

کنترل کردن سویدگی در استفاده کردن از اطلاعات

How much does your data exploration overfit?

Controlling bias via information usage

Daniel Russo and James Zou

گزارش پروژه‌ی درس تئوری اطلاعات

محمد رضا رحمانی

بهراد منیری

فهرست مطالب

۳	۱ مقدمه
۳	۱.۱ نمونه‌هایی از تحلیل تطبیقی داده
۴	۲.۱ مدل ریاضی تحلیل تطبیقی داده
۴	۳.۱ کارهای پیشین
۴	۲ کنترل کردن سویدگی از راه محدود کردن استفاده از اطلاعات
۴	۱.۲ معرفی یک کران بالا برای سویدگی به وسیله‌ی اطلاعات استفاده شده
۶	۲.۲ اطلاعات استفاده شده، به عنوان کرانی برای بقیه‌ی معیارهای سویدگی
۶	۳.۲ اطلاعات استفاده شده به عنوان یک کران پایین برای سویدگی
۷	۳ چه زمانی سویدگی بزرگ است؟
۷	۱.۳ رتبه‌بندی با وجود سیگنال در داده‌ها
۸	۲.۳ تفکیک تقریباً مستقل داده‌های غیر i.i.d.
۸	۴ محدود کردن سویدگی با تصادفی‌سازی
۹	۱.۴ رگولاریزیشن با انتخاب تصادفی
۹	۲.۴ تصادفی‌سازی در تحلیل‌های چند مرحله‌ای
۱۱	۵ پیشنهاد برای کارهای آتی
۱۳	آ پیش‌نیازهای موردنیاز جهت فهم اثبات‌های قضایا
۱۳	۱.آ متغیرهای تصادفی زیر-گاوسی
۱۷	۲.آ متغیرهای تصادفی زیر-نمایی
۱۸	۳.آ بیان دیگری از فاصله‌ی KL
۱۸	۱.۳.آ اثبات قضیه‌ی (۸.آ)
۱۹	ب اثبات قضایای بیان شده
۱۹	۱.ب قضایای فصل ۲
۱۹	۱.۱.ب اثبات قضیه‌ی (۱.۲)
۲۱	۲.۱.ب اثبات قضیه‌ی (۲.۲)
۲۲	۳.۱.ب اثبات قضیه‌ی (۳.۲)
۲۳	۴.۱.ب اثبات قضیه‌ی (۴.۲)
۲۴	۵.۱.ب اثبات قضیه‌ی (۵.۲)
۲۵	۶.۱.ب اثبات قضیه‌ی (۶.۲)
۲۷	۲.ب اثبات قضایای فصل ۳
۲۷	۱.۲.ب اثبات قضیه‌ی (۲.پ)
۲۸	۳.ب اثبات قضایای فصل ۴
۲۸	۱.۳.ب اثبات لم (۱.۴)
۲۸	۲.۳.ب اثبات قضیه‌ی (۱.۴)
۲۹	پ برخی دیگر از کاربردهای مسئله‌ی کاهش سویدگی از طریق کنترل اطلاعات استفاده شده
۲۹	۱.پ فیلتر کردن به کمک آماره‌های حاشیه‌ای
۳۱	۲.پ استفاده کردن از اطلاعات و مسئله‌ی طبقه‌بندی
۳۲	۳.پ مشاهده‌ی داده‌ها و تعداد کلاس‌ها در خوشه‌بندی
۳۲	۴.پ کنترل سویدگی از طریق کنترل FDR

در این پروژه به بررسی مقاله‌ی [۱] می‌پردازیم. در تحلیل داده‌های با بُعد بالا، عموماً آنالیز که باید بر روی داده انجام شود و متغیرهایی که باید تخمین زده شوند، قبل از مشاهده و بررسی ابتدایی داده‌ها برای تحلیل‌گر واضح نیستند. به همین سبب، در اکثر مواقع، تحلیل به صورت تطبیقی انجام می‌پذیرد، به این نحو که با انجام آزمایش‌هایی بر روی داده و به صورت تدریجی، آنالیزهای جالبی که باید بر روی داده انجام شوند کشف شده و سپس به کمک همین داده، این آنالیزها انجام می‌پذیرند. این تحلیل تطبیقی و مشاهده‌ی داده، قبل از انجام آنالیز نهایی، ممکن است باعث ایجاد سویدگی یا بایاس در نتیجه‌ی تحلیل نهایی شود. در آمار و استنتاج آماری، معمولاً فرض بر این است که قبل از انجام هر تحلیل، به داده هیچ «نگاهی» نشده است و انتخاب تحلیل‌های انجام شده مستقل از داده هستند. این فرض اساسی در آمار کلاسیک، تفاوتی چشمگیر با آنچه در عمل تحلیل‌گران داده انجام می‌دهند دارد. به عنوان مثال فرض کنید تحلیل گر داده قصد دارد داده‌های خود را به کمک الگوریتمی ساده مثل k -means به تعدادی خوشه تقسیم کند. معمولاً تعداد خوشه‌ها توسط تحلیل‌گر و بعد از رسم نمودارهای مختلفی از داده تعیین می‌شود. آیا این استفاده از داده برای تعیین آنالیز انجام شده، می‌تواند باعث خطای بزرگی در نتایج تحلیل شود؟ آیا راهی برای کمی‌سازی این بایاس و کنترل آن وجود دارد؟ در این مقاله به چنین سوالاتی پاسخ داده می‌شود.

در این مقاله، هدف ارائه‌ی چهارچوبی نظریه اطلاعاتی برای کران‌زدن سویدگی ناشی از نشت اطلاعات داده در حین فرآیند انتخاب یک آنالیز است. در این چهارچوب، برخلاف استنتاج کلاسیک آماری، فرض بر این است که از خود داده در فرآیند تعیین تحلیل استفاده شده است. این مقاله به کمک چهارچوب ارائه‌شده، سویدگی تعدادی از فرآیندهای معروف تحلیل داده‌ای که از خود داده برای انتخاب تحلیل مورد نظر استفاده می‌کنند را مطالعه می‌کند. در آخر نیز از ایده‌ی اضافه کردن نویز در حین تحلیل داده برای کاهش سویدگی استفاده کرده و به بررسی دسته‌ای از فرآیندهای تطبیقی تحلیل داده می‌پردازد.

۱.۱ نمونه‌هایی از تحلیل تطبیقی داده

در این بخش، دو مثال از تحلیل تطبیقی داده ارائه می‌شود.

