین با جایگزین نمودن آن در نتیجه

قضیه قبل داریم

$$m = O\left(\frac{1}{\epsilon} \log \frac{1}{\delta} + \frac{1}{\epsilon} \log \frac{1}{\epsilon}\right). \quad (29.3)$$

اما پیش از این برای همین فضای فرضیه نشان دادیم که

$$m \geq \frac{1}{\epsilon} \ln \frac{2}{\delta}. \quad (30.3)$$

دو نتیجه بالا نشان می دهند که کرانی که بوسیله بعد VC پیدا می شود بد نیست.

مثال ۱۰.۳ (ابرفضای در فضای n بعدی).

پیش از این نشان دادیم که بعد VC ابرفضای n بعدی برابر ۱ + n است. در نتیجه داریم

$$m = O\left(\frac{1}{\epsilon} \log \frac{1}{\delta} + \frac{n+1}{\epsilon} \log \frac{1}{\epsilon}\right). \quad (31.3)$$

این کران در واقع همان نتیجه گیری بدی را نشان می دهد که پیش از این در باره آن سخن به میان آمد. این رابطه نشان می دهد که پیچیدگی نمونه ای، یک تابع خطی از تعداد پارامترهای فرضیه یا همان تعداد ویژگی های نمونه است. حال مثال زیر را در نظر بگیرید.

مثال ۱۱.۳ (مستطیل های با اضلاع موازی محورها).

پیش از این نشان دادیم که برای فضای فرضیه مستطیل هایی که اضلاع آن موازی محورهای مختصات است $VC(H) = 4$ است. بنابراین با جایگزین نمودن آن در نتیجه قضیه قبل داریم

$$m = O\left(\frac{1}{\epsilon} \log \frac{1}{\delta} + \frac{4}{\epsilon} \log \frac{1}{\epsilon}\right). \quad (32.3)$$

اما پیش از این برای همین فضای فرضیه نشان دادیم که

$$m \geq \frac{4}{\epsilon} \ln \frac{4}{\delta}. \quad (33.3)$$

حال مثال زیر را در نظر بگیرید که نتیجه بد گذشته را نقض می کند.

مثال ۱۲.۳ (فضای فرضیه).

فرض کنید که H فضای فرضیه $\text{sgn}(\sin(\theta x))$ باشد. می توان نشان داد که بعد $VC(H) = \infty$ است. اما این فرضیه تنها یک پارامتر دارد و نتیجه بد گذشته را نقض می کند.

کران های ارایه شده در این بخش را می توان به حالت تحقق ناپذیر نیز گسترش داد. در این حالت، این کران به صورت زیر خواهد شد (در روابط به جای $\frac{1}{\epsilon}$ مقدار $\frac{1}{\epsilon^2}$ قرار می گیرد).

$$m = O\left(\frac{1}{\epsilon^2} \log \frac{1}{\delta} + \frac{d}{\epsilon^2} \log \frac{1}{\epsilon}\right). \quad (34.3)$$

رابطه (۲۸.۳) کران بالای پیچیدگی نمونه ای را نشان می دهد که به صورت خطی به بعد VC وابسته است. آیا این وابستگی خطی سفت

است؟ برای پاسخ به این پرسش می‌بایست کران پایین پیچیدگی نمونه‌ای را پیدا نمود. می‌توان نشان داد که دست کم $\frac{d}{4}$ نمونه لازم است تا برای هر $\frac{1}{8} < \delta < \epsilon$ ، با احتمال دست کم $1 - \delta$ نابرابری $\mathbf{R}(h) \leq \epsilon$ برقرار است. این کران پایین را می‌توان به صورت $\Omega(d/\epsilon)$ نشان داد. در نتیجه کران (۲۸.۳) سفت است و علاوه بر وابستگی خطی پیچیدگی نمونه‌ای به d ، به صورت خطی به $\frac{1}{\epsilon}$ نیز وابسته است. این نتیجه در قضیه زیر نشان داده شده است

قضیه ۵.۳. (۱)

فرض کنید $VC(H) = d$ باشد. برای هر الگوریتم یادگیری A ، یک فرضیه h و یک توزیع \mathcal{D} وجود دارد که اگر الگوریتم یادگیری A ، مجموعه آموزشی با $\frac{d}{4} \leq m \leq \frac{d}{\epsilon}$ نمونه که از توزیع \mathcal{D} به صورت مستقل نمونه برداری شده باشند و توسط تابع c برچسب گذاری شده باشند را دریافت و فرضیه h را تولید نماید آنگاه نابرابری زیر برقرار است.

$$\mathbb{P} \left[\mathbf{R}(h) > \frac{1}{8} \right] > \frac{1}{8}. \quad (۳۵.۳)$$

همچنین نشان داده شده است که یک فضای فرضیه H و یک توزیع \mathcal{D} وجود دارد بطوریکه $\left(\frac{1}{\epsilon} \left[d \ln \frac{1}{\epsilon} + \ln \frac{1}{\delta} \right] \right)$ نمونه نیاز است تا برای هر فرضیه سازگار $h \in H$ نابرابری $\mathbf{R}(h) \leq \epsilon$ برقرار باشد. همچنین نشان داده شده است که برای هر $k > 1$ ، الگوریتمی وجود دارد که تنها

$$O \left(\frac{1}{\epsilon} d \log^{(k)} \frac{1}{\epsilon} + \frac{1}{\epsilon} \log \frac{1}{\delta} \right) \quad (۳۶.۳)$$

نمونه نیاز دارد تا با احتمال دست کم $1 - \delta$ کران بالای خطای واقعی آن ϵ باشد که $\log^{(k)}(x) = \log(\log(\dots \log(x)))$ است که در رابطه k بار تکرار شده است. در رابطه بالا پیچیدگی نمونه‌ای به k وابسته است.

۳.۳ پیچیدگی رادامیچر

پیش از این مجموعه آموزشی $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$ و فرضیه $h : \mathcal{X} \mapsto \{0, 1\}$ را در نظر گرفتیم. در این حالت خطای آموزشی برابر است با

$$\hat{\mathbf{R}}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}[h(x_i) \neq y_i] \quad (۳۷.۳)$$

حال اگر مجموعه برچسب‌ها را به صورت $\mathcal{Y} = \{-1, +1\}$ در نظر بگیریم فرضیه h به صورت فرضیه $h : \mathcal{X} \mapsto \{-1, +1\}$ خواهد بود. براین اساس خطای آموزشی را می‌توانیم به صورت زیر تعریف نماییم.

$$\hat{\mathbf{R}}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}[h(x_i) \neq y_i] \quad (۳۸.۳)$$

$$= \frac{1}{m} \sum_{i=1}^m \begin{cases} 1 & \text{if } (h(x_i), y_i) = (+1, -1) \text{ or } (h(x_i), y_i) = (-1, +1) \\ 0 & \text{if } (h(x_i), y_i) = (+1, +1) \text{ or } (h(x_i), y_i) = (-1, -1) \end{cases} \quad (۳۹.۳)$$

$$= \frac{1}{m} \sum_{i=1}^m \frac{1 - y_i h(x_i)}{2} \quad (۴۰.۳)$$

$$= \frac{1}{2} - \frac{1}{2m} \sum_{i=1}^m y_i h(x_i) \quad (41.3)$$

عبارت $\frac{1}{m} \sum_{i=1}^m y_i h(x_i)$ را می توان همبستگی برجسب واقعی y_i با برجسب پیش بینی شده $h(x_i)$ دانست. بنابراین همبستگی به صورت زیر به خطای آموزشی وابسته است.

$$\text{Corr}(h) = 1 - 2\hat{R}(h) \quad (42.3)$$

برای اینکه فرضیه h را پیدا کنیم که خطای آموزشی را کمینه نماید می توانیم فرضیه h را پیدا کنیم که همبستگی را بیشینه نماید. یا به عبارتی

$$h = \operatorname{argmax}_{h \in H} \frac{1}{m} \sum_{i=1}^m y_i h(x_i). \quad (43.3)$$

حال فرض کنید که در یک آزمایش برجسب واقعی y_i با متغیر تصادفی رادمیچر σ_i که به صورت زیر تعریف می شود جایگزین گردد.

$$\sigma_i = \begin{cases} +1 & \text{با احتمال } \frac{1}{2} \\ -1 & \text{با احتمال } \frac{1}{2} \end{cases} \quad (44.3)$$

حال هدف این است که تابع زیر بیشینه گردد.

$$h = \operatorname{argmax}_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i). \quad (45.3)$$

این تابع هدف به جای انتخاب فرضیه h که همبستگی با برجسب را بیشینه نماید فرضیه h به گونه ای پیدا می شود که همبستگی با نویز (σ_i) را بیشینه نماید. از آنجایی که فرضیه h به متغیر تصادفی σ_i وابسته است برای اندازه گیری همبستگی فضای H نسبت به σ امید ریاضی به صورت زیر محاسبه می شود.

$$\mathbb{E}_{\sigma} \left[\max_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) \right] \quad (46.3)$$

مقدار بالا به نوعی قدرت توصیف یا غنای فضای فرضیه H را نشان می دهد. براساس دو حالت حدی، کران عبارت بالا را بدست می آوریم.

۱. اگر فضای فرضیه H دارای یک فرضیه باشد در اینصورت همبستگی برابر صفر است. زیرا عملگر \max دیگر معنایی نداشته و بنابراین تنها مجموع روی $\sigma_i h(x_i)$ می باشد. از آنجایی که متغیر تصادفی رادمیچر بطور میانگین در نیمی از حالت ها دارای مقدار $+1$ و در نیمی از حالت ها دارای مقدار -1 هستند بنابراین حاصل جمع به طور میانگین صفر است.

۲. اگر فضای فرضیه H دارای 2^m فرضیه باشد در اینصورت همبستگی برابر یک است. زیرا دنباله $\sigma = (\sigma_1, \sigma_2, \dots, \sigma_m)$ دارای 2^m حالت است و در این حالت برای هر دنباله σ یک فرضیه وجود دارد لذا فضای فرضیه توانایی دسته بندی هر کدام از دنباله ها را دارد و در نتیجه مقدار عبارت بالا برابر یک می شود.

بنابراین امید ریاضی همبستگی در بازه $[0, 1]$ قرار دارد.

حال فرض کنید به جای $h : \mathcal{X} \mapsto \{-1, +1\}$ با یک مجموعه توابع حقیقی $\mathcal{F} = \{f \mid f : \mathcal{Z} \mapsto \mathbb{R}\}$ کار کنیم. جایگزین نمودن

H با \mathcal{F} یک خانواده از توابع $h : \mathcal{Z} \mapsto \mathbb{R}$ را نتیجه می دهد. حال برای مجموعه نمونه های $S = (z_1, z_2, \dots, z_m)$ با $z_i \in \mathcal{Z}$ پیچیدگی

رادامیچر به صورت زیر تعریف می شود.

$$\hat{\mathcal{R}}_S(\mathcal{F}) = \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \sigma_i f(z_i) \right] \quad (47.3)$$

که $\hat{\mathcal{R}}_S(\mathcal{F})$ پیچیدگی تجربی رادامیچر نامیده می شود و بطور میانگین همبستگی تابع \mathcal{F} را با نویز نشان می دهد. ما می خواهیم همبستگی تابع \mathcal{F} را نسبت به توزیع \mathcal{D} روی \mathcal{X} را پیدا نماییم. برای این کار امید ریاضی $\hat{\mathcal{R}}_S(\mathcal{F})$ روی تمام مجموعه داده های با طول m که براساس به توزیع \mathcal{D} نمونه برداری شده اند را محاسبه می کنیم.

$$\mathcal{R}_m(h) = \mathbb{E}_{S \sim \mathcal{D}^m} [\hat{\mathcal{R}}_S(\mathcal{F})] \quad (48.3)$$

که $\mathcal{R}_m(h)$ پیچیدگی رادامیچر نامیده می شود محاسبه پیچیدگی تجربی رادامیچر برای برخی از فضاها یک مساله ان-پی سخت است. یکی از برتری های رابطه بالا این است که توانایی مدل سازی مدل های دیگر را دارد و کران هایی تولید می شود که سفت تر و مناسب تر هستند.

قضیه ۶.۳ (۱)

فرض کنید که H خانواده ای از توابع از \mathcal{X} به مجموعه $\{-1, +1\}$ باشند و توزیع \mathcal{D} روی فضای ورودی X باشد. آنگاه برای هر $\delta > 0$ با احتمال دست کم $1 - \delta$ روی یک مجموعه S با اندازه m که از توزیع \mathcal{D} نمونه برداری شده است برای هر $h \in H$ نامساوی های زیر برقرار است.

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \mathcal{R}_m(H) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \quad (49.3)$$

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \hat{\mathcal{R}}_S(H) + 3\sqrt{\frac{\log \frac{1}{\delta}}{2m}} \quad (50.3)$$

نامساوی دوم وابسته به داده های آموزشی است. محاسبه $\mathcal{R}_m(H)$ و $\hat{\mathcal{R}}_S(H)$ سخت است و در بخش های بعد تلاش می کنیم با روش های ترکیباتی این موضوع را پیوند بزنیم.

برهان

قضیه ۷.۳ (۱)

فرض کنید که \mathcal{F} خانواده ای از توابع از \mathcal{Z} به مجموعه $[0, 1]$ را نگاشت می نمایند و مجموعه $S = (z_1, z_2, \dots, z_m)$ که $z_i \sim \mathcal{D}$ باشد. اگر $\mathbb{E}[f] = \mathbb{E}_{z \sim \mathcal{D}} [f(z)]$ و $\hat{\mathbb{E}}_S[f] = \frac{1}{m} \sum_{i=1}^m f(z_i)$ تعریف شوند. با احتمال دست کم $1 - \delta$ برای تمامی $f \in \mathcal{F}$ نامساوی زیر برقرار است.

$$\mathbb{E}[f] \leq \hat{\mathbb{E}}_S[f] + 2\mathcal{R}_m(\mathcal{F}) + O\left(\sqrt{\frac{\ln \frac{1}{\delta}}{m}}\right) \quad (51.3)$$

$$\mathbb{E}[f] \leq \hat{\mathbb{E}}_S[f] + 2\hat{\mathcal{R}}_S(\mathcal{F}) + O\left(\sqrt{\frac{\ln \frac{1}{\delta}}{m}}\right) \quad (52.3)$$

برهان برای اثبات این قضیه، نخست کران $\mathbb{E}[f] - \hat{\mathbb{E}}_S[f]$ را برای تمامی $f \in \mathcal{F}$ و یا به طور معادل کران $\sup_{f \in \mathcal{F}} \{\mathbb{E}[f] - \hat{\mathbb{E}}_S[f]\}$ را پیدا می‌کنیم. $\phi(S)$ را به صورت زیر تعریف می‌کنیم

$$\phi(S) = \sup_{f \in \mathcal{F}} \left\{ \mathbb{E}[f] - \hat{\mathbb{E}}_S[f] \right\}. \quad (53.3)$$

در رابطه بالا $\phi(S)$ یک متغیر تصادفی است و به مجموعه آموزشی S وابسته است. فرض کنید که $S = (z_1, z_2, \dots, z_i, \dots, z_m)$ و $S' = (z_1, z_2, \dots, z'_i, \dots, z_m)$ دو مجموعه آموزشی باشند که تنها در یک نمونه با هم اختلاف داشته باشند. حال در گام های زیر قضیه را اثبات می‌کنیم. در گام نخست نشان می‌دهیم که با احتمال دست کم $1 - \delta$ نابرابری

$$\phi(S) \leq \mathbb{E}_S[\phi(S)] + \sqrt{\frac{\ln \frac{2}{\delta}}{2m}}. \quad (54.3)$$

برقرار است. برای اثبات قضیه نابرابری McDiarmid را به کار می‌بریم. در نابرابری McDiarmid اگر برای همه i ها نابرابری زیر برقرار باشد

$$|f(z_1, z_2, \dots, z_i, \dots, z_m) - f(z_1, z_2, \dots, z'_i, \dots, z_m)| \leq c_i \quad (55.3)$$

آنگاه نابرابری زیر برقرار خواهد بود.

$$\mathbb{P}[|f(S) - f(S')| \geq \epsilon] \leq 2 \exp\left(-\frac{2\epsilon^2}{\sum_{i=1}^m c_i^2}\right) \quad (56.3)$$

از تعریف $\phi(S)$ داریم:

$$\phi(S) = \sup_{f \in \mathcal{F}} \left\{ \mathbb{E}[f] - \hat{\mathbb{E}}_S[f] \right\} \quad (57.3)$$

$$= \sup_{f \in \mathcal{F}} \left\{ \mathbb{E}[f] - \frac{1}{m} \sum_{i=1}^m f(z_i) \right\}. \quad (58.3)$$