مثال ۱.۱. بهراد دیتاستی مربوط به تغییرات وزن و علل ژنتیکی آن در اختیار دارد. این دیتاست شامل داده‌ی ۱۰۰۰ نفر است و برای هر شخص، وزن در سه زمان مشخص اندازه‌گیری شده است. در این دیتاست، برای هر فرد نیز $expression\ value$ ۲۰۰۰ ژن نیز ثبت گشته است. سه تغییر وزن در این دیتاست می‌تواند مورد بررسی قرار بگیرد. تغییر وزن از زمان ۱ تا ۲، زمان ۲ تا ۳ یا از زمان ۱ تا ۳. بهراد، قبل از شروع هر آزمایشی، تصمیم گرفته است که تنها تغییرات از زمان ۱ تا ۳ را مورد بررسی قرار دهد. بهراد مقدار همبستگی هر ژن با تغییر وزن را محاسبه کرده و ژن با بیشترین همبستگی، به همراه R -Squared رابطه‌ی خطی تغییر وزن و آن میزان $expression$ ژن را گزارش می‌کند. مشاهده می‌شود که اگر بار دیگر این آزمایش به صورت مستقل انجام شود، همبستگی و ویژگی انتخاب شده توسط بهراد با تغییر وزن در داده‌های جدید، کمتر از همبستگی به دست آمده توسط بهراد در این آزمایش خواهد بود. این پدیده حتی در صورتی که بهراد ابتدا به کمک یک آزمون فرضیه، تعدادی از ژن‌ها را که همبستگی آن‌ها با وزن معنادار است انتخاب کرده و سپس ماکزیمم همبستگی را در بین آنها انتخاب کند نیز مشاهده می‌شود. این پدیده، *Winner's Curse Selection Bias* نام دارد. این سویدگی، به این دلیل ایجاد شده است که بهراد از همین دیتاست برای انتخاب ژن مدنظر استفاده کرده است.

مثال ۲.۱. محمدرضا همان داده‌ی بهراد را در اختیار دارد و چند آزمایش ساده با آن انجام می‌دهد. در ابتدا، برای هر زمان، میانگین $expression$ ژن‌ها در افراد مختلف محاسبه می‌کند. مشاهده می‌کند که این میانگین در زمان ۱ و ۲ تفاوت چندانی ندارد ولی در زمان ۳ به مراتب بزرگ‌تر شده است، بنابراین توجه خود را به زمان ۲ و ۳ جلب می‌کند. همچنین او مشاهده می‌کند که نیمی از ژن‌ها همواره $expression$ کمی دارند و لذا محمدرضا آن‌ها را به راحتی حذف می‌کند و تنها ۱۰۰۰ ژن را نگه می‌دارد. در آخر نیز از این ۱۰۰۰ ژن، ژنی را معرفی می‌کند که بیشترین همبستگی با تغییرات وزن در بازه‌ی زمانی ۲ تا ۳ را دارد. تحلیل محمدرضا به شدت تطبیقی است و نتایج کلاسیک آمار به راحتی قابل استفاده برای آن نیست. حدس می‌زنیم که اگر آزمایش مجدد تکرار شود، ژن انتخاب شده همبستگی کمتری با چیزی که محمدرضا به دست آورده داشته باشد. چگونه می‌توانیم این حدسیات را کمی کنیم؟

این دو مثال، نمونه‌هایی از تحلیل‌های تطبیقی و سویدگی ناشی از آن‌ها را در پردازش داده معرفی می‌کند. در این مثال‌ها می‌توان به دو شهود از این سویدگی رسید. اول این که به نظر می‌رسد که این سویدگی به توزیع خود داده ربط داشته باشد. به عنوان مثال، در تحلیل بهراد اگر یکی از ژن‌ها همبستگی خیلی بیشتر به نسبت سایر ژن‌ها داشته باشد، شهوداً با احتمال بالایی در هر تکرار مستقل آزمایش، همین ژن انتخاب می‌شود و مقدار سویدگی در این حالت کم است. دوم این که مقدار سویدگی ناشی از مرحله از

تحلیل می‌تواند با سویدگی ناشی از مراحل دیگر متفاوت باشد. به عنوان مثال در تحلیل محمدرضا، شهوداً مرحله‌ی انتخاب ژن با بیشترین expression در مرحله‌ی آخر، به نسبت سایر مراحل مولد سویدگی بیشتری است.

۲.۱ مدل ریاضی تحلیل تطبیقی داده

در مدل این مقاله از تحلیل داده، دیتاست $D \in \mathcal{D}$ با توزیع P از فضای تمام دیتاست‌های ممکن، D انتخاب شده است. تحلیل‌گر داده، تعداد زیاد m «تحلیل» مختلف را قبل از مشاهده‌ی دیتاست در نظر دارد، اما قصد دارد نتایج تنها یکی از آن‌ها را گزارش کند و این آنالیز را بعد از مشاهده‌ی D و یا آماره‌هایی از D انتخاب می‌کند. به صورت دقیق‌تر، تحلیل‌گر تعداد m تابع $\phi_1, \phi_2, \dots, \phi_m$ را در نظر گرفته که $\phi_i : \mathcal{D} \rightarrow \mathbb{R}$ است. بعد از مشاهده‌ی داده‌ها، تحلیل‌گر تابع $\phi_{T(D)}(D)$ را انتخاب کرده و آن را گزارش می‌کند. از آنجایی که خود T تابعی از D است، ممکن است سویدگی بزرگی در انتخاب $\phi_{T(D)}(D)$ وجود داشته باشد. به عنوان مثال اگر تمام ϕ_i ها دارای میانگین صفر باشد، ممکن است $\mathbb{E}[\phi_{T(D)}]$ به طرز معناداری مثبت باشد. در ادامه به بیان دو مثال (۱.۱) و (۲.۱) در قالب فرمول‌بندی مطرح شده خواهیم پرداخت.

مثال ۳.۱. در مثال (۱.۱)، دیتاست $D \in \mathbb{R}^{1000 \times 2003}$ است و شامل ۲۰۰۰ ویژگی و علاوه‌ی وزن در سه زمان مختلف است. این ویژگی‌ها نیز برای هزار نفر ثبت شده‌اند. توابع ϕ_i نیز همبستگی بین ژن‌ها و تغییر وزن از زمان ۱ تا ۳ هستند. بهراد نیز در نهایت $T = \operatorname{argmax}_i \phi_i$ را انتخاب می‌کند.

مثال ۴.۱. در مثال (۲.۱)، محمدرضا همان دیتاست بهراد است. توابع ϕ_i همبستگی هر ژن با ۳ تغییر وزن هستند. انتخاب T نیز به شیوه‌ی پیچیده‌ای انجام شده است.

میانگین واقعی تحلیل i ام را $\mu_i = \mathbb{E}[\phi_i]$ بنامید. برای یک دیتاست مشخص D ، اگر $T(D) = i$ باشد، خروجی برابر ϕ_T است. خروجی واقعی نیز μ_T است. مقدار $\phi_T - \mu_T$ خطای ناشی از تحلیل تطبیقی داده است. سویدگی تحلیل را برابر امیدریاضی این کمیت تعریف می‌کنیم، یعنی $\mathbb{E}[\phi_T - \mu_T]$ که در آن، امید ریاضی بر روی عدم قطعیت دیتاست D و روش انتخاب $T(D)$ است.

در این مقاله، کرانی برای $\mathbb{E}[\phi_T - \mu_T]$ بر مبنای «استفاده‌ی بد از اطلاعات» در انتخاب T معرفی می‌شود و نشان داده می‌شود که اگر در انتخاب T ، از اطلاعاتی از داده استفاده شود که در ϕ_i ها وجود ندارد، یعنی تنها از «اطلاعات خوب» استفاده شده باشد، این تحلیل تطبیقی منجر به سویدگی نخواهد شد و سویدگی تنها در حالتی رخ می‌دهد که در انتخاب T از اطلاعات گذشته در ϕ_i ها استفاده شود.

۳.۱ کارهای پیشین

در این بخش به بررسی اجمالی ادبیات پیشین این حوزه می‌پردازیم. خط مهمی از پژوهش در تئوری یادگیری ماشین، بررسی مفهوم پایداری الگوریتمی و استفاده از آن برای جلوگیری از over-fitting است [۲، ۳، ۴]. چهارچوب یادگیری PAC-Bayes نیز مفهومی مرتبط است که به ارائه‌ی کران‌های تعمیم الگوریتم‌های یادگیری ماشین بر حسب فاصله‌ی KL می‌پردازد [۵]. تفاوت جدی کار حاضر با کارهای پیشین این است که در روش‌های گذشته، مبتنی بر پایداری یا تئوری PAC، کران‌ها برای بدترین حالت ارائه شده و برای هر توزیع ورودی دلخواهی برقرار هستند و به طور مثال، پایداری یک الگوریتم ربطی به توزیع داده‌ای که با آن داده می‌شود ندارد. روش نظریه‌ی اطلاعاتی حاضر، از آن‌جا که بدترین حالت را در نظر نمی‌گیرد، قادر است کران‌های بهتری ارائه کند.

۲ کنترل کردن سویدگی از راه محدود کردن استفاده از اطلاعات

۱.۲ معرفی یک کران بالا برای سویدگی به وسیله‌ی اطلاعات استفاده شده

در این بخش، به معرفی یک کران بالا برای سویدگی با استفاده از یک کمیت مورد استفاده در نظریه‌ی اطلاعات می‌پردازیم. فرض کنید مجموعه داده‌ی D داده شده است و متغیرهای $\{\phi_1(D), \phi_2(D), \dots, \phi_m(D)\}$ محاسبه شده‌اند. ما یک اندیس $T(D)$ را انتخاب می‌کنیم و $\phi_{T(D)}$ را گزارش می‌کنیم. سؤال آنست که مقدار سویدگی این انتخاب $(\phi_{T(D)} - \mathbb{E}[\phi_{T(D)}])$ چقدر است؟

برای این کار، فرض کنید که Ω ، فضای نمونه‌ای مجموعه داده‌ی D باشد و $\phi = (\phi_1, \phi_2, \dots, \phi_m) : \Omega \rightarrow \mathbb{R}^m$ و $T : \Omega \rightarrow \{1, 2, \dots, m\}$ دو متغیر تصادفی باشند که روی فضای نمونه‌ای مشترک Ω تعریف شده‌اند. همچنین فرض کنید $\mu = (\mu_1, \mu_2, \dots, \mu_m) \triangleq \mathbb{E}[\phi]$ امید ریاضی بردار تصادفی ϕ باشد. در صورتی که برای هر $i \in \{1, 2, \dots, m\}$ یک متغیر تصادفی زیر-گاوسی باشد، می‌توان کران بالایی برای سوییچگی با استفاده از اطلاعات متقابل میان T و ϕ یافت.

قضیه ۱.۲. متغیر تصادفی برداری $\phi = (\phi_1, \dots, \phi_m)$ را در نظر بگیرید، اگر تعریف کنیم

$$\mu = (\mu_1, \mu_2, \dots, \mu_m) = \mathbb{E}[\phi]$$

و اگر برای هر $i \in \{1, \dots, m\}$ یک متغیر تصادفی زیر-گاوسی با پارامتر σ باشد، آن‌گاه:

$$|\mathbb{E}[\phi_T - \mu_T]| \leq \sigma \sqrt{2I(T; \phi)} \quad (۱)$$

کمیت $I(T; \phi)$ را «اطلاعات استفاده‌شده» می‌نامیم و بیان‌گر میزان وابستگی اندیس انتخابی T به مقادیر (ϕ_1, \dots, ϕ_m) است. از نظر شهودی، می‌توان ϕ_i ها را تخمین‌هایی از یک کمیت دانست و در این صورت، $I(T; \phi)$ وابستگی T به نویز موجود در این تخمین‌ها را نشان می‌دهد. همچنین می‌توان $I(T; \phi)$ را «استفاده‌ی T از اطلاعات نامید» که بیان‌گر آنست که چه میزان از نویز موجود در داده‌ها ($D \sim \mathcal{P}$) در انتخاب تخمین گزارش‌شده اثر می‌گذارد. به تعبیر دیگر، اگر ما به طمع رسیدن به بیشترین اطلاعات، T را به شدت به ϕ_i ها وابسته کنیم، هم‌زمان تأثیر اطلاعات T در انتخابی را هم افزایش داده‌ایم و این می‌تواند باعث ایجاد overfitting شود.

وقتی که T به طور کامل توسط ϕ_i ها تعیین شود، $I(T; \phi) = H(T)$ می‌شود که به معنای آنست که در تحقق‌های مختلف داده‌ها، T چگونه تغییر می‌کند.