از آنجایی که $f(z_i) \in [0, 1]$ است بنابراین تغییر هر یک از نمونه ها از z_i به z'_i در مجموعه آموزشی S مقدار $f(z_i)$ را حداکثر $\frac{1}{m}$ تغییر می‌دهد. لذا مقدار c_i در نابرابری McDiarmid برابر با $\frac{1}{m}$ است. اگر این مقدار را در نابرابری McDiarmid قرار دهیم خواهیم داشت.

$$\mathbb{P}[|\phi(S) - \mathbb{E}_S[\phi(S)]| \geq \epsilon] \leq 2 \exp\left(-\frac{2\epsilon^2}{\sum_{i=1}^m c_i^2}\right) \quad (59.3)$$

$$= 2 \exp\left(-\frac{2\epsilon^2}{\sum_{i=1}^m \left(\frac{1}{m}\right)^2}\right) \quad (60.3)$$

$$= 2 \exp(-2m\epsilon^2). \quad (61.3)$$

اگر مقدار ϵ را برابر $\sqrt{\frac{\log \frac{2}{\delta}}{2m}}$ قرار دهیم با ساده سازی نابرابری بالا خواهیم داشت.

$$\phi(S) \leq \mathbb{E}_S[\phi(S)] + \sqrt{\frac{\ln \frac{2}{\delta}}{2m}}. \quad (62.3)$$

در گام دوم ثابت می‌کنیم که نابرابری $\mathbb{E}_S[\phi(S)] \leq 2\mathcal{R}_m(\mathcal{F})$ برقرار است. براساس تعریف $\phi(S)$ داریم

$$\mathbb{E}_S[\phi(S)] = \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left(\mathbb{E}[f] - \hat{\mathbb{E}}_S[f] \right) \right] \quad (63.3)$$

$$= \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \mathbb{E}_{S'} \left[\mathbb{E}_{S'} [f] - \hat{\mathbb{E}}_S [f] \right] \right]. \quad (۶۴.۳)$$

از آنجایی که مجموعه S همانند مجموعه S' از توزیع \mathcal{D} نمونه برداری شده اند لذا $\mathbb{E}[f] = \mathbb{E}_{S'}[\hat{\mathbb{E}}_{S'}[f]]$ است و در نتیجه می توانیم خط دوم معادله بالا را بنویسیم. با توجه به اینکه تابع \sup یک تابع محدب است لذا می توانیم از نابرابری Jensen استفاده و معادله بالا را به صورت زیر بنویسیم.

$$\mathbb{E}_S[\phi(S)] \leq \mathbb{E}_{S, S'} \left[\sup_{f \in \mathcal{F}} \left(\hat{\mathbb{E}}_{S'}[f] - \hat{\mathbb{E}}_S[f] \right) \right] \quad (۶۵.۳)$$

$$= \mathbb{E}_{S, S'} \left[\sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m (f(z'_i) - f(z_i)) \right]. \quad (۶۶.۳)$$

اگر متغیر رادمیچر را به نابرابری بالا اضافه کنیم امید ریاضی آن تغییری نمی کند. زیرا در نابرابری بالا اگر مقدار $\sigma_i = 1$ باشد نابرابری بالا تغییری نمی کند و اگر $\sigma_i = -1$ باشد جای z_i و z'_i را عوض می کنیم چون امید ریاضی را روی تمام S و S' ها محاسبه می شود لذا تغییری در امید ریاضی نخواهیم داشت. بنابراین با افزودن متغیر رادمیچر به نابرابری بالا، خواهیم داشت.

$$\mathbb{E}_S[\phi(S)] = \mathbb{E}_{\sigma, S, S'} \left[\sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \sigma_i (f(z'_i) - f(z_i)) \right]. \quad (۶۷.۳)$$

اگر از ویژگی $\sup(a+b) \leq \sup(a) + \sup(b)$ استفاده کنیم خواهیم داشت.

$$\mathbb{E}_S[\phi(S)] \leq \mathbb{E}_{\sigma, S'} \left[\sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \sigma_i f(z'_i) \right] + \mathbb{E}_{\sigma, S} \left[\sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m -\sigma_i f(z_i) \right]. \quad (۶۸.۳)$$

با توجه به اینکه متغیرهای رادمیچر σ_i و $-\sigma_i$ دارای توزیع یکسانی هستند لذا خواهیم داشت.

$$\mathbb{E}_S[\phi(S)] \leq 2 \mathbb{E}_{\sigma, S} \left[\sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \sigma_i f(z_i) \right] \quad (۶۹.۳)$$

$$= 2\mathcal{R}_m(\mathcal{F}). \quad (۷۰.۳)$$

با جایگذاری نابرابری (۷۰.۳) در نابرابری (۵۴.۳) داریم

$$\phi(S) \leq 2\mathcal{R}_m(\mathcal{F}) + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}}. \quad (۷۱.۳)$$

که با جایگذاری، قضیه اثبات می گردد.

نابرابری (۵۲.۳) نیز به صورت مشابه اثبات می گردد. در ادامه این بخش، نخست به بررسی کاربردی از پیچیدگی رادمیچر می پردازیم.

لم ۳.۳ (کران رادمیچر برای خطای ۰-۱)

فرض کنید که فرضیه $\{ -1, +1 \} : \mathcal{X} \mapsto h$ یک دسته بند دودویی و H فضای فرضیه و همچنین $f_h(x, y) = \mathbb{I}[h(x) \neq y]$ تابع خطای صفر-یک و $\mathcal{F}_H = \{ f_h \mid h \in H \}$ باشد. آنگاه برای هر مجموعه $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$ فرض کنید $S_X = (x_1, \dots, x_m)$ تصویر S روی X باشد. آنگاه رابطه زیر بین پیچیدگی تجربی رادمیچر H و \mathcal{F} برقرار است.

$$\hat{\mathcal{R}}_S(\mathcal{F}_H) = \frac{1}{\sqrt{m}} \hat{\mathcal{R}}_{S_X}(H) \quad (۷۲.۳)$$

برهان برای اثبات این لم $Z = X \times \{-1, +1\}$ را در نظر بگیرید. براساس تعریف فضای \mathcal{F}_H ، هر تابع $f_h \in \mathcal{F}_H$ متناظر با برخی از $h \in H$ است. از تعاریف مربوط به خطا داریم

$$\mathbf{R}(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathbb{I}[h(x) \neq y]] = \mathbb{E}[f_h] \quad (۷۳.۳)$$

$$\hat{\mathbf{R}}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}[h(x_i) \neq y_i] = \hat{\mathbb{E}}_S[f_h] \quad (۷۴.۳)$$

لذا مطابق قضیه ۷.۳، کران $\mathbf{R}(h) - \hat{\mathbf{R}}(h)$ را پیدا می‌کنیم.

$$\hat{\mathcal{R}}_S(F_H) = \mathbb{E}_\sigma \left[\sup_{f_h \in \mathcal{F}_H} \frac{1}{m} \sum_{i=1}^m \sigma_i f_h(x_i, y_i) \right] \quad (۷۵.۳)$$

$$= \mathbb{E}_\sigma \left[\sup_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i \left(\frac{1 - y_i h(x_i)}{2} \right) \right] \quad (۷۶.۳)$$

$$= \mathbb{E}_\sigma \left[\sup_{h \in H} \frac{1}{2m} \sum_{i=1}^m \sigma_i + \sup_{h \in H} \frac{1}{2m} \sum_{i=1}^m (-y_i \sigma_i) h(x_i) \right] \quad (۷۷.۳)$$

$$= \frac{1}{2m} \sum_{i=1}^m \mathbb{E}_\sigma[\sigma_i] + \frac{1}{2} \mathbb{E}_\sigma \left[\sup_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) \right] \quad (۷۸.۳)$$

با توجه به ویژگی‌های متغیرهای تصادفی رادمیچر، جمله اول برابر با صفر و همچنین توزیع $-y_i \sigma_i$ همانند توزیع σ_i است لذا معادله بالا به صورت زیر خواهد شد.

$$\hat{\mathcal{R}}_S(F_H) = \frac{1}{2} \mathbb{E}_\sigma \left[\sup_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) \right] \quad (۷۹.۳)$$

$$= \frac{1}{2} \hat{\mathcal{R}}_{S_X}(H). \quad (۸۰.۳)$$

که بدین ترتیب اثبات لم کامل می‌گردد.



اگر نتیجه لم ۳.۳ را در نتایج قضیه ۷.۳ جایگزین نمایم کران خطای زیر را بدست خواهیم آورد.

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \hat{\mathcal{R}}_{S_X}(H) + O\left(\sqrt{\frac{\ln \frac{1}{\delta}}{m}}\right). \quad (۸۱.۳)$$

در ادامه رابطه بین پیچیدگی رادمیچر و معیارهای دیگری که پیش از این برای توصیف غنای فضای فرضیه بررسی نمودیم را بدست می‌آوریم.

قضیه ۸.۳ (۰).

برای هر فضای فرضیه $|H| < \infty$ نابرابری زیر برقرار است.

$$\hat{\mathcal{R}}_S(H) \leq \sqrt{\frac{2 \ln |H|}{m}}. \quad (۸۲.۳)$$

پیش از اثبات این قضیه، نخست لم زیر که به لم ماسارت معرف است را ارایه می کنیم.

لم ۴.۳ (لم ماسارت).

فرض کنید $A \subset \mathbb{R}^m$ مجموعه متناهی از بردارها باشد به گونه ای که برای هر بردار $\mathbf{a} \in A$ ، نابرابری $\|\mathbf{a}\| \leq 1$ برقرار باشد. آنگاه نابرابری زیر برقرار است.

$$\mathbb{E}_{\sigma} \left[\max_{\mathbf{a} \in A} \sum_{i=1}^m \sigma_i a_i \right] \leq \sqrt{2 \ln |A|} \quad (۸۳.۳)$$

که σ_i متغیرهای تصادفی رادمیچر و a_1, \dots, a_m مولفه های بردار \mathbf{a} هستند.

اثبات قضیه ۸.۳ اگر فضای A را به صورت زیر تعریف کنیم.

$$A = \left\{ \frac{1}{\sqrt{m}} (h(x_1), h(x_2), \dots, h(x_m)) \right\}.$$

آنگاه A مجموعه ای از بردارها در فضای \mathbb{R}^m است که برای همه $\mathbf{a} \in A$ داریم $\|\mathbf{a}\| = 1$. لذا از تعریف پیچیدگی تجربی رادمیچر داریم

$$\hat{\mathcal{R}}_S(H) = \mathbb{E}_{\sigma} \left[\sup_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) \right] \quad (۸۴.۳)$$

$$= \mathbb{E}_{\sigma} \left[\sup_{\mathbf{a} \in A} \frac{\sqrt{m}}{m} \sum_{i=1}^m \sigma_i a_i \right] \quad (۸۵.۳)$$

$$= \frac{1}{\sqrt{m}} \mathbb{E}_{\sigma} \left[\max_{\mathbf{a} \in A} \sum_{i=1}^m \sigma_i a_i \right] \quad (۸۶.۳)$$

$$\leq \frac{1}{\sqrt{m}} \sqrt{2 \ln |A|} \quad (۸۷.۳)$$

$$= \sqrt{\frac{2 \ln |A|}{m}}. \quad (۸۸.۳)$$

نابرابری بالا براساس لم ماسارت بدست آمده است. با توجه به اینکه مجموعه A تمام دسته بندی های ممکن را برای مجموعه S را ارایه می

نماید لذا $A \subseteq H$ و در نتیجه $|A| \leq |H|$ را داریم و بنابراین قضیه اثبات می گردد. ■

مشکل کران بدست آمده در قضیه بالا این است که قید $|H| < \infty$ وجود دارد و برای بسیاری از فضاهای فرضیه کارایی ندارد. از

آنجایی که $\hat{\mathcal{R}}_S(H)$ علاوه بر H تنها به S وابسته است می خواهیم مشخص نماییم که چگونه رفتار همه فرضیه های این فضا به مجموعه آموزشی وابسته هستند. در فضایی زیر این قید حذف شده و کران های سفت تری بدست می آوریم.

قضیه ۹.۳ (۰).

برای هر فضای فرضیه H داریم

$$\hat{\mathcal{R}}_S(H) \leq \sqrt{\frac{2 \ln \Pi_H(S)}{m}} \leq \sqrt{\frac{2 \ln \Pi_H(m)}{m}}. \quad (۸۹.۳)$$

برهان در اثبات این قضیه می‌خواهیم رفتار فرضیه‌ها به مجموعه آموزشی را مشخص نماییم. برای هر مجموعه آموزشی، فضای فرضیه H' را زیر مجموعه ای H تعریف نموده که برای هر رفتار S یک نماینده از H در نظر می‌گیرد. دقت شود که H' زیر مجموعه ای از H است و دارای ویژگی زیر است.

$$|H'| = \Pi_H(S) \leq \Pi_H(m) \leq 2^m < \infty. \quad (90.3)$$

حال براساس تعریف پیچیدگی تجربی رادمیچر داریم

$$\hat{\mathcal{R}}_S(H) = \mathbb{E}_\sigma \left[\sup_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) \right] \quad (91.3)$$

از آنجایی که برای هر $h \in H$ که مقدار $\hat{\mathcal{R}}_S(H)$ را بیشینه نماید یک $h' \in H'$ وجود دارد که همان مقدار را نتیجه می‌دهد لذا داریم

$$\hat{\mathcal{R}}_S(H) = \mathbb{E}_\sigma \left[\sup_{h' \in H'} \frac{1}{m} \sum_{i=1}^m \sigma_i h'(x_i) \right] \quad (92.3)$$

$$= \hat{\mathcal{R}}_S(H'). \quad (93.3)$$

رابطه بالا نشان می‌دهد که \sup روی H بیشتر از \sup روی H' نخواهد بود و همچنین \sup روی H' بیشتر از \sup روی H نخواهد بود و بنابراین هر دو \sup با هم برابر خواهند بود لذا از قضیه ۸.۳ داریم

$$\hat{\mathcal{R}}_S(H) = \hat{\mathcal{R}}_S(H') \quad (94.3)$$

$$\leq \sqrt{\frac{2 \ln |H'|}{m}} \quad (95.3)$$

$$= \sqrt{\frac{2 \ln \Pi_H(S)}{m}} \quad (96.3)$$

و به این ترتیب اثبات قضیه کامل می‌گردد.

حال رابطه پیچیدگی تجربی رادمیچر را با بعد VC بررسی می‌کنیم.

قضیه ۱۰.۳ ()

اگر $VC(H) = d$ باشد آنگاه برای هر $m \geq d \geq 1$ داریم

$$\hat{\mathcal{R}}_S(H) \leq \sqrt{\frac{2d \ln \left(\frac{em}{d}\right)}{m}} \quad (97.3)$$

برهان از لم ۲.۳ داشتیم

$$\Pi_H(m) \leq \left(\frac{em}{d}\right)^d. \quad (98.3)$$

با جایگزینی نابرابری بالا در نتیجه قضیه ۹.۳ داریم

$$\hat{\mathcal{R}}_S(H) \leq \sqrt{\frac{2 \ln \Pi_H(m)}{m}} \quad (99.3)$$

$$\leq \sqrt{\frac{2 \ln \left(\frac{em}{d}\right)^d}{m}} \quad (100.3)$$

$$= \sqrt{\frac{2d \ln \left(\frac{em}{d}\right)}{m}} \quad (101.3)$$

$$= \sqrt{\frac{2 \ln \left(\frac{em}{d}\right)}{\left(\frac{m}{d}\right)}}. \quad (102.3)$$

و بدین ترتیب اثبات قضیه کامل می شود.

۴.۳ یادگیرهای عمومی

پیش از این دیدیم که داده های آموزشی می توانند یادگیر را بد راهنمایی کنند و منجر به بیش برآزش شود. یک راه برای حل این مشکل محدود کردن فضای جستجو به یک فضای فرضیه H است. این فضای فرضیه را می توان به صورت دانش پیشین دانست که در چنین فضایی یک فرضیه با خطای کم موجود است. آیا وجود چنین دانش پیشینی برای موفقیت یادگیر لازم است؟ یا به عبارتی دیگر می توان بدون دانش پیشین، یک یادگیر موفق داشت؟ به این یادگیر که از هیچ گونه دانش پیشینی استفاده نمی کند یادگیر عمومی می گویند.

در واقع پرسش را می توان این گونه مطرح نمود آیا برای یک مجموعه آموزشی با اندازه m ، یک الگوریتم یادگیری A وجود دارد که برای هر توزیع D ، اگر الگوریتم A مجموعه آموزشی با m نمونه که از روی توزیع D به صورت مستقل نمونه برداری شده اند دریافت و با احتمال بالا فرضیه h را تولید نماید که خطای کمی را تولید کند؟

قضیه No free lunch بیان می دارد که چنین یادگیر عمومی وجود ندارد. یا به بیان دقیق تر برای هر مساله دسته بندی، برای هر یادگیر یک توزیع احتمال وجود دارد که یادگیر روی آن توزیع توانایی یادگیری را نداشته و شکست می خورد. واژه شکست به این معنا است که یادگیر با دریافت مجموعه آموزشی با ویژگی های یاد شده، فرضیه ای را تولید می نماید که خطای زیادی (برای نمونه خطای $2/3$) دارد در حالیکه برای این توزیع یادگیر دیگری وجود دارد که فرضیه ای با خطای کم را تولید می نماید. بنابراین در یک مساله یادگیری که نمونه ها با توزیع D نمونه برداری می شوند به دانش پیشین نیاز خواهیم داشت. این دانش پیشین می تواند گونه های مختلف زیر را داشته باشد.

- فرض کنیم که توزیع D یک شکل پارامتری داشته که پارامترهای آن نامشخص است. بنابراین می توانیم نخست پارامترهای توزیع را بدست آورده و سپس دسته بند را طراحی نماییم. به این گونه مدل ها، مدل های مولد می گویند.
- فرض کنیم یک فرضیه $h \in H$ وجود دارد به گونه ای که $\mathbf{R}(h) = 0$. یک قید نرم تر روی دانش پیشین این است که $\min_{h \in H} \mathbf{R}(h)$ کوچک باشد.

قضیه ۱۱.۳ (No free lunch)

فرض کنید که A یک الگوریتم یادگیری برای دسته بندی دودویی با تابع خطای $1/2$ در دامنه \mathcal{X} باشد. همچنین فرض کنید که m عددی کوچک تر از $\frac{1}{4}$ باشد. آنگاه یک توزیع D روی $\{0, 1\} \times \mathcal{X}$ وجود دارد به گونه ای که

۱. یک فرضیه $h : \mathcal{X} \rightarrow \{0, 1\}$ وجود دارد که $\mathbf{R}(h) = 0$
۲. با احتمال دست کم $\frac{1}{4}$ روی مجموعه آموزشی S با اندازه m که از روی توزیع D نمونه برداری شده است داریم $\mathbf{R}(A(S)) \geq \frac{1}{8}$.

شیوه اثبات این قضیه به این صورت است که یک توزیع D را به گونه ای می سازیم که خطای فرضیه تولید شده زیاد باشد. نکته ای

که در این قضیه وجود دارد این است که هر الگوریتم که تنها نیمی از داده های فضا را ببیند نمی تواند در باره شیوه برچسب زنی داده ها داوری نماید. رابطه قضیه NFL با دانش پیشین چیست؟ برای بررسی این پرسش، فرض کنید $H = \{h \mid h : \mathcal{X} \mapsto \{0, 1\}\}$ مجموعه تمام توابع A و یک الگوریتم کمینه سازی خطای تجربی باشد. در واقع این نشان می دهد که هیچ گونه دانش پیشین نداریم. براین اساس نتیجه زیر حاصل می شود.

نتیجه ۴.۳.

فرض کنید \mathcal{X} فضای نمونه نامتناهی و $H = \{h \mid \mathcal{X} \mapsto \{0, 1\}\}$ مجموعه تمامی فرضیه های ممکن باشد آنگاه H قابلیت یادگیری احتمالا تقریبا درست را ندارد.

حال این پرسش مطرح می شود که چگونه می توان جلوی شکست الگوریتم یادگیری را گرفت؟ یک راه برای جلوگیری شکست الگوریتم یادگیری، محدود کردن فضای فرضیه H است. اما چگونه می توان فضای فرضیه خوب انتخاب نمود؟ پیش از این دیدیم که انتخاب فضای فرضیه پیچیده راه حل مناسب نیست زیرا بیش برآزش روی می دهد. از سویی دیگر، انتخاب فضای فرضیه ساده هم راه حل مناسب نیست چون خطای آموزشی زیادی دارد. بنابراین یک مصالحه وجود دارد. برای پاسخ به پرسش یاد شده، خطای الگوریتم کمینه سازی خطا را به دو بخش تجزیه می کنیم. برای این تجزیه خطا، فرض کنید الگوریتم کمینه سازی خطای تجربی با دریافت مجموعه آموزشی S فرضیه H را تولید می کند. بنابراین خطای واقعی h را می توان به صورت زیر نوشت.

$$\mathbf{R}(h) = \mathbf{R}_{app}(h) + \mathbf{R}_{est}(h) \quad (۱۰۳.۳)$$

که $\mathbf{R}_{app}(h) = \min_{h \in H} \mathbf{R}(h)$ خطای تقریب و $\mathbf{R}_{est}(h) = \mathbf{R}(h) - \mathbf{R}_{app}(h)$ خطای برآورد نامیده می شود. خطای تقریب که عبارت است از کمترین خطا در فضای فرضیه H و براساس فضای فرضیه انتخاب شده مشخص می گردد و به مجموعه آموزشی وابسته نیست. خطای تقریب نشان می دهد چه مقدار از خطا براساس فرض های صورت گرفته برای محدود کردن به فضای فرضیه انتخاب شده ایجاد می شود. خطای برآورد عبارت است از اختلاف بین خطای تولید شده و خطای بهترین فرضیه در فضای فرضیه انتخاب شده. این خطا هم به مجموعه آموزشی و اندازه آن وابسته است و هم به پیچیدگی فضای فرضیه انتخاب شده. پیش از این نشان داده شد که این خطا برای فضای فرضیه متناهی به صورت لگاریتمی با $|H|$ افزایش و با افزایش m کاهش می یابد.

بهترین فضای فرضیه، فضایی است که تنها یک دسته بند دارد و آن دسته بند بهینه بیز می باشد. این دسته بند به توزیع D وابسته است که آن را نمی دانیم. یکی از اهداف نظریه یادگیری این است که چگونه بتوان یک فضای فرضیه غنی و در عین حال خطای برآورد را در حد معقولی داشت.

یادگیری نایکنواخت

نماد یادگیری احتمالا تقریبا درست که پیش از این در باره آن سخن گفتیم اجازه می‌دهد که اندازه مجموعه آموزشی به پارامترهای دقت و اطمینان وابسته باشد اما این اندازه مجموعه آموزشی برای همه فرضیه‌ها و همه توزیع‌های داده برقرار است. کلاس‌های فرضیه‌هایی که در این مدل یادگرفتنی هستند محدود است و این کلاس‌ها می‌بایست دارای بعد VC متناهی باشند. در این بخش، یک مدل ضعیف تر و سست شده قابلیت یادگیری را مطرح می‌کنیم و ویژگی‌های این مدل یادگیری که یادگیری نایکنواخت نام دارد را بررسی می‌کنیم.

در ادامه، نخست مدل یادگیری نایکنواخت بررسی می‌شود که در این مدل، اندازه مجموعه آموزشی علاوه بر پارامترهای اطمینان و دقت به فرضیه نیز وابسته است. سپس ویژگی‌های این مدل را بررسی و نشان می‌دهیم که مدل نسخه سست شده مدل یادگیری احتمالا تقریبا درست بدون پیش فرض است و در ادامه شرط لازم برای اینکه فضای فرضیه به صورت نایکنواخت یادگرفتنی باشد را بیان می‌کنیم و یک الگوریتم یادگیری با نام کمینه سازی خطای ساختاری بیان می‌شود که توانایی یادگیری نایکنواخت را دارا است. همچنین روش یادگیری کمینه طول توصیف برای فضای فرضیه شمارا بیان می‌شود. در پایان مدل یادگیری سازگاری را که مدل ضعیف تر یادگیری است ارایه و برتری‌ها و کاستی‌های مدل‌های مختلف یادگیری که در فصل‌های پیشین ارایه شد را بررسی می‌کنیم.

۱.۴ قابلیت یادگیری نایکنواخت

پیش از این در باره همگرایی یکنواخت سخن گفتیم. واژه یکنواخت به این معنا به کار برده شد که با یک مجموعه آموزشی با اندازه ثابت و برای همه فرضیه‌های درون مجموعه فرضیه‌ها و برای همه توزیع‌ها، می‌توانستیم کران خطا یا پیچیدگی نمونه‌ای را تعیین نماییم. به بیانی دیگر همگرایی یکنواخت را می‌توان به صورت رسمی به صورت زیر بیان نمود.

تعریف ۱.۴ همگرایی یکنواخت کلاس فرضیه H ویژگی همگرایی یکنواخت را دارد اگر یک تابع $m_H^{UC} : (0, 1)^2 \mapsto \mathbb{N}$ و الگوریتم یادگیری A وجود داشته باشد به گونه‌ای که برای هر $\epsilon, \delta \in (0, 1)$ و برای هر توزیع D روی $\mathcal{X} \times \mathcal{Y}$ اگر مجموعه آموزشی با اندازه $m > m_H^{UC}(\epsilon, \delta)$ به صورت مستقل از توزیع D از روی فضای نمونه، نمونه برداری شده باشند. آنگاه با احتمال دست کم $1 - \delta$ نایبربری زیر برای همه فرضیه‌های درون H برقرار است.

$$\mathbf{R}(A(S)) \leq \mathbf{R}(h) + \epsilon. \quad (1.4)$$

فرض کنید که می‌خواهیم یک دسته‌بند طراحی نماییم که با احتمال بالا دارای خطای کمی باشد. اگر یک دسته‌بند خطی را در نظر بگیریم ممکن است با احتمال بالا خطای آن کم نباشد. حال اگر دسته‌بند یک چندجمله‌ای از درجه ۲ باشد ممکن است خطای مناسب را نداشته باشد. اگر فضای فرضیه‌ها شامل همه چندجمله‌ای‌ها باشد به احتمال بسیار زیاد دسته‌بندی پیدا می‌کنیم که خطای آن کم است و لذا مناسب خواهد بود. اما بعد VC این فضا برابر ∞ است و بنابراین کران‌هایی که پیش از این بدست آوردیم دیگر ارزشی ندارد. همان گونه که پیش از این بیان شد کلاس فرضیه‌هایی که قابلیت یادگیری با ویژگی همگرایی یکنواخت را داشته باشند بسیار محدود است و این کلاس فرضیه‌ها باید دارای بعد VC متناهی باشند. در ادامه شرایط بیان شده بالا را تعدیل می‌کنیم و اجازه می‌دهیم که اندازه مجموعه آموزشی علاوه بر ϵ و δ به h نیز وابسته باشد. به بیانی دیگر، اندازه مجموعه آموزشی به صورت یکنواخت برای همه فرضیه‌ها یکسان نیست. به این ویژگی، ویژگی همگرایی نایکنواخت یا یادگیری نایکنواخت می‌گوییم و به صورت زیر تعریف می‌شود.

تعریف ۲.۴ یادگیری نایکنواخت کلاس فرضیه H قابلیت یادگیری نایکنواخت را دارد اگر یک تابع $m_H^{NUL} : (0, 1)^2 \times H \mapsto \mathbb{N}$ و یک الگوریتم یادگیری A وجود داشته باشد به گونه‌ای که برای هر $\epsilon, \delta \in (0, 1)$ و برای هر $h \in H$ اگر مجموعه آموزشی با اندازه $m > m_H^{NUL}(\epsilon, \delta, h)$ به صورت مستقل از توزیع D از روی فضای نمونه، نمونه برداری شده باشند. آنگاه با احتمال دست کم $1 - \delta$ نایبربری زیر برای همه فرضیه‌های درون H برقرار است.

$$\mathbf{R}(A(S)) \leq \mathbf{R}(h) + \epsilon. \quad (2.4)$$

هرچه فضای H بزرگتر باشد فرضیه‌های این فضا انعطاف پذیرتر و در نتیجه هزینه این انعطاف پذیری بزرگ شدن مجموعه آموزشی برای یادگیری آن خواهد شد. در یادگیری نایکنواخت فضای فرضیه گروه‌بندی می‌شود بطوریکه هرچه یک گروه پیچیده‌تر باشد پیچیدگی نمونه‌ای یادگیری از آن گروه بزرگتر خواهد شد.

قضیه ۱.۴ (یادگیری نایکنواخت).

کلاس فرضیه H شامل دسته بندهای $h : \mathcal{X} \mapsto \{0, 1\}$ قابلیت یادگیری نایکنواخت را دارد اگر و تنها اگر از اجتماع شمارای کلاس فرضیه‌هایی که قابلیت یادگیری احتمالا تقریباً درست بدون پیش فرض را داشته باشند تشکیل شده باشد.

اثبات قضیه بالا از قضیه زیر نتیجه گیری می‌شود.

قضیه ۲.۴ (یادگیری نایکنواخت).

اگر $H = \bigcup_{n \in \mathbb{N}} H_n$ باشد و هر کدام از کلاس‌های فرضیه H_n دارای ویژگی همگرایی یکنواخت باشند آنگاه کلاس فرضیه H یادگیری نایکنواخت است.

اثبات قضیه ۱.۴ برای اثبات بخش اگر، فرض کنید که $H = \bigcup_{n \in \mathbb{N}} H_n$ تعریف شده باشد که هر کدام از H_n ها دارای قابلیت یادگیری احتمالا تقریبا درست بدون پیش فرض باشد. از قضیه مهم یادگیری داشتیم که هر کدام از H_n ها دارای ویژگی همگرایی یکنواخت می باشند. بر اساس قضیه ۲.۴ نتیجه می گیریم که H دارای ویژگی یادگیری نا یکنواخت است.

برای اثبات بخش تنها اگر، فرض کنید که H با استفاده از الگوریتم A یادگیر یکنواخت باشد. آنگاه برای هر $n \in \mathbb{N}$ ، فرض کنید که H_n را به صورت زیر تعریف کنیم.

$$H_n = \left\{ h \in H \mid m_H^{NUL} \left(\frac{1}{\lambda}, \frac{1}{\sqrt{V}}, h \right) \leq n \right\}. \quad (۳.۴)$$

روشن است که $H = \bigcup_{n \in \mathbb{N}} H_n$. به علاوه بر اساس تعریف m_H^{NUL} ، می دانیم که برای هر توزیع \mathcal{D} که شرط تحقق پذیری^۱ را بر اساس H_n ارضا نماید با احتمال دست کم $\frac{\epsilon}{V}$ روی $\mathcal{D}^m \sim S$ داریم $\mathbf{R}(A(S)) \leq \frac{1}{\lambda}$. بر اساس قضیه اساسی یادگیری، این نتیجه می دهد که بعد $VC(H_n)$ باید متناهی باشد و در نتیجه H_n قابلیت یادگیری احتمالا تقریبا درست بدون پیش فرض را دارد.

مثال زیر نشان می دهد که یادگیری نایکنواخت نسخه راحت شده یادگیری احتمالا تقریبا درست بدون پیش فرض است. یعنی فضاهای فرضیه ای وجود دارند که قابلیت یادگیری نایکنواخت را دارند اما قابلیت یادگیری احتمالا تقریبا درست بدون پیش فرض را ندارند.

مثال ۱.۴ (قابلیت یادگیری نایکنواخت خانواده چند جمله ای ها).

مساله دسته بندی دودویی با فضای نمونه $X = \mathbb{R}$ را در نظر بگیرید که برای هر $n \in \mathbb{N}$ ، فرضیه H_n چند جمله ای از درجه n باشد یعنی

$$H_n = \{h_p \in H \mid h_p(x) = \text{sgn}(P_n(x))\}. \quad (۴.۴)$$

که $P_n(x)$ چند جمله ای از درجه n است. فرض کنید که فضای فرضیه به صورت $H = \bigcup_{n \in \mathbb{N}} H_n$ باشد. آنگاه H مجموعه تمام چند جمله ای های ممکن است. به راحتی می توان نشان داد که $VC(H_n) = n + 1$ و $VC(H) = \infty$ است. در نتیجه H دارای قابلیت یادگیری احتمالا تقریبا درست بدون پیش فرض نیست اما بر اساس قضیه ۱.۴ مجموعه فرضیه های H دارای قابلیت یادگیری نایکنواخت را دارد.

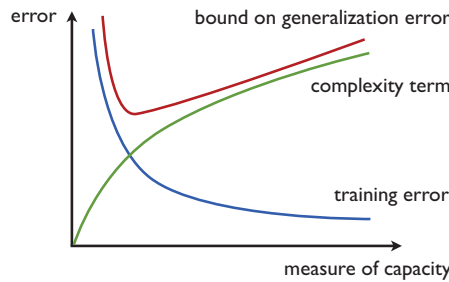
حال این پرسش مطرح است که آیا الگوریتمی برای یادگیری نایکنواخت وجود دارد؟ در ادامه به بررسی چنین الگوریتمی می پردازیم.