در صورت قضیه‌ی (۱.۲)، فرض شده‌است که همه‌ی $\phi_i - \mu_i$ ها، زیر-گاوسی با پارامتر مشترک $\sigma_i = \sigma$ هستند. در حالتی که σ_i ها برابر نباشند، یک تعمیم از قضیه‌ی (۱.۲) می‌تواند به صورت $|\mathbb{E}[\phi_T - \mu_T]| \leq \max_i \{\sigma_i\} \sqrt{2I(T; \phi)}$ بیان شود، ولی کران بهتری هم وجود دارد که در قضیه‌ی (۲.۲) بیان شده‌است.

قضیه ۲.۲. متغیر تصادفی برداری $\phi = (\phi_1, \dots, \phi_m)$ را در نظر بگیرید، اگر تعریف کنیم $\mu = \mathbb{E}[\phi]$ و اگر برای هر $i \in \{1, \dots, m\}$ یک متغیر تصادفی زیر-گاوسی با پارامتر σ_i باشد، آن‌گاه:

$$|\mathbb{E}[\phi_T - \mu_T]| \leq \sqrt{\mathbb{E}[\sigma_T^2]} \sqrt{2I(T; \phi)} \quad (۲)$$

همچنین اگر فرض زیر-گاوسی بودن $\phi_i - \mu_i$ ها را کمی ضعیف کنیم و به زیر-نمایی بودن تقلیل دهیم، می‌توان قضیه‌ی مشابهی را بیان کرد.

قضیه ۳.۲. متغیر تصادفی برداری $\phi = (\phi_1, \dots, \phi_m)$ را در نظر بگیرید، اگر تعریف کنیم $\mu = \mathbb{E}[\phi]$ و اگر برای هر $i \in \{1, \dots, m\}$ یک متغیر تصادفی زیر-نمایی با پارامترهای (σ, b) باشد، آن‌گاه:

$$\mathbb{E}[\phi_T - \mu_T] \leq bI(T; \phi) + \frac{\sigma^2}{2b} \quad (۳)$$

علاوه بر این، اگر $b < 1$ باشد، داریم:

$$\mathbb{E}[\phi_T - \mu_T] \leq \sqrt{b}I(T; \phi) + \frac{\sigma^2}{2\sqrt{b}} \quad (۴)$$

برای فهم بهتر قضیه‌ی (۱.۲)، چند مثال را در نظر می‌گیریم.

مثال ۱.۲. فرض کنید $\phi_i = \frac{1}{n} \sum_{j=1}^n f_i(X_j)$ که داده‌های X_i به صورت $i.i.d.$ از توزیع $p_X(x)$ برداشته شده‌اند. در این صورت، اگر $f_i(X_j) - \mathbb{E}[f_i(X_j)]$ یک متغیر تصادفی زیر-گاوسی با پارامتر σ باشد، $\phi_i - \mu_i$ یک متغیر تصادفی زیر-گاوسی با پارامتر $\frac{\sigma}{\sqrt{n}}$ است و در نتیجه:

$$|\mathbb{E}[\phi_T - \mu_T]| \leq \sigma \sqrt{\frac{2I(T; \phi)}{n}} \quad (۵)$$

مثال ۲.۲. فرض کنید T مستقل از ϕ انتخاب شود، این حالت وقتی رخ می‌دهد که انتخاب T زودتر از مشاهده‌ی داده‌ها رخ داده و نمی‌تواند با مشاهده‌ی داده‌ها و ϕ_i ها تغییر کند. همچنین اگر داده‌ها را به دو بخش مستقل تقسیم کنیم و از یک بخش برای تعیین T استفاده کنیم و سپس ϕ_T را با استفاده از بخش دوم داده‌ها محاسبه و گزارش کنیم، در این حالت هم T از ϕ مستقل است. در این حالت $I(T; \phi) = 0$ و در نتیجه $\mathbb{E}[\phi_T] = \mathbb{E}[\mu_T]$.

مثال ۳.۲. فرض کنید که هرکدام از ϕ_i ها، یک متغیر تصادفی گاوسی با توزیع $\mathcal{N}(0, \sigma^2)$ باشند. اگر فرض کنیم $T = \arg \max_{1 \leq i \leq m} \phi_i$ ، در این صورت داریم $I(T; \phi) = H(T) = \log(m)$ ، زیرا ϕ_i ها تقارن دارند و با احتمال برابر، هرکدام می‌توانند از بقیه بیشتر باشند. در نتیجه داریم:

$$\mathbb{E}[\phi_T - \mu_T] = \mathbb{E}[\phi_T] \leq \sigma \sqrt{2 \log(m)} \quad (۶)$$

می‌توان نشان داد که با افزایش m ، این نامساوی به تساوی تبدیل می‌شود. یک حالت کلی‌تر را نیز می‌توان در نظر گرفت. فرض کنیم ابتدا ϕ_i ها از بزرگ‌ترین تا کوچک‌ترین مرتب می‌شوند، و سپس یکی از $m_0 < m$ اندیس مربوط به ϕ_i های بزرگ‌تر، با توزیع یکنواخت انتخاب می‌شود. در این حالت، داریم $I(T; \phi) = H(T) - H(T|\phi)$ و به دلیل تقارن ϕ_i ها، $H(T) = \log(m)$ و به دلیل توزیع یکنواخت T روی m_0 اندیس مربوط به ϕ_i های بزرگ‌تر، $H(T|\phi) = \log(m_0)$. در نتیجه:

$$\mathbb{E}[\phi_T - \mu_T] = \mathbb{E}[\phi_T] \leq \sigma \sqrt{2 \log\left(\frac{m}{m_0}\right)} \quad (۷)$$

۲.۲ اطلاعات استفاده شده، به عنوان کرانی برای بقیه‌ی معیارهای سویدگی

در بخش قبل، مشاهده کردیم که اطلاعات استفاده شده می‌تواند به عنوان کران بالایی برای $|\mathbb{E}[\phi_T - \mu_T]|$ به کار رود. گاهی وقت‌ها ما به معیارهای دیگری برای سویدگی، نیاز داریم، مانند $\mathbb{E}[(\phi_T - \mu_T)^2]$ و $\mathbb{E}[\phi_T - \mu_T]$. در این بخش می‌خواهیم به کمک $I(T; \phi)$ و $\sqrt{I(T; \phi)}$ ، کرانی برای بقیه‌ی معیارهای سویدگی هم بیابیم.