۲.۴ کمینه سازی هزینه ساختاری

پیش از این دیدیم که با کمک دانش پیشین می توانیم کلاس فرضیه ها را پیدا کنیم یعنی در واقع امیدواریم که H دارای فرضیه ای باشد که برای مساله داده شده دارای خطای کمی باشد. یک روش دیگر برای به کار بردن دانش پیشین، مشخص نمودن ترجیح در کلاس فرضیه های H است که در کمینه سازی هزینه ساختاری به کار می رود. روش کمینه سازی هزینه ساختاری، دنباله ای از فرضیه ها را به صورت افزایشی اندازه آنها به صورت زیر در نظر می گیرید.

$$H_0 \subseteq H_1 \subseteq H_2 \subseteq \dots \subseteq H_n \subseteq \dots \quad (۵.۴)$$

^۱Realizability



شکل ۱.۴ کمینه سازی هزینه ساختاری

و تلاش می کند پاسخ خطای تجربی کمینه برای هر H_n را پیدا نماید و سپس از این کلاس های مختلف یک فرضیه را پیدا و به عنوان خروجی اعلام نماید. فرضیه انتخاب شده فرضیه ای است که تابع هدف زیر را که به صورت مجموع خطای تجربی و پیچیدگی فرضیه تعریف می شود کمینه باشد.

$$h_m^{SRM} = \operatorname{argmin}_{h \in H_n, n \in \mathbb{N}} [\hat{\mathbf{R}}(h) + \text{Complexity}(H_n, m)] \quad (۶.۴)$$

که تابع $\text{Complexity}(H_n, m)$ به قدرت توصیف فضای فرضیه که ظرفیت یا غنای فضای فرضیه را مشخص می نماید و اندازه مجموعه آموزشی وابسته است. شکل ۱.۴ خطای کمینه سازی هزینه ساختاری را نشان می دهد.

علیرغم اینکه روش کمینه سازی هزینه ساختاری از پشتوانه نظری بسیار غنی برخوردار است اما دارای هزینه محاسباتی بسیار زیادی است و می بایست الگوریتم کمینه سازی هزینه تجربی را برای هر خانواده از H_n اجرا نماید. آیا تعداد مسایلی که باید توسط الگوریتم کمینه سازی خطا حل نماییم بی نهایت است؟ این تعداد مسایل بینهایت نیست زیرا اگر برای یک $n > n_0$ خطای تجربی برابر صفر باشد دیگر نیاز به حل آن مسایل نیست و در نتیجه تعداد این مسایل بی نهایت نیست.

همان گونه که پیش از این بیان شد، از دانش پیشین می توانیم برای مشخص نمودن ترجیح در کلاس فرضیه های H استفاده نماییم. در کمینه سازی هزینه ساختاری، در ابتدا فرض می کنیم که H را می توان به صورت $H = \bigcup_{n \in \mathbb{N}} H_n$ نوشت و سپس تابع وزنی را مشخص نمود که $w : \mathbb{N} \rightarrow [0, 1]$ به هر کلاس فرضیه H_n داده می شود به گونه ای که مقدار بیشتر وزن به معنای ترجیح بیشتر است. برای هر $n \in \mathbb{N}$ فرض می کنیم که H_n دارای ویژگی همگرایی یکنواخت است و پیچیدگی نمونه ای آن برابر $m_{H_n}^{UC}(\epsilon, \delta)$ است. در ادامه بیان می کنیم که چگونه می توانیم چنین دانش پیشینی را برای یادگیری به کار ببریم. اگر تابع $\epsilon_n : \mathbb{N} \times (0, 1) \rightarrow (0, 1)$ را به صورت زیر تعریف نماییم

$$\epsilon_n(m, \delta) = \min\{\epsilon \mid m_{H_n}^{UC}(\epsilon, \delta) \leq m\}. \quad (۷.۴)$$

یعنی برای هر مجموعه آموزشی با طول ثابت m ، ما به کمترین کران بین خطای واقعی و خطای تجربی علاقه مند هستیم. از تعریف همگرایی یکنواخت و ϵ_n می توان نتیجه گرفت که برای هر m و δ با احتمال دست کم $1 - \delta$ روی انتخاب $S \sim \mathcal{D}^m$ نابرابری زیر برای همه $h \in H_n$ برقرار است.

$$|\mathbf{R}(h) - \hat{\mathbf{R}}(h)| \leq \epsilon_n(m, \delta). \quad (۸.۴)$$

فرض کنید که $w : \mathbb{N} \rightarrow [0, 1]$ تابع وزن روی کلاس فرضیه ها باشد به گونه ای که $\sum_{n=1}^{\infty} w(n) \leq 1$. این تابع وزن اهمیت یا پیچیدگی کلاس فرضیه را نشان می دهد. اگر H مجموعه متناهی از M فضای فرضیه H_1, H_2, \dots, H_M باشد آنگاه می توانیم به هر کدام از کلاس

ها وزن $\frac{1}{M}$ تخصیص بدهیم. این وزن یکنواخت نشان دهنده نداشتن دانش پیشین (ترجیح) است. اما اگر H از اجتماع شمارای کلاس فرضیه های H_n تشکیل شده باشد. آنگاه وزن یکنواخت امکان پذیر نیست اما روش های دیگر وزن دهی امکان پذیر است. برای نمونه توابع زیر دو نمونه از توابع وزن ممکن هستند.

$$w(n) = \frac{6}{(\pi n)^2} \quad (۹.۴)$$

$$w(n) = 2^{-n} \quad (۱۰.۴)$$

الگوریتم کمینه سازی ساختاری تلاش می کند کران را کمینه نماید. یعنی هدف پیدا کردن فرضیه ای است که کران بالای خطای واقعی را کمینه نماید. در ادامه روش سراسری را جهت تعریف تابع وزن معرفی می نماییم.

قضیه ۳.۴ (کران خطا برای کمینه سازی هزینه ساختاری).

فرض کنید $w : \mathbb{N} \rightarrow [0, 1]$ یک تابع وزن باشد به گونه ای که $\sum_{n=1}^{\infty} w(n) \leq 1$ و همچنین فرض کنید H روی کلاس فرضیه ای باشد که بتوان به صورت $H = \bigcup_{n \in \mathbb{N}} H_n$ نوشت که هر کدام از H_n ها دارای ویژگی همگرایی یکنواخت با پیچیدگی نمونه ای $m_{H_n}^{UC}$ باشد اگر ϵ_n به صورت $\epsilon_n(m, \delta) = \min\{\epsilon \mid m_{H_n}^{UC}(\epsilon, \delta) \leq m\}$ تعریف شده باشد. آنگاه برای هر $\epsilon, \delta \in (0, 1)$ و هر توزیع \mathcal{D} با احتمال دست کم $1 - \delta$ برای مجموعه آموزشی $S \sim \mathcal{D}^m$ نابرابری زیر برای همه $n \in \mathbb{N}$ و همه $h \in H_n$ برقرار است.

$$|\mathbf{R}(h) - \hat{\mathbf{R}}(h)| \leq \epsilon_n(m, w(n) \times \delta) \quad (۱۱.۴)$$

و به بیانی دیگر برای هر $\delta \in (0, 1)$ و توزیع \mathcal{D} ، با احتمال دست کم $1 - \delta$ نابرابری زیر برای همه $h \in H$ برقرار است.

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \min_{h \in H_n, n \in \mathbb{N}} \epsilon_n(m, w(n) \times \delta) \quad (۱۲.۴)$$

برهان برای هر n مقدار δ_n را به صورت $\delta_n = w(n) \times \delta$ تعریف می کنیم. با توجه به اینکه ویژگی همگرایی یکنواخت برای همه n ها مطابق رابطه (۸.۴) برقرار است. اگر n را ثابت در نظر بگیریم لذا با احتمال دست کم $1 - \delta$ روی انتخاب $S \sim \mathcal{D}^m$ نابرابری زیر برای همه $h \in H_n$ برقرار است.

$$|\mathbf{R}(h) - \hat{\mathbf{R}}(h)| \leq \epsilon_n(m, \delta_n) \quad (۱۳.۴)$$

با بکارگیری کران اجتماعات^۲ روی $n = 1, 2, 3, \dots$ با احتمال دست کم $1 - \delta \sum_n w(n) = 1 - \delta$ نابرابری زیر برای همه n ها و همه $h \in H_n$ برقرار است.

$$|\mathbf{R}(h) - \hat{\mathbf{R}}(h)| \leq \epsilon_n(m, \delta_n) \quad (۱۴.۴)$$

و بدین ترتیب اثبات قضیه کامل می شود.

اگر $n(h)$ را به صورت زیر تعریف کنیم

$$n(h) = \min\{n \mid h \in H_n\} \quad (۱۵.۴)$$

^۲Union bound

Structural Risk Minimization (SRM)**prior knowledge:**

$\mathcal{H} = \bigcup_n \mathcal{H}_n$ where \mathcal{H}_n has uniform convergence with $m_{\mathcal{H}_n}^{UC}$

$w : \mathbb{N} \rightarrow [0, 1]$ where $\sum_n w(n) \leq 1$

define: ϵ_n as in Equation (7.1); $n(h)$ as in Equation (7.4)

input: training set $S \sim \mathcal{D}^m$, confidence δ

output: $h \in \operatorname{argmin}_{h \in \mathcal{H}} [L_S(h) + \epsilon_{n(h)}(m, w(n(h))) \cdot \delta]$

شکل ۲.۴ الگوریتم کمینه سازی هزینه ساختاری

آنگاه رابطه (۱۲.۴) نتیجه زیر را می دهد.

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \epsilon_{n(h)}(m, w(n(h))) \times \delta \quad (16.4)$$

روش کمینه سازی هزینه ساختاری فضا را برای پیدا کردن فرضیه ای که کران بالا را کمینه نماید جستجو می نماید و این روش در الگوریتم ۲.۴ ارایه شده است.

همچنین می توان نشان داد که روش کمینه سازی هزینه ساختاری می تواند برای یادگیری نایکنواخت هر فضای فرضیه ای که از اجتماع شمارای فضای فرضیه هایی که ویژگی همگرایی یکنواخت را دارند به کار رود.

قضیه ۴.۴ (قابلیت یادگیری نایکنواخت با روش کمینه سازی هزینه ساختاری).

فرض کنید $H = \bigcup_{n \in \mathbb{N}} H_n$ بطوریکه H_n دارای ویژگی همگرایی یکنواخت با پیچیدگی نمونه ای $m_{H_n}^{UC}$ باشد. همچنین تابع وزن $w : \mathbb{N} \rightarrow [0, 1]$ به صورت $w(n) = \frac{\epsilon}{\pi^2 n^2}$ تعریف شده باشد. آنگاه H قابلیت یادگیری نایکنواخت توسط روش کمینه سازی هزینه ساختاری با نرخ زیر را دارد.

$$m_H^{NUL}(\epsilon, \delta, h) \leq m_{H_{n(h)}}^{UC} \left(\frac{\epsilon}{2}, \frac{\epsilon \delta}{(\pi n(h))^2} \right) \quad (17.4)$$

برهان فرض کنید الگوریتم A ، یک الگوریتم کمینه سازی هزینه ساختاری با تابع وزن $w(n)$ باشد. آنگاه برای هر $\epsilon, h \in H$ و δ فرض کنید نابرابری $(\frac{\epsilon}{2}, w(n(h))\delta)$ برقرار باشد. از آنجایی که $\sum_n w(n) = 1$ ، با استفاده از قضیه ۳.۴، با احتمال دست کم $1 - \delta$ برای مجموعه آموزشی $S \sim \mathcal{D}^m$ نابرابری زیر برای همه فرضیه های $h \in H$ برقرار است.

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \epsilon_{n(h)}(m, w(n(h)))\delta. \quad (18.4)$$

نابرابری بالا برای فرضیه $A(S)$ که توسط الگوریتم کمینه سازی هزینه ساختاری تولید می شود برقرار است. از تعریف کمینه سازی هزینه ساختاری می توانیم نابرابری زیر را بدست آوریم.

$$\mathbf{R}(A(S)) \leq \min_h \hat{\mathbf{R}}(h') + \epsilon_{n(h')} (m, w(n(h'))\delta) \quad (19.4)$$

$$\leq \hat{\mathbf{R}}(h) + \epsilon_{n(h)}(m, w(n(h)))\delta \quad (20.4)$$

در پایان اگر $m \geq m_{H_{n(h)}}^{UC}(\epsilon/2, w(n(h))\delta)$ باشد نابرابری $\epsilon_{n(h)}(m, w(n(h)))\delta \leq \epsilon/2$ هم برقرار است. از طرفی ویژگی همگرایی یکنواخت، هر H_n با احتمال دست کم $1 - \delta$ نابرابری $\hat{\mathbf{R}}(h) \leq \mathbf{R}(h) + \epsilon/2$ برقرار است با ترکیب این نابرابری با نابرابری (۲۰.۴)،

نابرابری زیر بدست می آید

$$\mathbf{R}(A(S)) \leq \mathbf{R}(h) + \epsilon. \quad (21.4)$$

که اثبات قضیه تکمیل می گردد.

پیش از این نشان دادیم که هر فضای فرضیه که از اجتماع شمارای مجموعه فرضیه هایی با بعد VC متناهی تشکیل شده باشند قابلیت یادگیری نایکخواخت را دارد. می توان نشان داد که برای هر فضای نامتناهی \mathcal{X} ، مجموعه همه دسته بند های دو دسته ای روی فضای \mathcal{X} را نمی توان اجتماع شمارای مجموعه فرضیه هایی با بعد VC متناهی تشکیل داد. این بدان معنا است که قضیه No free lunch برای یادگیری نایکخواخت نیز برقرار است. هنگامی که فضای نمونه نامتناهی است هیچ یادگیر نایکخواختی برای مجموعه همه دسته بندهای ممکن دو دسته ای وجود ندارد. هرچند برای هریک از این دسته بندها، یک الگوریتم یادگیری وجود دارد.

نکته دیگری که در یادگیری نایکخواخت وجود دارد این است که نقش دانش پیشین برای یادگیری نایکخواخت چیست و آیا این دانش پیشین سبب افزایش پیچیدگی نمونه ای می شود یا کاهش آن؟ در قابلیت یادگیری یکنواخت، معمولاً دانش پیشین به صورت انتخاب فضای فرضیه در الگوریتم یادگیری وارد می شود در حالیکه در یادگیری نایکخواخت این دانش پیشین به صورت ترجیح یا تابع وزن در الگوریتم یادگیری وارد می شود. بدیهی است این نوع دانش پیشین ضعیف تر از گونه ای است که در قابلیت یادگیری یکنواخت است و هزینه این ضعیف کردن دانش پیشین، افزایش پیچیدگی نمونه ای هر کدام از فضاها H_n است. اگر دانش پیشین در باره n داشته باشیم می توانیم الگوریتم کمینه سازی خطای تجربی را روی فضای H_n اجرا نماییم که در این حالت پیچیدگی نمونه ای برابر است با $C \frac{n + \log(\frac{1}{\delta})}{\epsilon^2}$ که $VC(H_n) = n$ و C یک مقدار ثابت است. اگر دانش پیشینی در باره n نداشته باشیم الگوریتم کمینه سازی هزینه ساختاری را روی فضای H اجرا نموده و در این حالت با استفاده از قضیه ۴.۴ پیچیدگی نمونه ای یادگیر نایکخواخت برابر است با

$$m_H^{NUL}(\epsilon, \delta, h) \leq m_{H_n}^{UC}(\epsilon/2, w(n)\delta) = O\left(\frac{n + \log\left(\frac{1}{w(n)}\right) + \log\left(\frac{1}{\delta}\right)}{\epsilon^2}\right) \quad (22.4)$$

برای نمونه اگر تابع وزن $w(n) = \frac{1}{n^2}$ در نظر گرفته شود و تفاضل پیچیدگی نمونه ای را محاسبه نماییم خواهیم داشت.

$$\begin{aligned} m_H^{NUL}(\epsilon, \delta, h) - m_{H_n}^{UC}(\epsilon/2, w(n)\delta) &= O\left(\frac{n + \log\left(\frac{1}{w(n)\delta}\right)}{\epsilon^2}\right) - O\left(\frac{n + \log\left(\frac{1}{\delta}\right)}{\epsilon^2}\right) \\ &= O\left(\frac{\log\left(\frac{1}{w(n)}\right)}{\epsilon^2}\right) \end{aligned} \quad (23.4)$$

$$= O\left(\frac{\log n}{\epsilon^2}\right) \quad (24.4)$$

در نتیجه هزینه ضعیف کردن دانش پیشین افزایش پیچیدگی نمونه ای یادگیر است. در نتیجه برای یک مساله مشخص، یادگیر نایکخواخت به تعداد بیشتری نمونه نیاز دارد تا یادگیر یکنواخت. حال فرض کنید که $w : H \mapsto [0, 1]$ تابع وزن به صورت $\sum_{h \in H} w(h) \leq 1$ باشد. در قضیه نشان می دهیم که یادگیری فضای فرضیه H که وزن آن به صورت بالا تعریف شده باشد قابلیت یادگیری نایکخواخت را دارد.