قضیه ۴.۲. متغیر تصادفی برداری $\phi = (\phi_1, \dots, \phi_m)$ را در نظر بگیرید، اگر تعریف کنیم $\mu = \mathbb{E}[\phi]$ و اگر برای هر $i \in \{1, \dots, m\}$ یک متغیر تصادفی زیر-گاوسی با پارامتر σ باشد، آن‌گاه:

$$\mathbb{E}[|\phi_T - \mu_T|] \leq \sigma + c\sigma \sqrt{2I(T; \phi)} \quad (۸)$$

که $c < 36$ یک ثابت است.

قضیه ۵.۲. متغیر تصادفی برداری $\phi = (\phi_1, \dots, \phi_m)$ را در نظر بگیرید، اگر تعریف کنیم $\mu = \mathbb{E}[\phi]$ و اگر برای هر $i \in \{1, \dots, m\}$ یک متغیر تصادفی زیر-گاوسی با پارامتر σ باشد، آن‌گاه:

$$\mathbb{E}[(\phi_T - \mu_T)^2] \leq 1.25\sigma^2 + c_2\sigma^2 I(T; \phi) \quad (۹)$$

که $c \leq 10$ یک ثابت است.

۳.۲ اطلاعات استفاده شده به عنوان یک کران پایین برای سویدگی

در این بخش می‌خواهیم از اطلاعات متقابل به عنوان کران پایینی برای سویدگی استفاده کنیم. فرض کنید $\phi \sim \mathcal{N}(\mu, I)$ و $T = \arg \max_i \phi_i$ ، از آنجا که T یک تابع یقینی از ϕ است، $I(T; \phi) = H(T)$. در نتیجه از قضیه‌ی (۵.۲) می‌توان نوشت:

$$\mathbb{E}[(\phi_T - \mu_T)^2] \leq \sigma^2(1.25 + 10I(T; \phi)) = \sigma^2(1.25 + 10H(T)) \quad (۱۰)$$

در این قسمت نشان می‌دهیم که $H(T)$ می‌تواند به عنوان یک کران پایین از $\mathbb{E}[(\phi_T - \mu_T)^2]$ هم به کار رود.

قضیه ۶.۲. فرض کنید $T = \arg \max_i \phi_i$ که در آن، $\phi \sim \mathcal{N}(\mu, I)$ در این صورت ثابت‌های $c_2 < 2.5$ ، $c_1 = \frac{1}{8}$ و $c_3 = 10$ و $c_4 = 1.5$ وجود دارند، به قسمی که:

$$c_1 H(T) - c_2 \leq \mathbb{E}[(\phi_T - \mu_T)^2] \leq c_3 H(T) + c_4 \quad (۱۱)$$

برای به دست آوردن یک دید شهودی از این رابطه، یک حالت ساده را در نظر می‌گیریم: فرض کنید $m = 2$ ، $\phi_1 = x$ که x یک مقدار ثابت است و $\phi_2 \sim \mathcal{N}(0, 1)$. همچنین فرض کنید $T = \arg \max_i \phi_i$ اگر $x \gg 0$ باشد، داریم:

$$\log\left(\frac{1}{\mathbb{P}[T=2]}\right) = \log\left(\frac{1}{\mathbb{P}[\phi_2 \geq x]}\right) \approx \frac{x^2}{2}$$

در نتیجه در حالتی که $x \rightarrow \infty$ می‌توان نوشت:

$$H(T_x) \sim \mathbb{P}[T_x = 2] \log\left(\frac{1}{\mathbb{P}[T_x = 2]}\right) \sim \mathbb{P}[T_x = 2] x^2 \sim \mathbb{E}[(\phi_{T_x} - \mu_{T_x})^2]$$

که T_x به معنای اندیس انتخاب شده است، وقتی که $\phi_1 = x$ باشد و $f(x) \sim g(x)$ به معنای آنست که اگر $x \rightarrow \infty$ آنگاه $\frac{f(x)}{g(x)} \rightarrow 1$

۳ چه زمانی سویدگی بزرگ است؟

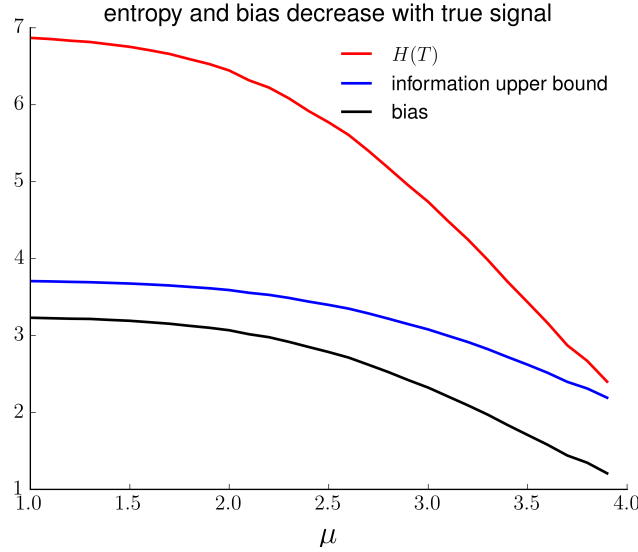
در این بخش، به بررسی چند روش ساده، ولی معمول و پر استفاده در انتخاب ویژگی و تخمین پارامتر می‌پردازیم. در بسیاری از کاربردها، انتخاب تحلیل و تخمین آن بر روی همان دیتاست اصلی انجام می‌شوند. روش بهره‌برداری از اطلاعات، یک چهارچوب مناسب برای بررسی میزان سویدگی نتایج حاصل از چنین مطالعاتی است. در این بخش بررسی می‌کنیم که هر روش، تحت چه شرایطی منجر به سویدگی شده و تحت چه شرایطی سویدگی آن قابل صرف نظر و کوچک است.