قضیه ۵.۴ (۰)

اگر H فضای فرضیه ها و $w : H \mapsto [0, 1]$ تابع وزن به صورت $\sum_{h \in H} w(h) \leq 1$ باشد. آنگاه برای هر $h \in H$ و هر $\delta > 0$

و هر توزیع \mathcal{D} و برای هر مجموعه آموزشی S با m عضو نابرابری زیر برقرار است.

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \sqrt{\frac{\ln\left(\frac{1}{w(h)}\right) + \ln(2/\delta)}{2m}} \right] \geq 1 - \delta. \quad (25.4)$$

برهان بطور معادل نابرابری زیر را اثبات می کنیم.

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[\exists h \in H \mid |\hat{\mathbf{R}}(h) - \mathbf{R}(h)| \geq \sqrt{\frac{\ln\left(\frac{1}{w(h)}\right) + \ln(2/\delta)}{2m}} \right] \leq \delta. \quad (26.4)$$

اگر برای هر $h \in H$ مقدار ϵ_h را به صورت زیر تعریف کنیم.

$$\epsilon_h = \sqrt{\frac{\ln\left(\frac{1}{w(h)}\right) + \ln(2/\delta)}{2m}}. \quad (27.4)$$

با استفاده از قاعده کران اجتماعات خواهیم داشت

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[\exists h \in H \mid |\hat{\mathbf{R}}(h) - \mathbf{R}(h)| \geq \epsilon_h \right] \leq \sum_{h \in H} \mathbb{P}_{S \sim \mathcal{D}^m} \left[|\hat{\mathbf{R}}(h) - \mathbf{R}(h)| \geq \epsilon_h \mid \cdot \right] \quad (28.4)$$

با استفاده از نابرابری هافدینگ می توانیم کران بالا را محاسبه نماییم.

$$\begin{aligned} \mathbb{P}_{S \sim \mathcal{D}^m} \left[|\hat{\mathbf{R}}(h) - \mathbf{R}(h)| \geq \epsilon_h \right] &\leq 2 \exp[-2m\epsilon_h^2] \\ &= 2 \exp\left[-\ln\left(\frac{2}{w(h)\delta}\right)\right] \\ &= 2 \exp\left[\ln\left(\frac{w(h)\delta}{2}\right)\right] \\ &= w(h)\delta \\ &\leq \delta. \end{aligned} \quad (29.4)$$

با جایگزینی نابرابری بالا در نابرابری (28.4) خواهیم داشت.

$$\begin{aligned} \mathbb{P}_{S \sim \mathcal{D}^m} \left[\exists h \in H \mid |\hat{\mathbf{R}}(h) - \mathbf{R}(h)| \geq \sqrt{\frac{\ln\left(\frac{1}{w(h)}\right) + \ln(2/\delta)}{2m}} \right] &\leq \sum_{h \in H} \mathbb{P}_{S \sim \mathcal{D}^m} \left[|\hat{\mathbf{R}}(h) - \mathbf{R}(h)| \geq \epsilon_h \right] \\ &\leq \delta \sum_{h \in H} w(h) \quad (30.4) \\ &\leq \delta. \quad (31.4) \end{aligned}$$

نابرابری آخر با توجه به قید $\sum_{h \in H} w(h) \leq 1$ نتیجه گرفته شده است.

قضیه ۶.۴.۰)

فرض کنید H فضای فرضیه و $w : H \mapsto [0, 1]$ تابع وزن به صورت $\sum_{h \in H} w(h) \leq 1$ باشد. اگر الگوریتم یادگیری A با دریافت مجموعه آموزشی S ، فرضیه

$$\hat{h} = \operatorname{argmin}_{h \in H} \hat{\mathbf{R}}(h) + \sqrt{\frac{\ln\left(\frac{1}{w(\hat{h})}\right) + \ln(2/\delta)}{2m}} \quad (32.4)$$

را تولید نماید آنگاه برای هر $h \in H$ و هر $\delta > 0$ و هر توزیع D ، اگر مجموعه آموزشی دست کم $m \geq \frac{2 \ln\left(\frac{1}{w(h)}\right) + \ln(2/\delta)}{\epsilon^2}$ نمونه داشته باشد. آنگاه نابرابری زیر برقرار است.

$$\mathbb{P}_{S \sim D^m} [\mathbf{R}(\hat{h}) \geq \mathbf{R}(h) + \epsilon] \leq \delta. \quad (33.4)$$

برهان از قضیه ۵.۴ نتیجه می‌گیریم که با احتمال دست کم $1 - \delta$ دو نابرابری زیر برقرار هستند.

$$\mathbf{R}(\hat{h}) \leq \hat{\mathbf{R}}(\hat{h}) + \sqrt{\frac{\ln\left(\frac{1}{w(\hat{h})}\right) + \ln(2/\delta)}{2m}} \quad (34.4)$$

$$\hat{\mathbf{R}}(h) \leq \mathbf{R}(h) + \sqrt{\frac{\ln\left(\frac{1}{w(h)}\right) + \ln(2/\delta)}{2m}} \quad (35.4)$$

اگر m را از صورت قضیه در رابطه بالا جایگزین نماییم خواهیم داشت.

$$\mathbf{R}(\hat{h}) \leq \hat{\mathbf{R}}(h) + \frac{\epsilon}{2}. \quad (36.4)$$

با توجه به اینکه \hat{h} کران قید شده در صورت قضیه را کمینه می‌کند خواهیم داشت.

$$\mathbf{R}(\hat{h}) \leq \hat{\mathbf{R}}(\hat{h}) + \sqrt{\frac{\ln\left(\frac{1}{w(\hat{h})}\right) + \ln(2/\delta)}{2m}} \quad (37.4)$$

$$\hat{\mathbf{R}}(h) \leq \mathbf{R}(h) + \sqrt{\frac{\ln\left(\frac{1}{w(h)}\right) + \ln(2/\delta)}{2m}} \quad (38.4)$$

$$\leq \hat{\mathbf{R}}(h) + \frac{\epsilon}{2} \quad (39.4)$$

$$\leq \mathbf{R}(h) + \frac{\epsilon}{2} + \frac{\epsilon}{2} \quad (40.4)$$

$$= \mathbf{R}(h) + \epsilon. \quad (41.4)$$

که اثبات قضیه کامل می‌گردد.

۳.۴ کمینه طول توصیف

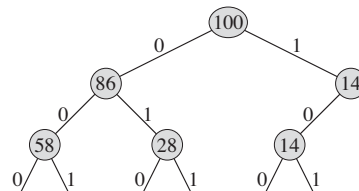
فرض کنید که H فضای شمارایی از فرضیه ها باشد در این صورت می توانیم آن را به صورت $H = \bigcup_{n \in \mathbb{N}} \{h_n\}$ بنویسیم. پیش از این با کمک نابرابری هافدینگ نشان دادیم که یک فضای تک فرضیه ای دارای ویژگی همگرایی یکخواخت با نرخ $m^{UC}(\epsilon, \delta) = \frac{\log(2/\delta)}{\epsilon^2}$ است. در نتیجه مقدار ϵ_n برابر است با $\epsilon_n(\epsilon, \delta) = \sqrt{\frac{\log(2/\delta)}{\epsilon^2}}$ و قاعده کمینه سازی هزینه ساختاری به صورت زیر است.

$$\operatorname{argmin}_{h_n \in H} \left[\hat{\mathbf{R}}(h_n) + \sqrt{\frac{\log(1/w(n)) + \log(2/\delta)}{2m}} \right] \quad (42.4)$$

به طور معادل می توانیم تابع وزن را به صورت $w : H \mapsto [0, 1]$ در نظر گیریم و قاعده کمینه سازی هزینه ساختاری به صورت زیر می شود.

$$\operatorname{argmin}_{h \in H} \left[\hat{\mathbf{R}}(h_n) + \sqrt{\frac{-\log w(h) + \log(2/\delta)}{2m}} \right] \quad (43.4)$$

در این حالت دانش پیشین به وسیله $w(h)$ که به هر فرضیه انتساب داده می شود مشخص می گردد. در این حالت به فرضیه هایی که فکر می کنیم بهتر هستند وزن بیشتری داده می شود و الگوریتم یادگیری به فرضیه هایی که وزن بیشتری دارند ترجیح بیشتری می دهد. در این بخش، یک تابع وزن روی H تعریف می کنیم که به طول توصیف فرضیه ها وابسته است. حال چگونه یک فرضیه را توصیف می کنیم؟ اگر $\Sigma = \{0, 1\}$ مجموعه الفبا و Σ^* همه رشته های ممکن باشد. یک زبان برای توصیف فضای فرضیه H تابعی به صورت $\mathcal{L} : H \mapsto \Sigma^*$ است که هر عضو $h \in H$ را به یک توصیف $\mathcal{L}(h)$ نگاشت می کند و $\mathcal{L}(h)$ توصیف h نامیده می شود. در ادامه طول این توصیف را با $|h|$ نشان می دهیم. زبان توصیف h باید ویژگی پیشوند-آزاد^۳ را داشته باشد. برای نمونه درخت هافمن که در شکل زیر نشان داده شده است یک زبان پیشوند-آزاد است.



شکل ۳.۴ درخت هافمن به عنوان نمونه ای از زبان پیشوند-آزاد.

هر مجموعه از رشته های پیشوند-آزاد دارای ویژگی زیر است که به نابرابری کرافت^۴ معروف می باشد.

لم ۱.۴ (نابرابری کرافت).

اگر $s \subset \{0, 1\}^*$ یک مجموعه از رشته های پیشوند-آزاد باشند آنگاه نابرابری زیر برقرار است.

$$\sum_{\sigma \in s} \frac{1}{2^{|\sigma|}} \leq 1. \quad (44.4)$$

برهان برای اثبات این لم، یک سکه سالم را به کار می بریم. یک رشته با طول صفر در ابتدا در نظر می گیریم. سکه را پرتاب می کنیم اگر شیر آمد نماد ۱ و اگر خط آمد نماد ۰ را به رشته اضافه می کنیم. با افزودن هر بیت به رشته بررسی می کنیم که آیا رشته عضوی از S است یا نه. اگر عضوی از S باشد کار تولید این رشته تمام شده است و رشته ای دیگر را آغاز می کنیم. از آنجا که احتمال تولید یک رشته σ برابر است با $2^{-|\sigma|}$ و رشته ها مستقل از هم تولید شده و برخی از رشته ها نیز عضو S نیستند. بنابراین لم اثبات می گردد.

³Prefix-free

⁴Kraft inequality

با توجه به نابرابری کرافت، می توانیم تابع وزن $w(h)$ را به صورت زیر تعریف کنیم.

$$w(h) = \frac{1}{\Psi^{|h|}}. \quad (۴۵.۴)$$

با توجه به تعریف بالا، می توانیم قضیه زیر را داشته باشیم.

قضیه ۷.۴.۰)

فضای فرضیه H و زبان توصیف پیشوند-آزاد $\{0, 1\}^*$ از $\mathcal{L} : H \mapsto \{0, 1\}^*$ برای توصیف این فضا را در نظر بگیرید. آنگاه برای هر مجموعه آموزشی $S \sim \mathcal{D}^m$ ، هر توزیع \mathcal{D} و هر $\delta > 0$ با احتمال دست کم $1 - \delta$ نابرابری زیر برای همه $h \in H$ برقرار است.

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \sqrt{\frac{|h| + \ln(\Psi/\delta)}{\Psi m}} \quad (۴۶.۴)$$

که $|h|$ طول توصیف $\mathcal{L}(h)$ است.

برهان اگر $w(h) = \frac{1}{\Psi^{|h|}}$ تابع وزن باشد و با استفاده از قضیه ۳.۴ و جایگزینی $\epsilon_n(m, \delta) = \sqrt{\frac{\ln(\Psi/\delta)}{\Psi m}}$ و استفاده از نابرابری $|h| \ln \Psi < |h|$ خواهیم داشت

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \sqrt{\frac{|h| \ln \Psi + \ln(\Psi/\delta)}{\Psi m}} \quad (۴۷.۴)$$

$$\leq \hat{\mathbf{R}}(h) + \sqrt{\frac{|h| + \ln(\Psi/\delta)}{\Psi m}}. \quad (۴۸.۴)$$

و بدین ترتیب قضیه اثبات می‌گردد.

روش یادگیری که کران $\hat{\mathbf{R}}(h) + \sqrt{\frac{|h| + \ln(\Psi/\delta)}{\Psi m}}$ را کمینه می‌کند روش کمینه سازی طول توصیف^۵ می‌گویند که الگوریتم آن در شکل زیر نشان داده شده است. در واقع این الگوریتم مصالحه‌ای بین خطای تجربی و طول توصیف برقرار می‌کند.

Minimum Description Length (MDL)

prior knowledge:

\mathcal{H} is a countable hypothesis class

\mathcal{H} is described by a prefix-free language over $\{0, 1\}$

For every $h \in \mathcal{H}$, $|h|$ is the length of the representation of h

input: A training set $S \sim \mathcal{D}^m$, confidence δ

output: $h \in \operatorname{argmin}_{h \in \mathcal{H}} \left[L_S(h) + \sqrt{\frac{|h| + \ln(2/\delta)}{2m}} \right]$

شکل ۴.۴ الگوریتم کمینه سازی طول توصیف

قضیه ۷.۴ بیان می‌دارد که دو فرضیه‌ای که خطای تجربی یکسانی دارند، فرضیه‌ای دارای کران خطای واقعی کمتری است که دارای طول توصیف کمتری باشد. نتیجه این قضیه با آنچه که اصل اوکم^۶ بیان می‌دارد که عبارت است از «توصیف کوتاه‌تر نسبت به توصیف طولانی‌تر ترجیح دارد» پیام یکسانی دارد. این قضیه همچنین بیان می‌دارد که فرضیه‌ای که طول توصیف طولانی‌تری دارد برای اینکه کران

^۵Minimum description length

^۶Occam razor

خطای واقعی که فرضیه با طول توصیف کمتری دارد برسد نیاز به مجموعه آموزشی با اندازه بزرگتری نیاز دارد. مساله‌ای که وجود دارد این است که طول توصیف یک فرضیه به زبان توصیف فرضیه وابسته است. اگر زبان توصیف فرضیه مستقل از مجموعه آموزشی و پیش از دیدن مجموعه آموزشی انتخاب شود آنگاه کران یاد شده در این قضیه برقرار است.

۴.۴ سازگاری

نمادهای یادگیری که تا کنون بررسی کردیم پیچیدگی نمونه‌ای را تابعی از ϵ, δ, h در نظر می‌گرفتند. سازگاری^۷ یک قاعده یادگیری است که علاوه بر اینکه پیچیدگی نمونه‌ای را تابعی از ϵ, δ, h در نظر می‌گیرد تابعی از توزیع \mathcal{P} که نمونه‌های آموزشی و آزمایشی بر اساس آن نمونه‌برداری شده‌اند در نظر می‌گیرد. سازگاری به صورت زیر تعریف می‌شود.

تعریف ۳.۴ سازگاری فرض کنید که $Z = \mathcal{X} \times \mathcal{Y}$ فضای دامنه، \mathcal{P} مجموعه‌ای از توزیع‌ها روی Z و H فضای فرضیه باشد. یک قاعده یادگیری A نسبت به H و \mathcal{P} سازگار است اگر یک تابع $m_H^{CON} : (\epsilon, \delta) \times H \times \mathcal{P} \mapsto \mathbb{N}$ وجود داشته باشد به طوری که برای هر $\delta \in (0, 1)$ ، هر $\epsilon \in (0, 1)$ ، هر $h \in H$ ، هر $D \in \mathcal{P}$ ، اگر مجموعه آموزشی با اندازه $m \geq m_H^{CON}(\epsilon, \delta, h, D)$ انتخاب شود آنگاه با احتمال دست کم $1 - \delta$ روی مجموعه آموزشی $S \sim D^m$ نابرابری زیر برقرار باشد

$$\mathbf{R}(A(S)) \leq \mathbf{R}(h) + \epsilon. \quad (۴۹.۴)$$

اگر \mathcal{P} تمام توزیع‌های ممکن باشد آنگاه A را نسبت به H سازگار عمومی^۸ می‌گویند. همان‌گونه که تعریف سازگاری نشان می‌دهد، سازگاری نسخه راحت شده یادگیری نایکناخت است. یعنی همه یادگیرهای نایکناخت، یادگیرهای سازگار هستند اما همه یادگیرهای سازگار، یادگیرهای نایکناخت نیستند. برای نمونه، مثال زیر را در نظر بگیرید.

مثال ۲.۴ (یادگیرهای تنبل).

یک الگوریتم یادگیری تنبل را در نظر بگیرید که نمونه‌های آموزشی را ذخیره می‌کند و برچسب هر نمونه آزمون x با رای اکثریت روی برچسب نمونه‌های x که در مجموعه آموزشی وجود دارد تعیین می‌گردد و در صورتی که نمونه x در مجموعه آموزشی وجود نداشت یک برچسب پیش‌فرض انتخاب می‌گردد. می‌توان نشان داد برای هر دامنه شمارای \mathcal{X} و مجموعه برچسب متناهی \mathcal{Y} و هزینه صفر-یک الگوریتم تنبل یاد شده سازگار عمومی است اما یک یادگیر نایکناخت نیست.

۵.۴ جمع بندی نمادهای مختلف یادگیری

تاکنون سه نماد مختلف یادگیری را بررسی نمودیم و اکنون به بررسی سودمندی این نمادها می‌پردازیم. سودمندی یک تعریف ریاضی به آنچه آن تعریف نیاز دارد وابسته است. در ادامه به بررسی سودمندی این تعاریف از نگاه اهداف زیر می‌پردازیم.

۱. هزینه فرضیه یادگرفته شده چیست؟

نخستین هدف پیدا کردن کران‌ها، پیدا کردن خطای واقعی فرضیه یادگرفته شده است. روش‌های یادگیری احتمالا تقریبا درست و یادگیری نایکناخت، کران خطای واقعی را براساس خطای تجربی پیدا می‌کنند. اما مدل سازگاری چنین کرانی را پیدا نمی‌کند. اما

⁷Consistency

⁸Universally consistent

همیشه می‌توان خطای واقعی فرضیه یادگرفته شده را با کمک مجموعه اعتبارسنجی تخمین زد.

۲. برای اینکه فرضیه یادگرفته شده به خوبی بهترین فرضیه باشد چند نمونه آموزشی مورد نیاز است؟

هنگامی که می‌خواهیم یک مساله یادگیری را حل نماییم نخستین پرسش این است که چند نمونه آموزشی جمع آوری نماییم؟ برای روش‌های یادگیری احتمالا تقریبا درست پاسخ روشن است اما برای روش‌های یادگیری نایک‌نواخت تعداد نمونه‌های آموزشی به بهترین فضای فرضیه و برای روش‌های یادگیری سازگار تعداد نمونه‌های آموزشی علاوه بر بهترین فضای فرضیه به توزیع نمونه‌برداری نیز وابسته است. از این سو مدل یادگیری احتمالا تقریبا درست تنها مدل قابل استفاده برای قابلیت یادگیری است. از سویی دیگر، هرچند ممکن است خطای برآورد فرضیه یادگرفته شده کم باشد اما اگر خطای تقریب H بزرگ باشد ممکن است خطای واقعی آن زیاد باشد. برای پاسخ به این پرسش که «برای اینکه فرضیه یادگرفته شده به خوبی دسته بند بهینه بیز باشد چند نمونه آموزشی مورد نیاز است؟» برای این پرسش، مدل یادگیری احتمالا تقریبا درست پاسخ روشنی ندارد. این به این دلیل است که سودمندی روش‌های مبتنی بر مدل یادگیری احتمالا تقریبا درست به کیفیت دانش پیشین وابسته است.

هنگامی که فرضیه یادگرفته شده خطای واقعی زیادی داشته باشند، کران‌های بدست آمده برای مدل یادگیری احتمالا تقریبا درست کمک می‌کنند که چه کاری باید انجام دهیم. این بدان دلیل است که می‌توانیم کرانی از خطا که ریشه آن از بخش برآورد خطا است را تعیین و در نتیجه مشخص نماییم ریشه چه مقدار از خطا از بخش تقریب است. اگر خطای تقریب بزرگ باشد باید فضای فرضیه را تغییر بدهیم. به طور مشابه اگر یک یادگیر نایک‌نواخت شکست بخورد می‌توانیم توابع وزن مختلف را در نظر بگیریم. در حالیکه اگر یک روش یادگیری سازگار شکست بخورد نمی‌دانیم که آیا بدلیل خطای برآورد است یا خطای تقریب. اگر هم بدانیم که خطای برآورد زیاد باشد نمی‌دانیم چند نمونه باید داشته باشیم تا خطای برآورد کوچک داشته باشیم.

۳. چگونه یاد بگیریم؟ چگونه دانش پیشین را نمایش دهیم؟

یکی از مهم‌ترین ویژگی‌های نظریه یادگیری این است که به پرسش «چگونه باید یاد گرفت؟» پاسخ می‌دهد. تعریف یادگیری احتمالا تقریبا درست محدودیت یادگیری را به کمک قضیه No free lunch و همچنین لزوم دانش پیشین را مشخص می‌نماید. این مدل روشی روشن برای نمایش دانش پیشین با کمک انتخاب فضای فرضیه را ارائه می‌نماید پس از انتخاب فضای فرضیه، یک روش عمومی یادگیری که همان کمینه سازی خطای تجربی است را ارائه می‌نماید

مدل یادگیری نایک‌نواخت نیز روشی روشن برای نمایش دانش پیشین با کمک تابع وزن روی فضای فرضیه را ارائه می‌نماید پس از انتخاب تابع وزن، یک روش عمومی یادگیری که همان کمینه سازی خطای ساختاری است را ارائه می‌نماید. روش یادگیری کمینه سازی خطای ساختاری این برتری را دارد که در انتخاب مدل کمک می‌کند به ویژه زمانی که دانش پیشین جزئی در اختیار داشته باشیم.

اما مدل سازگاری بر عکس دو مدل پیشین، هیچ روشی برای نمایش دانش پیشین ارائه نمی‌کند. در بسیاری از حالت‌ها، نیازی به دانش پیشین نیست. برای نمونه، الگوریتم مثال ۲.۴ یک الگوریتم سازگار برای هر فضای نمونه شمارا با تعداد متناهی برچسب است. این اشاره می‌کند که سازگاری یک نیازمندی خیلی ضعیف است.

۴. باید کدام الگوریتم یادگیری را ترجیح دهیم؟

ممکن است دلیل آورده شود که با وجود اینکه سازگاری یک نیازمندی ضعیف است، خیلی خوب است که یک الگوریتم نسبت به مجموعه همه توابع از \mathcal{L} به \mathcal{L} سازگار باشد. این روش تضمین می‌نماید که اگر به اندازه کافی نمونه آموزشی داشته باشیم دسته بند یادگرفته شده همیشه به خوبی دسته بند بهینه بیز خواهد بود. در نتیجه اگر دو الگوریتم که یکی سازگار باشد و دیگری سازگار نباشد الگوریتم سازگار را ترجیح می‌دهیم. هر چند این نتیجه‌گیری به دو دلیل مشکل ساز است. نخست آنکه در بسیاری از کاربردها و

برای بسیاری از توزیع‌ها، پیچیدگی نمونه‌ای الگوریتم‌های سازگار خیلی زیاد است و در عمل نمی‌توانیم این مقدار نمونه آموزشی جمع‌آوری نماییم. دوم اینکه، خیلی سخت نیست که الگوریتم‌های احتمالا تقریبا درست و نایکناخت را نسبت به مجموعه همه توابع از \mathcal{X} به \mathcal{Y} الگوریتم‌های سازگار تبدیل نماییم. از آنجایی که به سادگی می‌توان هر الگوریتم یادگیری را به الگوریتم سازگار تبدیل نمود بنابراین خردمندانه نیست که الگوریتم‌های سازگار را به الگوریتم‌های ناسازگار ترجیح دهیم.

پرسش‌ها

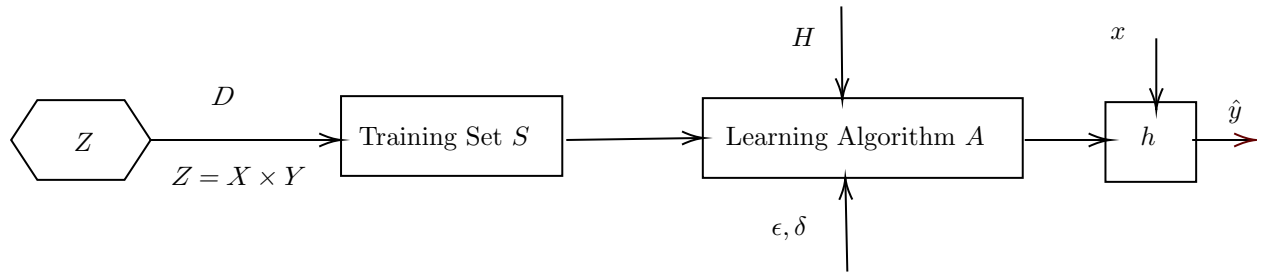
- ۱.۴ هدف یادگیری فضای فرضیه‌ای است به گونه‌ای که $\mathcal{X} = \mathbb{R}$ و به ازای تعداد متناهی نقطه برچسب $+1$ و برای بقیه نقاط برچسب -1 تولید نماید. آیا این فضا قابلیت یادگیری نایکناخت را دارد؟ ادعای خود را اثبات نمایید.
- ۲.۴ آیا قضیه NFL برای یادگیری نایکناخت برقرار است؟ ادعای خود را اثبات نمایید.

پیچیدگی محاسباتی الگوریتم‌های یادگیری

در بخش‌های پیشین در باره ویژگی‌های آماری روش‌های یادگیری که عبارت است از تعداد نمونه‌های لازم برای یادگیری، سخن گفته شد و به بیانی دیگر، در باره اطلاعات مورد نیاز یادگیری سخن گفته شد. در کاربردهای عملی، هدف داشتن یادگیری خودکار است که در آن منابع محاسباتی بسیار مهم هستند. هنگامی که نمونه‌های آموزشی جمع‌آوری شد می‌بایست محاسباتی انجام گیرد تا فرضیه نهایی تولید شود. پس از تولید این فرضیه نیاز به محاسبات است تا برچسب نمونه‌های آزمون تعیین گردند. این منابع محاسباتی در کاربردهای عملی یادگیری ماشین بسیار مهم هستند. منابع لازم برای یادگیری دو دسته هستند. (۱) نمونه‌های برچسب خورده مورد نیاز که توسط پیچیدگی نمونه‌ای مشخص می‌شوند و (۲) منابع محاسباتی لازم که توسط پیچیدگی محاسباتی مشخص می‌شوند. در این بخش پیچیدگی محاسباتی و به‌طور خاص زمان اجرای الگوریتم‌های یادگیری را بررسی می‌کنیم.

زمان اجرای یک الگوریتم به ماشین استفاده، زبان برنامه نویسی، الگوریتم طراحی شده و همچنین پیاده سازی آن وابسته است. برای جلوگیری از وابستگی زمان اجرا به ماشین (به طور مشابه برای عوامل دیگر)، زمان اجرای الگوریتم‌ها را به صورت مجانبی محاسبه می‌کنیم. برای نمونه زمان اجرای الگوریتم مرتب‌سازی ادغامی از مرتبه $O(n \log n)$ است که n اندازه ورودی است. این بدان معنا است که زمان اجرای این الگوریتم در هر ماشینی به صورت زیر تعریف می‌شود. «ثابت‌های c و n_0 وجود دارند که به ماشین وابسته هستند به گونه‌ای که برای هر مقدار $n > n_0$ زمان اجرای الگوریتم مرتب‌سازی ادغامی (به ثانیه) حداکثر $cn \log n$ است.» یک الگوریتم را کارا می‌گویند اگر در زمان از مرتبه $O(p(n))$ اجرا شود که $p(n)$ یک چند جمله‌ای براساس n است. این نوع تحلیل، زمان اجرای الگوریتم را بر اساس اندازه ورودی (n) محاسبه می‌نماید. برای نمونه در الگوریتم مرتب‌سازی ادغامی، تعداد اعضا آرایه ورودی به عنوان اندازه ورودی در نظر گرفته می‌شود.

در ادامه این فصل نخست پیچیدگی محاسباتی الگوریتم به‌طور رسمی تعریف می‌شود و سپس چندین مثال و شیوه محاسبه زمان اجرا



شکل ۱.۵ یک الگوریتم یادگیری نمونه

در آنها بررسی می‌گردد و همچنین نشان می‌دهیم که چگونه می‌توان با تغییر فضای فرضیه مسایل یادگیری را حل و در پایان سختی مسایل یادگیری بررسی می‌شود.

۱.۵ پیچیدگی محاسباتی یادگیری

یک الگوریتم یادگیری به فضای نمونه $Z = \mathcal{X} \times \mathcal{Y}$ ، فضای فرضیه‌های H ، تابع هزینه ℓ و نمونه‌های آموزشی که بر اساس توزیع ناشناخته D و به صورت مستقل از فضای نمونه Z ، نمونه برداری شده‌اند دسترسی دارد. برای هر مقدار از پارامترهای ϵ و δ الگوریتم یادگیری باید فرضیه h را تولید نماید که با احتمال دست کم $1 - \delta$ نابرابری زیر برقرار باشد.

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}(h) + \epsilon. \quad (1.5)$$

به بیانی دیگر یک الگوریتم یادگیری ماشین را می‌توان به صورت شکل ۱.۵ نشان داد.

برای تحلیل الگوریتم‌های یادگیری از روش‌های استاندارد تحلیل الگوریتم‌ها استفاده می‌کنیم که در آنها (۱) یک ماشین انتزاعی مانند ماشین تورینگ در نظر می‌گیریم و (۲) تحلیل مجانبی الگوریتم را انجام می‌دهیم. نخستین مساله‌ای که می‌بایست مشخص شود تعیین اندازه ورودی است. بر اساس اینکه اندازه ورودی را چه در نظر بگیریم تحلیل‌های متفاوتی خواهیم داشت. برای نمونه حالت‌های زیر را در نظر بگیرید.

۱. اندازه ورودی برابر با اندازه مجموعه آموزشی S باشد.

اندازه مجموعه آموزشی را نمی‌توان به تنهایی به عنوان اندازه ورودی در نظر گرفت چون اگر اندازه مجموعه آموزشی از $m(\epsilon, \delta)$ بزرگتر باشد الگوریتم یادگیری به سادگی می‌تواند از نمونه‌های اضافی صرف‌نظر نماید و آنها را استفاده نکند. بنابراین مجموعه آموزشی بزرگتر باعث افزایش زمان یادگیری نخواهد شد و مساله یادگیری را پیچیده تر نخواهد کرد.

۲. اندازه ورودی برابر با اندازه فضای فرضیه باشد

همان گونه که پیش از این دیدیم، اندازه فضای فرضیه باعث پیچیده‌تر شدن مساله نخواهد شد و باعث افزایش زمان اجرای الگوریتم یادگیری نخواهد شد. فضای فرضیه آستانه را می‌توان به صورت کارا یادگرفت در حالی که اندازه این فضا بی‌نهایت است.

۳. اندازه ورودی برابر با تعداد ابعاد فضای دامنه (تعداد ویژگی) باشد.

۴. اندازه ورودی برابر با ϵ یا δ یا تابعی از ϵ یا δ مانند $f(\epsilon, \delta)$ باشد. برای نمونه $p(\frac{1}{\epsilon}, \frac{1}{\delta})$ که یک تابع چندجمله‌ای است.

۵. یک پارامتر که پیچیدگی فضای فرضیه را نشان دهد. به همین دلیل است که در تعریف یادگیری احتمالا تقریباً درست از پارامتری که پیچیدگی یا قدرت توصیف فضای فرضیه را نشان بدهد استفاده می‌کنیم.

خروجی الگوریتم یادگیری یک تابع به فرم $\mathcal{Y} \rightarrow \mathcal{X} : h$ است. اینکه چگونه این تابع را نمایش دهیم بسیار مهم است. برای نمونه الگوریتمی را در نظر بگیرید که ورودی آن مجموعه آموزشی S است و خروجی آن توصیفی به شکل «الگوریتم کمینه‌سازی خطا را روی S اجرا کن» باشد. این توصیف قابل قبول است اما همه محاسبات را به زمان آزمون منتقل می‌کند. برای رفع مشکل زمان اجرای یک الگوریتم را به صورت زیر تعریف می‌کنیم.