۱.۳ رتبه‌بندی با وجود سیگنال در داده‌ها

مثالی از مسئله‌ی رتبه‌بندی، انتخاب K مقدار i با بیشترین ϕ_i است. مثالی خاصی از این مسئله، یافتن ماکزیمم مقادیر ϕ_i می‌باشد. در این بخش نشان می‌دهیم وجود سیگنال (در مقابل نویز) در داده، باعث کم شدن سویدگی در مسئله‌ی رتبه‌بندی می‌شود. اطلاعات متقابل $I(T; \phi)$ کران دار است: $I(T; \phi) \leq H(T) \leq \log(m)$. به صورت شهودی، هنگامی که در داده، سیگنال بزرگی از این که کدام T باید انتخاب شود در بر دارد، توزیع T از توزیع یونیفرم دور است و اطلاعات متقابل از $\log(m)$ کمتر است. مثال زیر را در نظر بگیرید:

$$\phi_i \sim \begin{cases} \mathcal{N}(\mu, \sigma^2) & i = I^* \\ \mathcal{N}(0, \sigma^2) & i \neq I^* \end{cases} \quad (۱۲)$$

که $\mu > 0$ است. تحلیل‌گر داده قصد دارد I^* را کشف کرده و مقدار ϕ_{I^*} را گزارش کند. برای این کار $T = \arg \max_i \phi_i$ اگر $\mu = 0$ باشد، سیگنالی در داده وجود ندارد و در نتیجه تحلیل‌گر به صورت یکنواخت یکی از مقادیر $i \in \{1, 2, \dots, m\}$ را انتخاب می‌کند، در این حالت $I(T; \phi) = H(T) = \log(m)$ است. با افزایش μ ، متغیر T بر روی I^* متمرکز می‌شود و این منجر به کاهش $H(T)$ و در نتیجه کاهش سویدگی $\mathbb{E}[\phi_T - \mu_T]$ می‌شود. در یک شبیه‌سازی، $m = 1000$ نمونه از ϕ_i ها تولید شده‌اند که همه‌ی آن‌ها به جز یکی از توزیع نرمال استاندارد تولید شده‌اند، اما یکی از آن‌ها توزیع $\phi_{I^*} \sim \mathcal{N}(\mu, 1)$ دارد. مقدار μ در بازه‌ی $[0, 4]$ تغییر داده شده است و میانگین نتیجه در ۱۰۰۰ تکرار در شکل (۱) آمده است. این نتایج، به وضوح درستی بحث فوق را تایید می‌کند. با افزایش μ کران بالای سویدگی و مقدار واقعی سویدگی کاهش یافته است.



شکل ۱: نمودار میزان سویدگی ϕ_T و کران‌های آن بر حسب μ

۲.۳ تفکیک تقریباً مستقل داده‌های غیر i.i.d.

در این مسئله، تحلیل‌گر داده، n نمونه‌ی زمانی از یک زنجیره‌ی مارکوفی را در اختیار دارد. این تحلیل‌گر قصد دارد از این داده به همان نحوی که از داده‌های i.i.d. استفاده می‌کند، بهره‌برد. برای این کار، او داده‌ها را به سه قسمت تقسیم می‌کند: از s_1, \dots, s_{n_1} برای انتخاب تحلیل استفاده می‌کند، داده‌های $s_{n_1+1}, \dots, s_{n_2}$ را دور می‌اندازد و از s_{n_2+1}, \dots, s_n برای تخمین تحلیل انتخاب‌شده با استفاده از بخش اول داده‌ها بهره می‌برد. اگر داده‌ها i.i.d. می‌بودند، هیچ نشت اطلاعاتی رخ نمی‌داد و این تحلیل هیچ سویدگی‌ای ایجاد نمی‌کرد. اما در این حالت که در دیتا روابط زمانی وجود دارد، این گزاره صادق نیست. از نظر شهودی انتظار داریم که اگر $n_2 - n_1$ به قدر کافی بزرگ باشد، سویدگی تحلیل به قدر دلخواه کم شود. فرض کنید زنجیره‌ی مارکوف ایستان بوده و توزیع ایستان آن را π بنامید. همچنین فرض کنید:

$$\forall \tau \in \mathbb{N} \quad \max_s D(\mathbf{P}(s_\tau = \cdot | s_1 = s) || \pi) \leq c_0 e^{-c_1 \tau} \quad (۱۳)$$

یعنی توزیع به شرط حالت اولیه نیز به توزیع ایستان میل می‌کند. از آنجایی که این فرآیند، یک فرآیند ایستان است، داریم $\mathbf{P}(s_t = s) = \pi$ با این فرض و با کمک نامساوی (۱۳) می‌توان نوشت:

$$I(s_{t+\tau}; s_t) = \sum_s \mathbf{P}(s_t = s) D(\mathbf{P}(s_{t+\tau} = \cdot | s_t = s) || \mathbf{P}(s_t = \cdot)) \leq c_0 e^{-c_1 \tau}$$

با توجه به نامساوی پردازش اطلاعات، داریم:

$$I(T; \phi) \leq I(s_1, \dots, s_{n_1}; s_{n_2+1}, \dots, s_n) \leq I(s_{n_1}; s_{n_2+1}) \leq c_0 e^{-c_1(n_2-n_1)} \quad (۱۴)$$

با ترکیب این کران با قضیه‌ی (۱.۲)، شهود ما اثبات می‌شود یعنی هرچقدر $n_2 - n_1$ بزرگتر باشد، سویدگی کمتر می‌شود.

۴ محدود کردن سویدگی با تصادفی‌سازی

در این بخش، نشان می‌دهیم که حتی اگر فرآیند انتخاب T ، فرآیندی پیچیده بوده و یا ناشناخته باشد، می‌توان با افزودن نویز در مرحله‌ی انتخاب، تا حد زیادی سویدگی را کم کرد. روش‌های مبتنی بر تصادفی‌سازی، تاکنون به شدت در ادبیات یادگیری ماشین و آمار تکرار شده‌اند. در این بخش قصد داریم به تحلیل این روش‌ها بر اساس معیار «اطلاعات بد» معرفی شده در فصل‌ها بپردازیم.

۱.۴ رگولاریزیشن با انتخاب تصادفی

فرض کنید m تحلیل مختلف $\phi_1, \phi_2, \dots, \phi_m$ در اختیار داریم. قصد داریم تحلیلی را بیابیم که بزرگترین مقدار ممکن را دارد. برای این کار، یک راه اولیه، انتخاب $T = \operatorname{argmax}_i \phi_i$ است. با توجه به (۱.۲)، می‌دانیم که اگر $I(T; \phi)$ کوچک باشد، سویدگی تخمین کم می‌شود. داریم $I(T; \phi) = H(T) - H(T|\phi)$. با توجه به این رابطه، به نظر می‌رسد که می‌توان با افزایش $H(T|\phi)$ یا به عبارتی، تصادفی‌سازی انتخاب T ، میزان سویدگی را کنترل کرد. با انتخاب T به صورت تصادفی و مستقل از داده، $H(T|\phi)$ را بیشینه می‌کند، اما از آنجایی که $H(T)$ را هم زیاد می‌کند، به کمترین $I(T; \phi)$ منتج نمی‌شود. یک ایده این است که علاوه بر افزایش $H(T|\phi)$ ، مقدار ϕ_T را هم تا حد امکان بزرگ انتخاب کنیم. بعد از مشاهده ϕ ها، اندیس T را به صورت تصادفی و با توزیع π از $\{1, 2, \dots, m\}$ انتخاب می‌کنیم. از آنجا که قصد داریم ϕ_T بزرگ باشد، منطقی است که توزیع π ای را انتخاب کنیم که جواب مسئله‌ی بهینه‌سازی زیر باشد:

$$\begin{aligned} & \underset{\pi \in \mathbb{R}_+^m}{\text{maximize}} && H(\pi) \\ & \text{subject to} && \sum_{i=1}^k \pi_i \phi_i \geq b \text{ and } \sum_{i=1}^k \pi_i = 1. \end{aligned}$$

شرط $\sum_{i=1}^k \pi_i \phi_i \geq b$ تضمین می‌کند که توزیع انتخاب‌شده به نحوی باشد که به صورت متوسط، مقدار ϕ_T بزرگ باشد. این یک مسئله‌ی کلاسیک در تئوری اطلاعات است و می‌دانیم پاسخ آن، توزیع گیس است یعنی توزیع $\pi^* = Ae^{\beta \phi_i}$. در این توزیع، β به نحوی انتخاب شده است که شرط $\sum_{i=1}^k \pi_i \phi_i = b$ برقرار باشد. اگر دو دیتاست در نظر بگیریم که تنها تفاوت کوچکی با یکدیگر دارند، ممکن است $T = \operatorname{argmax}_i \phi_i$ در این دو دیتاست تفاوت بزرگی داشته باشد، اما روش تصادفی فوق، احتمالاً جواب‌های نزدیکی را در این دو دیتاست باز خواهد گرداند. به کمک یک شبیه‌سازی، نشان می‌دهیم این روش از سویدگی کم می‌کند. فرض کنید دو دسته ϕ_i داریم. برای تحلیل‌های $1 \leq i \leq N_1$ داریم $\mu_i = \mu > 0$ و برای تحلیل‌های $N_1 + 1 \leq i \leq m$ نیز $\mu_i = 0$ است. در شبیه‌سازی گرفته‌ایم $N_1 = 1000$ و $m - N_1 = 100000$ ، یعنی تعداد تحلیل‌ها با میانگین صفر بسیار بیشتر از تعداد تحلیل‌های با میانگین ناصفر است. مقدار μ را نیز در بازه‌ی 1 تا 5 تغییر می‌دهیم. در این تحلیل به جای گزارش کردم ماکزیمم، K تحلیل با مقدار بزرگتر را گزارش می‌کنیم. در این شبیه‌سازی به دو کمیت زیر علاقه‌مندیم:

○ سویدگی: $\frac{1}{K} \sum_{i=1}^K \phi_{T_i} - \mu_{T_i}$

○ دقت: $\frac{|\{T_i: T_i \leq N_1\}|}{K}$

نتایج این شبیه‌سازی در نمودار (۲) آمده است. در این نمودار دیده می‌شود که در رژیم با μ کوچک، مطابق انتظار هر دو روش دقت کمی دارند. در رژیم با μ بزرگ، هر دو روش دقتی نزدیک به ۱ دارند ولی روش تصادفی Max Entropy سویدگی کمتری را تجربه کرده است. در رژیم $1 < \mu < 4$ روش Max Entropy خطای بیشتر و در عین حال، سویدگی کمتری دارد.

تذکر ۱.۴. توجه نمایید که الگوریتم فوق، لزوماً منجر به کم شدن سویدگی نمی‌شود، زیرا در عین حالی که $H(T|\phi)$ را افزایش می‌دهد، مقدار $H(T)$ را نیز زیاد می‌کند.

۲.۴ تصادفی‌سازی در تحلیل‌های چند مرحله‌ای

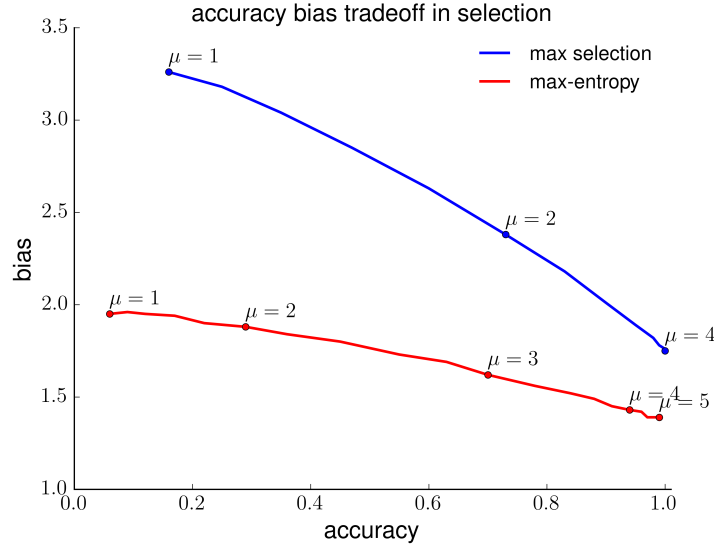
در ابتدا مدلی از تحلیل داده‌ی تطبیقی ارائه می‌کنیم در آن، تحلیل‌گر داده با انجام آنالیزهای پی‌درپی و استفاده از نتایج آنالیزهای قبلی، سعی می‌کند به شناختی از داده برسد.

○ در قدم اول، تحلیل‌گر تابع ϕ_{T_1} را انتخاب می‌کند. نتیجه‌ی این تحلیل $Y_{T_1} \in \mathbb{R}$ است.

○ در تکرار k ام، تحلیل‌گر تابع ϕ_{T_k} را انتخاب می‌کند و در این انتخاب خود، از متغیرهای

$$\{Y_{T_1}, Y_{T_2}, \dots, Y_{T_{k-1}}, T_1, \dots, T_{k-1}\}$$

استفاده می‌کند.