تعریف ۱.۵ زمان اجرای الگوریتم یادگیری عبارت از بیشترین زمانی که برای تولید فرضیه

یا تولید برچسب یک نمونه مورد نیاز است. به عبارتی دیگر، زمان اجرای یک الگوریتم یادگیری عبارت است از

$$\max\{h \text{ فرضیه } x, \text{ زمان مورد نیاز برای تولید فرضیه } h\}.$$

این تعریف به این دلیل ارایه شده است که فرضیه تولید شده به اندازه کافی کارا باشد و بتوان از فرضیه تولید شده استفاده کرد در غیر اینصورت زمان اجرای الگوریتم‌های یادگیری معتبر نخواهد بود.

حال مساله یادگیری مربع مستطیل‌ها در فضای \mathbb{R}^n را در نظر بگیرید. که ورودی الگوریتم کمینه‌سازی خطای تجربی یک مجموعه آموزشی با اندازه m است و خروجی فرضیه‌ای است که روش کمینه‌سازی خطای تجربی آن را تولید می‌کند.

برای نمونه می‌توانیم در مساله یادگیری مستطیل با اضلاع موازی محورها، حالت‌های زیر را در نظر بگیریم.

۱. مقدار متغیرهای ϵ یا δ را ثابت در نظر گرفته و مقدار بعدها (n) را تغییر دهیم.

۲. مقدار متغیرهای ϵ یا n را ثابت در نظر گرفته و بازه اطمینان (δ) را تغییر دهیم.

۳. مقدار متغیرهای δ یا n را ثابت در نظر گرفته و مقدار خطا (ϵ) را تغییر دهیم.

در هر سه حالت یک دنباله از مسایل خواهیم داشت. هنگامی که اندازه ورودی و پارامترهای دیگر انتخاب شدند می‌توانیم تحلیل مجانبی الگوریتم را انجام دهیم.

بر اساس مدل یادگیری احتمالا تقریباً درست، روش معقول برای تحلیل پیچیدگی محاسباتی یک الگوریتم یادگیری، استفاده از پیچیدگی نمونه‌ای آن است. به بیانی دیگر پیچیدگی نمونه‌ای را می‌توان بعنوان اندازه مفید ورودی الگوریتم دانست. همان‌گونه که پیش از این بیان شد پیچیدگی نمونه‌ای به مقدار دقت (پارامتر ϵ)، بازه اطمینان (پارامتر δ) و پیچیدگی فضای فرضیه وابسته است. بنابراین باید بگوییم که فضای فرضیه H در زمان چند جمله‌ای قابل یادگیری است اگر یک الگوریتم یادگیری A وجود داشته باشد که فضای H را با خطای ϵ و اطمینان δ یاد می‌گیرد و پیچیدگی محاسباتی آن یک تابع چندجمله‌ای براساس پیچیدگی نمونه‌ای اش باشد. تعریف بالا فرض می‌کند که اندازه هر نمونه ثابت است. متداول است که به صورت صریح براساس ابعاد ویژگی تحلیل پیچیدگی محاسباتی را انجام داد. در این حالت یک ضریب $\theta(n)$ که n تعداد بعدها ورودی است در عبارت زمان اجرا ضرب می‌شود.

تعریف بالا فریبنده است بدین معنا که الگوریتم یادگیری در زمان آموزش نمونه‌ها را در حافظه ذخیره نماید و در زمانی که نمونه آزمون وارد می‌شود با روش کمینه‌سازی خطای تجربی، فرضیه را پیدا نموده و سپس برچسب نمونه آزمون را مشخص می‌نماید. برای جلوگیری از چنین فریبندگی، نیاز داریم که فرضیه تولید شده نیز در زمان چندجمله‌ای براساس پیچیدگی نمونه‌ای، برچسب نمونه‌های آزمون را تعیین نماید.

تعریف ۲.۵ پیچیدگی محاسباتی الگوریتم‌های یادگیری فرض کنید H فضای فرضیه روی فضای نمونه n بعدی، ϵ پارامتر دقت

و δ پارامتر اطمینان باشد. همچنین فرض کنید $m_H(\epsilon, \delta)$ پیچیدگی نمونه‌ای با پارامترهای دقت ϵ و اطمینان δ باشد. آنگاه پیچیدگی

محاسباتی یادگیری احتمالا تقریباً درست فضای فرضیه H از مرتبه $O(nm_H(\epsilon, \delta))$ است اگر یک الگوریتم یادگیری وجود داشته

باشد که فضای H را در زمان $O(nm_H(\epsilon, \delta))$ احتمالاً تقریباً درست یاد بگیرد و زمان اجرای فرضیه تولید شده توسط این الگوریتم برای تعیین برچسب نمونه x از مرتبه $O(nm_H(\epsilon, \delta))$ باشد.

در تعریف بالا، همه پارامترهای مرتبط همانند فضای فرضیه H ، ابعاد فضای نمونه n ، پارامتر دقت ϵ و پارامتر اطمینان δ به کار رفته اند. در نتیجه اندازه مفید ورودی $n m_H(\epsilon, \delta)$ نیز یک عدد ثابت است. با توجه به تعریف بالا، حال یادگیری در زمان چندجمله‌ای را تعریف می‌نماییم. برای این تعریف ما در باره دنباله‌ای از مسایل با رشد اندازه مفید ورودی سخن می‌گوییم. از آنجایی که اندازه مفید ورودی به فضای فرضیه H ، ابعاد فضای نمونه n ، پارامتر دقت ϵ و پارامتر اطمینان δ وابسته است لذا ما باید به صورت دقیق مشخص نماییم که کدام یک از اینها ثابت و کدام یک متغیر هستند. برای روشن شدن مطلب، یادگیری فضای فرضیه متناهی را در نظر بگیرید. این مساله را می‌توانیم با استفاده از جستجوی کورکورانه روی H با استفاده از مجموعه آموزشی با اندازه $m_H(\epsilon, \delta)$ در زمان $O(nm_H(\epsilon, \delta))$ یاد بگیریم. اگر $|H|$ را ثابت فرض کنیم و بقیه پارامترها را متغیر فرض کنیم این الگوریتم در زمان چند جمله‌ای اجرا می‌شود. اگر بقیه پارامترها را ثابت فرض کنیم و زمان اجرا را بر اساس $|H|$ تحلیل کنیم آنگاه زمان اجرای این الگوریتم نمایی است در حالیکه $m_H(\epsilon, \delta)$ به صورت لگاریتمی با $|H|$ رشد می‌کند. زیرا $|H|$ بر اساس n یک تابع نمایی است. با توجه به مطالب بیان شده، یادگیری کارا به صورت زیر تعریف می‌شود.

تعریف ۳.۵ یادگیری کارا دنباله‌ای از مسایل یادگیری $(n_k, H_k, \epsilon_k, \delta_k)_{k=1}^{\infty}$ به صورت کارا قابل یادگیری است اگر یک چند جمله‌ای p وجود داشته باشد به گونه‌ای که برای هر k یک الگوریتم یادگیری وجود داشته باشد به گونه‌ای که فضای H_k را به صورت احتمالاً تقریباً درست با پارامترهای ϵ_k و δ_k در زمان $p(n_k m_{H_k}(\epsilon_k, \delta_k))$ یاد بگیرد.

برای نمونه در یادگیری فضای فرضیه متناهی که پیش از این بیان شد دو دنباله از مسایل یادگیری زیر را در نظر بگیرید.

۱. در دنباله نخست پارامترهای n_k, H_k, δ_k را ثابت فرض می‌کنیم که با n, H, δ نمایش می‌دهیم و مقدار $\epsilon_k = \frac{1}{k}$ در نظر می‌گیریم که متغیر است. در این حالت چندجمله‌ای به صورت $p(x) = |H|x$ و روشن است که الگوریتم کمینه سازی خطای تجربی در زمان $p(n_k m_{H_k}(\epsilon_k, \delta_k))$ اجرا می‌شود و بنابراین این دنباله از مسایل یادگیری به صورت کارا قابل یادگیری هستند.
۲. در دنباله دوم پارامترهای $n_k, \epsilon_k, \delta_k$ را ثابت فرض می‌کنیم که با n, ϵ, δ نمایش می‌دهیم و مقدار $\log |H_k| = k$ در نظر بگیریم که متغیر است. در این حالت هیچ الگوریتم کمینه سازی خطای تجربی نمی‌توان پیدا کرد که در زمان $p(n_k m_{H_k}(\epsilon_k, \delta_k))$ مساله را حل نماید و در نتیجه روش کمینه سازی خطای تجربی نمی‌تواند این مساله را در زمان چند جمله‌ای حل نماید. حال با توجه به دو تعریفی که در این بخش گفته شد پیچیدگی محاسباتی الگوریتم‌های یادگیری را به صورت زیر تعریف می‌کنیم.

تعریف ۴.۵ پیچیدگی محاسباتی الگوریتم‌های یادگیری پیچیدگی محاسباتی الگوریتم‌های یادگیری در دو گام نخست مشخص می‌شوند. نخست پیچیدگی محاسباتی یک مساله یادگیری ثابت که با سه تایی (Z, H, ℓ) مشخص می‌شود را محاسبه می‌کنیم که Z فضای نمونه، H فضای فرضیه و ℓ تابع هزینه است. سپس در گام دوم نرخ تغییرات پیچیدگی محاسباتی را با دنباله‌ای از این مسایل در نظر می‌گیریم.

۱. تابع $f: (0, 1)^2 \rightarrow \mathcal{N}$ ، مساله یادگیری (Z, H, ℓ) و الگوریتم یادگیری A را در نظر بگیرید. می‌گوییم الگوریتم A مساله یادگیری را در زمان $O(f)$ حل می‌کند اگر یک ثابت c وجود داشته باشد به گونه‌ای که برای هر توزیع D روی Z و ورودی‌های $\epsilon, \delta \in (0, 1)$ هنگامی که الگوریتم A مجموعه آموزشی که به صورت مستقل از توزیع D نمونه برداری شده‌اند را دریافت نماید

(آ) الگوریتم A پس از زمان حداکثر $cf(\epsilon, \delta)$ پایان می‌یابد.

(ب) خروجی الگوریتم A که با h_A نمایش داده می‌شود را می‌توان برای پیش‌بینی برچسب نمونه‌های جدید به کار برد به گونه‌ای که زمان اجرای آن حداکثر $cf(\epsilon, \delta)$ باشد.

(ج) خروجی الگوریتم A احتمالاً تقریباً درست باشد یعنی با احتمال دست‌کم $1 - \delta$ نابرابری زیر برقرار است.

$$\mathbf{R}(h_A) \leq \min_{h' \in H} \mathbf{R}(h') + \epsilon. \quad (2.5)$$

۲. دنباله‌ای از مسایل یادگیری $(Z_k, H_k, \ell_k)_{k=1}^{\infty}$ را در نظر بگیرید به گونه‌ای که مساله k ام با فضای دامنه Z_k ، فضای فرضیه H_k و تابع هزینه ℓ_k مشخص می‌گردد. فرض کنید که الگوریتم A برای حل مسایلی از این دست طراحی شده باشد. یک تابع $\mathcal{N} \mapsto \mathcal{N} \times (0, 1)^2$ را g در نظر بگیرید. می‌گوییم زمان اجرای الگوریتم A برای این دنباله از مسایل $O(g)$ است اگر برای همه h, k ، الگوریتم A مساله (Z_k, H_k, ℓ_k) را در زمان $O(f_k)$ حل می‌کند به گونه‌ای که $\mathcal{N} \mapsto (0, 1)^2$ به صورت $f_k(\epsilon, \delta) = g(n, \epsilon, \delta)$ تعریف می‌شود.

می‌گوییم که الگوریتم A نسبت به دنباله (Z_k, H_k, ℓ_k) یک الگوریتم کارا است اگر زمان اجرای آن از مرتبه $O(p(n, 1/\epsilon, 1/\delta))$ باشد که p یک چندجمله‌ای براساس پارامترهایش است.

از تعریف بالا روشن است که کارایی الگوریتم به شیوه شکستن مساله به دنباله‌ای از مسایل وابسته است. برای نمونه، یادگیری یک فضای فرضیه متناهی را در نظر بگیرید. همان گونه که پیش از این نشان دادیم روش کمینه سازی خطای تجربی روی فضای فرضیه H یادگیری را تضمین می‌کند اگر اندازه مجموعه آموزشی دست کم $\epsilon^{-2} \log(|H|/\delta)$ باشد. فرض کنید که زمان که زمان ارزیابی هر فرضیه به زمان ثابتی نیاز داشته باشد، الگوریتم جستجوی کورکورانه در زمان $O(|H|m_H(\epsilon, \delta))$ مساله را یاد می‌گیرد. برای هر فضای فرضیه متناهی اگر دنباله مسایل را به صورت $|H_k| = k$ تعریف کنیم الگوریتم A کارا است. اما اگر دنباله مسایل را به صورت $|H_k| = 2^k$ تعریف کنیم آنگاه پیچیدگی نمونه‌ای هنوز به یک تابع چندجمله‌ای براساس پارامترهایش است اما پیچیدگی محاسباتی الگوریتم A به صورت نمایی رشد با k رشد می‌کند بنابراین دیگر الگوریتمی کارا نیست.

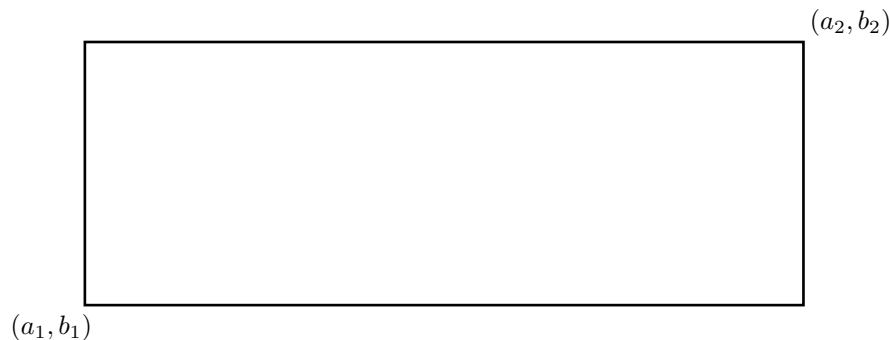
۲.۵ پیاده سازی قاعده کمینه سازی خطای تجربی

برای هر فضای فرضیه H ، قاعده کمینه سازی خطای تجربی یک رویکرد متداول یادگیری است. برای نمونه، پیش از این دیدیم که اگر برای یک مساله دسته بندی دودویی، یادگیری امکان‌پذیر باشد با قاعده کمینه سازی خطای تجربی نیز امکان‌پذیر است. در ادامه این بخش، پیچیدگی محاسباتی پیاده سازی قاعده کمینه سازی خطای تجربی برای فضاهای فرضیه مختلف را بررسی می‌کنیم. برای فضای دامنه Z ، فضای فرضیه H و تابع هزینه ℓ قاعده کمینه سازی خطای تجربی به صورت زیر تعریف می‌شود. برای هر مجموعه آموزشی متناهی $S \sim \mathcal{D}^m$ ، فرضیه‌ای را تولید می‌کند که خطای تجربی $\frac{1}{|S|} \sum_{z \in S} \ell(h, z)$ را کمینه نماید.

فضای فرضیه متناهی

۱.۲.۵

محدود کردن فضای فرضیه به یک فضای متناهی، بسیار منطقی است زیرا هنگامی که می‌خواهیم یک فرضیه را در کامپیوتر نشان دهیم باید گسسته سازی صورت پذیرد. برای نمونه فضای خطوط $ax + b$ را در نظر بگیرید. برای نشان دادن پارامترهای این خط باید آن‌ها را در حافظه کامپیوتر ذخیره‌سازی نماییم برای نمونه هرکدام از پارامترهای a و b را با ۱۶ بیت نمایش دهیم در این صورت تعداد خطوطی که با این حافظه



شکل ۲.۵ نمایش یک مستطیل در فضای دوبعدی

می‌توان نمایش داد متناهی خواهد بود. نمونه دیگر فضای فرضیه‌ای را در نظر بگیرید که بتوان با ۱۰۰۰ بیت برنامه $C + +$ آن‌ها را نوشت. همان گونه که در بخش‌های پیشین نشان دادیم، کران بالای پیچیدگی نمونه‌ای یادگیری فضای فرضیه متناهی برابر است با $m_H(\epsilon, \delta) = \log(c|H|/\delta)/\epsilon^c$ که برای حالت تحقق‌پذیر $c = 1$ و برای حالت تحقق‌ناپذیر $c = 2$ است. برای نمونه برای مثال پیاده‌سازی با زبان $C + +$ ، تعداد فرضیه‌ها برابر است با 2^{1000} در حالی که پیچیدگی نمونه‌ای برابر است با $m_H(\epsilon, \delta) = c(1000 + \log \frac{c}{\delta})/\epsilon^c$. یک رویکرد سراسر برای پیاده‌سازی قاعده کمیته‌سازی خطای تجربی برای فضای فرضیه متناهی جستجوی کورکورانه است. اگر فرض کنید که محاسبه مقدار $\ell(h, z)$ برای یک نمونه در زمان ثابت k انجام پذیرد. زمان اجرای الگوریتم کمیته‌سازی قاعده تجربی برابر است با $k|H|m$ که m اندازه مجموعه آموزشی است. اگر مقدار m برابر کران بالای پیچیدگی محاسباتی باشد در این صورت زمان اجرای این الگوریتم برابر است با $c k |H| \log(c|H|/\delta)/\epsilon^c$. این رابطه نشان می‌دهد که زمان اجرای الگوریتم به صورت خطی به اندازه فضای فرضیه $|H|$ وابسته است و این وابستگی خطی سبب می‌شود که این الگوریتم برای فضاهای فرضیه بزرگ مناسب نباشد. اگر دنباله مسایل $(Z_k, H_k, \ell_k)_{k=1}^{\infty}$ را به گونه‌ای تعریف کنیم که $\log(|H_k|) = k$ ، زمان اجرای جستجوی کورکورانه از مرتبه‌نمایی است. به همین جهت معمولاً از فضاهای فرضیه دیگر مانند فضای ابرصفحه‌ها استفاده می‌شود. نکته‌ای که بسیار مهم است این است که کارایی پایین یک روش پیاده‌سازی قاعده کمیته‌سازی خطای تجربی به این معنا نیست که نمی‌توان این قاعده را به صورت کارا پیاده‌سازی نمود. در ادامه این بخش مثال‌های دیگری را بررسی می‌کنیم که می‌توان این قاعده را به صورت کارا پیاده‌سازی نمود.

۲.۲.۵ مستطیل‌هایی با اضلاع موازی محورهای مختصات

پیش از این درباره فضای مستطیل‌های با اضلاع موازی محورهای مختصات سخن به میان آمد. با توجه به شکل زیر می‌توانیم یک مستطیل در فضای n بعدی را به $2n$ پارامتر نشان داد.

بنابراین H_n را که عبارت است از فضای همه مستطیل‌های با اضلاع موازی محورهای مختصات در فضای n بعدی را به صورت زیر در تعریف می‌کنیم.

$$H_n = \{h_{(a_1, a_2, \dots, a_n, b_1, b_2, \dots, b_n)} : \forall i, a_i \leq b_i\} \quad (۳.۵)$$

که

$$h_{(a_1, a_2, \dots, a_n, b_1, b_2, \dots, b_n)}(x) = \begin{cases} 1 & \text{if } \forall i, x_i \in [a_i, b_i] \\ 0 & \text{otherwise.} \end{cases} \quad (۴.۵)$$

نخست حالت سازگار را در نظر بگیرید یعنی در فضای فرضیه یک فرضیه $h \in H_n$ وجود دارد به گونه ای که برای همه نمونه های آموزشی $(x, y) \in S$ شرط $h(x) = y$ برقرار است. حال پیاده سازی قاعده کمینه سازی خطای تجربی برای این فضا را به این صورت در نظر بگیرید که a_i برابر باشد با کوچکترین مقدار بعد i ام نمونه های مثبت و b_i برابر باشد با بزرگترین مقدار بعد i ام نمونه های مثبت. به بیانی دیگر برای همه i ها داشته باشیم

$$a_i = \min\{x_i : (x_i, 1) \in S\}$$

$$b_i = \max\{x_i : (x_i, 1) \in S\}$$

به سادگی می توان نشان داد که خطای تجربی برابر صفر است و فرضیه تولید شده با داده های آموزشی سازگار است. برای پیدا کردن هر کدام از a_i ها و b_i ها نیاز به زمان $O(m)$ است و چون n بعد داریم لذا زمان اجرای الگوریتم $O(nm)$ است. حال حالتی را در نظر بگیرید که فرضیه سازگار وجود نداشته باشد و به دنبال فرضیه ای هستیم که خطای آموزش را کمینه نماید. نشان داده شده است که کمینه سازی خطای تجربی در این حالت یک مساله NP-Hard است و بنابراین راه حل چند جمله ای ندارد. اما اگر فضای فرضیه را محدود نماییم برای نمونه بعد n را ثابت در نظر بگیریم الگوریتم کارا برای حل مساله وجود دارد. یعنی الگوریتم یادگیری احتمالا تقریباً درست بدون پیش فرض وجود دارد که در زمان چند جمله ای با $\frac{1}{6}$ و $\frac{1}{6}$ اجرا می شود اما تابعی چند جمله ای از n نیست. برای نمونه برای پیدا کردن مستطیل در n بعد حداکثر به $2n$ نمونه است بنابراین می توانیم همه زیر مجموعه های $2n$ عضوی از مجموعه m عضوی آموزش را در نظر بگیریم و با استفاده از جستجوی کورکورانه مساله را حل نماییم. بدیهی است که الگوریتم فرضیه با کمینه خطای تجربی را پیدا می کند. زمان اجرای این الگوریتم برابر است با $m^{O(n)}$. در نتیجه این الگوریتم در زمان چند جمله ای نسبت به m و در زمان نمایی نسبت به n اجرا می شود.

۳.۲.۵ ترکیب عطفی تعدادی ویژگی دودویی

یک ترکیب عطفی تعدادی ویژگی دودویی نگاشتی از $X = \{0, 1\}^n$ به $Y = \{0, 1\}$ است و می توان به صورت گزاره $x_{i_1} \wedge x_{i_2} \wedge \dots \wedge x_{i_k} \wedge \bar{x}_{j_1} \wedge \bar{x}_{j_2} \wedge \dots \wedge \bar{x}_{j_r}$ نشان داد که $i_1, i_2, \dots, i_k, j_1, j_2, \dots, j_r \in \{1, 2, \dots, n\}$. تابعی که این گزاره نشان می دهد برابر است با

$$h(x) = \begin{cases} 1 & \text{if } x_{i_1} = \dots = x_{i_k} = 1 \text{ and } x_{j_1} = \dots = x_{j_r} = 0 \\ 0 & \text{otherwise.} \end{cases} \quad (5.5)$$

حال مجموعه چنین توابعی را با H_n نشان می دهیم. نخست پیاده سازی قاعده کمینه سازی خطا را در حالتی که فرضیه سازگار وجود دارد در نظر بگیرید. نشان می دهیم که برای این حالت الگوریتمی کارا وجود دارد. این الگوریتم نخست فرضیه $h = (x_1 \wedge \bar{x}_1) \wedge (x_2 \wedge \bar{x}_2) \wedge \dots \wedge (x_n \wedge \bar{x}_n)$ در نظر می گیرد. بدیهی است که این فرضیه به همه نمونه هایی که دریافت می کند بر حسب صفر انتساب می دهد. این الگوریتم با دریافت نمونه های مثبت، حروفی (ویژگی ها یا نقیض آنها) که با نمونه سازگار نباشند را از فرضیه حذف می کند و با دریافت نمونه های منفی کاری انجام نمی دهد. زمان اجرای این الگوریتم $O(nm)$ است. اما برای حالتی که فرضیه سازگار موجود نباشد الگوریتمی با زمان اجرای چند جمله ای بر اساس n و m وجود ندارد.

۴.۲.۵ یادگیری کلاس مفاهیم 3-term DNF

در این مثال نشان می دهیم که تعمیم ترکیب عطفی ویژگی های دودویی منجر به بغرنج شدن حل مساله کمینه سازی خطای تجربی می شود. فضای فرضیه های H_n شامل همه 3-term DNF ها (یا ترکیب سه عبارت با طول دلخواه) باشد. برای نمونه تابع $h(x_1, \dots, x_n) =$

یک 3-term DNF $(x_1 \wedge x_2 \wedge x_3) \vee (x_1 \wedge x_4) \vee (x_2 \wedge x_4)$ است. تعداد توابع 3-term DNF برابر است با 3^{3n} و بنابراین اندازه H_n برابر است با 3^{3n} (یا $\log|H_k| = \theta(n)$) و در نتیجه کران بالای پیچیدگی نمونه‌ای آن برابر است با $3n \log(3/\delta)/\epsilon$. در این حالت $h(x)$ را می‌توانیم به صورت $h(x) = A_1(x) \vee A_2(x) \vee A_3(x)$ نمایش دهیم و زمانی $h(x) = 1$ است که دست کم یکی از عبارت‌های $A_1(x)$ یا $A_2(x)$ یا $A_3(x)$ صفر نباشند. نشان داده شده است که هیچ الگوریتم چندجمله‌ای برای حل این مساله وجود ندارد.

۳.۵ یادگیری مستقل از نمایش، سخت نیست

در بخش پیش نشان دادیم که هیچ الگوریتمی با پیچیدگی زمانی چند جمله‌ای برای حل مساله 3-term DNF وجود ندارد. در این بخش نشان می‌دهیم که این فضا را با کمک الگوریتم کمینه سازی خطا با استفاده از فضای بزرگتری می‌توانیم یاد بگیریم. این نتیجه تناقضی با نتیجه بخش پیشین ندارد. زیرا الگوریتم یادگیری فرضیه‌ای را تولید می‌کند که به فضای فرضیه تعریف شده اصلی تعلق ندارد. در اینجا فضای توابع 3-term DNF را با فضای بزرگتری که با کارایی بیشتری قابل یادگیری است جایگزین می‌کنیم. یعنی فرضیه‌ای که الگوریتم یادگیری تولید می‌نماید الزاما یک 3-term DNF نیست. به همین دلیل نام یادگیری مستقل از نمایش یا یادگیری نامناسب^۱ بر روی آن گذاشته شده است. برای حل مساله از این واقعیت استفاده می‌کنیم که هر 3-term DNF یک 3-term DNF معادل دارد. برای پیدا کردن این معادل، \vee را روی \wedge توزیع می‌کنیم. با توزیع \vee روی \wedge ، فضای ویژگی از n به $(2n)^3$ تغییر می‌کند. با توجه به اینکه فرض کردیم فرضیه سازگار وجود دارد لذا با الگوریتم کمینه سازی خطای تجربی می‌توانیم این فرضیه را پیدا کنیم. کران بالای پیچیدگی نمونه‌ای آن برابر است با $n^3 \log(1/\delta)/\epsilon$ و در نتیجه زمان اجرا چندجمله‌ای براساس n است. این مثال نشان می‌دهد که می‌توانیم با تغییر نمایش مساله راحت‌تر یاد گرفته شود.

۴.۵ سختی یادگیری

پیش از این در باره پیچیدگی محاسباتی یادگیری سخن به میان آمد. همچنین نشان دادیم که سختی پیاده سازی الگوریتم کمینه سازی خطای تجربی برای فضای H الزاما به این معنا نیست که فضای فرضیه H قابل یادگیری نیست. حال پرسش این است که چگونه می‌توانیم نشان دهیم که یک مساله یادگیری از نظر محاسباتی یک مساله سخت است. یک رویکرد استفاده از فرضیات مورد استفاده در رمزنگاری است. از یک رویکرد، می‌توان یادگیری را مخالف رمزنگاری دانست. در یادگیری ما فرض می‌کنیم اگر به مجموعه آموزشی که توسط یک مفهوم برچسب خورده است دسترسی داشته باشیم می‌توانیم با خطای قابل قبولی مفهوم را بیاد بگیریم و مفهوم را استخراج نماییم در حالیکه در رمزنگاری هدف این است که یک فرد نمی‌تواند با دریافت تعدادی پیام، محتوی آن‌ها را استخراج نماید. بنابراین می‌توانیم نتایج موجود در رمزنگاری را به حوزه یادگیری منتقل نماییم. بدبختانه در حال حاضر کسی نمی‌تواند ادعا نماید که یک پروتکل رمزنگاری امن است. در ادامه این بخش ما از نتایج بدست آمده در رمزنگاری برای نشان دادن یادگیری استفاده می‌کنیم. در رمزنگاری برای اثبات اینکه یک پروتکل رمزنگاری امن است از فرضیات رمزنگاری استفاده می‌شود. بسیاری از روش‌های رمزنگاری فرض می‌کنند که یک تابع یک طرفه $f : \{0, 1\}^n \rightarrow \{0, 1\}^n$ وجود دارد که به سادگی محاسبه می‌شود. اما محاسبه وارون آن سخت است. بطور رسمی تابع f را می‌توان در زمان چند جمله‌ای براساس n محاسبه نمود اما برای هر الگوریتم تصادفی با زمان اجرای چندجمله‌ای A و برای هر چندجمله‌ای $p(\cdot)$

¹Improper learning

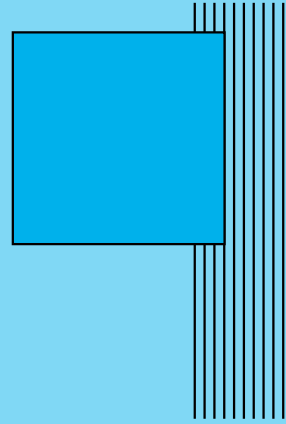
نابرابری زیر برقرار است.

$$\mathbb{P}[f(A(f(x))) = f(x)] < \frac{1}{p(n)}, \quad (۶.۵)$$

که احتمال براساس انتخاب تصادفی x براساس توزیع یکنواخت روی $\{0, 1\}^n$ و تصادفی بودن الگوریتم A محاسبه می‌شود. برای اینکه محاسبه وارون آن سخت باشد از یک کلید استفاده می‌شود و شخصی که کلید s_n را داشته باشد به سادگی می‌تواند تابع f را محاسبه و شخصی که کلید s_n را نداشته باشد به سختی می‌تواند وارون تابع f را محاسبه نماید.

فرض کنید که F_n خانواده ای از توابع رمزنگاری روی $\{0, 1\}^n$ باشد که در زمان چند جمله ای قابل محاسبه هستند. اگر ما الگوریتم رمزنگاری را ثابت فرض کنیم و یک کلید s_n که یک تابع درون F_n را نشان می‌دهد و رشته x به آن بدهیم در زمان چندجمله‌ای می‌تواند $f(x)$ را محاسبه نماید. حال مساله یادگیری وارون این فضا $H_n^F = \{f^{-1} : f \in F_n\}$ را در نظر بگیرید. از آنجایی که وارون هر تابع در این فضا را می‌توان بوسیله کلید s_n محاسبه نمود. اندازه این کلید یک چندجمله ای براساس n (که با $p(n)$ نمایش داده می‌شود) است. بنابراین فضای H_n^F را می‌توان براساس کلید پارامتری نموده و اندازه آن حداکثر $2^{p(n)}$ است و در نتیجه پیچیدگی نمونه‌ای آن چندجمله‌ای است.

می‌توان نشان داد که هیچ یادگیر کارایی برای حل این مساله وجود ندارد. اگر چنین یادگیری وجود داشت می‌توانستیم از روی رشته های فضای $\{0, 1\}^n$ با توزیع یکنواخت نمونه برداری نمود و تابع $f(x)$ را برای هر نمونه x محاسبه نمود و سپس یک مجموعه آموزشی با اندازه چندجمله‌ای از زوج‌های $(f(x), x)$ تولید نمود و در پایان با تقریب تابع f^{-1} را بدست آورد که یکطرفه بودن f را نقض می‌کند.



نمایه

ب	بیش برآزش، ۴۵
پ	پیچیدگی نمونه‌ای، ۱۳
ت	توابع عطفی عطف یکنوا، ۷ توابع فصلی فرم نرمال دو لفظی، ۱۱
خ	خصوصیات، ۶ خطا خطای برآورد، ۴۶ خطای تجربی، ۱۳ خطای تقریب، ۴۶ خطای واقعی، ۱۲.۶
د	دانش پیشین، ۴۵ دسته بند، ۶ دسته‌ای، ۲
ف	فرضیه، ۶ فضای نمونه، ۶
ک	کران اجتماعات، ۱۸ کلاس مفاهیم، ۷
م	مجموعه آموزشی، ۶ مجموعه برجسب ها، ۶ مدل احتمالاً تقریباً درست، ۱۲ مفهوم، ۷.۶ اندازه مفهوم، ۷
ن	نمونه، ۶ نمونه های آزمایشی، ۶
و	ویژگی ها، ۶ ویژگی همگرایی نایکنواخت، ۴۸
ی	یادگیر عمومی، ۴۵ یادگیری نایکنواخت، ۴۸
K	توابع عطفی فرم نرمال عطفی درجه، ۹