شکل ۲: نتایج شبیه‌سازی با تغییر μ با دو روش تصادفی این بخش و روش انتخاب ماکزیمم‌ها

ما در حالت کلی اجازه می‌دهیم که نتیجه‌ی آنالیز k ام، Y_{T_k} با ϕ_{T_k} متفاوت باشد. به عنوان مثال، ممکن است این نتیجه به صورت $Y_{T_k} = \phi_{T_k} + \text{noise}$ باشد. این سناریو زمانی رخ می‌دهد که تحلیل‌گر داده برای کاستن سویدگی، به نتایج نویز اضافه کرده و از نتیجه‌ی نویزی در مراحل بعد تحلیل خود استفاده نماید. در این مدل، بین متغیرهای مسئله، زنجیره‌ی ماکوف زیر برقرار است:

$$T_{k+1} - H_k = \{T_1, Y_{T_1}, T_2, Y_{T_2}, \dots, T_k, Y_{T_k}\} - D - \phi$$

با نامسای پرداختش اطلاعات، می‌توان گفت که $I(T_{k+1}) \leq I(H_k; \phi)$ است، بنابراین فرآیندی که در آن اطلاعات متقابل سابقه و فیدبک سیستم (H_k) و نتایج آنالیزها ϕ کنترل شده باشد، بایاس کمی دارد.

لم ۱۰۴. با تعریف $H_k = \{T_1, Y_{T_1}, T_2, Y_{T_2}, \dots, T_k, Y_{T_k}\}$ داریم:

$$I(T_{k+1}; \phi) \leq I(H_k; \phi) = \sum_{i=1}^k I(Y_{T_i}; \phi_{T_i} | H_{i-1}, T_i)$$

از لم (۱۰۴) می‌توان دریافت که با محدود کردن $I(Y_{T_i}; \phi_{T_i} | H_{i-1}, T_i)$ در هر مرحله، می‌توان $I(T_{k+1}; \phi)$ را کنترل کرد. در نتیجه انگاریک «بودجه‌ی اطلاعات استفاده‌شده» مانند I_b داریم که در هر مرحله، تحلیل‌گر برای انتخاب اندیس فعلی براساس اندیس‌های قبلی از آن خرج می‌کند. در نتیجه فرآیند را می‌توان تا جایی ادامه داد که بودجه تمام شود. اگر بتوانیم در هر مرحله، $I(Y_{T_i}; \phi_{T_i} | H_{i-1}, T_i)$ را کم کنیم، می‌توان تعداد مراحل بیشتری را پیمود. یک راه برای دستیابی به این هدف، اضافه کردن نویز گاوسی است. فرض کنید $\phi_i \sim \mathcal{N}(\mu_i, \frac{\sigma^2}{n})$ و برای هر k ، (ϕ_1, \dots, ϕ_k) مشترکاً گاوسی باشند. همچنین فرض کنید که در j امین مرحله، داریم $Y_{T_j} = \phi_{T_j} + W_j$ که W_j ها مستقلند و داریم $W_j \sim \mathcal{N}(0, \frac{\omega_j^2}{n})$ در این حالت می‌توان $\mathbb{E}[|Y_{T_{k+1}} - \mu_{T_{k+1}}|]$ را محدود کرد.

قضیه ۱۰۴. فرض کنید $\phi_i \sim \mathcal{N}(\mu_i, \frac{\sigma^2}{n})$ و برای هر k ، (ϕ_1, \dots, ϕ_k) مشترکاً گاوسی باشند. همچنین فرض کنید که به ازای هر j ، داریم $Y_{T_j} = \phi_{T_j} + W_j$ که W_j ها مستقلند و داریم $W_j \sim \mathcal{N}(0, \frac{\sigma^2 \sqrt{j}}{n})$. در این صورت برای هر $k \in \mathcal{N}$ داریم:

$$\mathbb{E}[|Y_{T_{k+1}} - \mu_{T_{k+1}}|] \leq c \left(\frac{\sigma k^{\frac{1}{4}}}{n^{\frac{1}{2}}} \right) \quad (15)$$

که c یک ثابت است و مستقل از σ, k, n .

اگر دنباله‌ی (T_1, T_2, \dots) به صورت غیر تطبیقی تولید می‌شدند و نویزی نیز اضافه نمی‌شد، به کران $\mathbb{E}[|Y_{T_{k+1}} - \mu_{T_{k+1}}|] \leq \frac{\sigma}{n}$ می‌رسیدیم. در مدل تطبیقی و در مراحل ابتدایی نیز می‌توانستیم به این سویدگی برسیم، اما به مرور میزان سویدگی افزایش پیدا می‌کرد. ضرب $k^{\frac{1}{4}}$ را می‌توان به عنوان هزینه‌ی پردازش تطبیقی در نظر گرفت. در این جا ذکر این نکته جالب است که می‌توان مثال‌هایی تطبیقی یافت که بدون افزودن نویز، در آن‌ها سویدگی از حالت غیر تطبیقی بیشتر است، یعنی $\mathbb{E}[|Y_{T_{k+1}} - \mu_{T_{k+1}}|] = \Omega(\sigma \frac{k}{n})$.

۵ پیشنهاد برای کارهای آتی

برای انجام کارهای آتی چندین مسیر امیدوارکننده به نظر می‌رسد که در این بخش، به صورت اجمالی به آن‌ها اشاره می‌کنیم.

- در مدل‌سازی این مقاله، تحلیل‌گر داده همواره تنها یک تابع را بر می‌گزیند و مقدار آن را گزارش می‌کند. در کاربردهای تجربی، بسیار پیش می‌آید که تعدادی تحلیل انتخاب شده و نتایج آن گزارش شوند. به عنوان مثال، ممکن است این تعداد با توجه به داده، تعداد تحلیل‌های جالب مشخص شود. این انتخاب تصادفی تعداد داده‌های گزارش شده نیز می‌تواند باعث نشر «اطلاعات بد» از دیتاست شده و باعث افزایش سویدگی گردد. تعمیم چهارچوب مطرح شده در این مقاله برای حالتی که تحلیل‌گر تعدادی تحلیل را گزارش می‌کند و خود این تعداد نیز یک متغیر تصادفی است، از اهمیت تجربی و تئوری بالایی برخوردار است.

- در این مقاله و در مدل‌سازی تحلیل تطبیقی داده‌ها، همواره فرض بر این است که دیتاست در طول زمان ثابت است. این فرض در بسیاری از کاربردها از واقعیت به دور است. فرض به‌ترین است که خود دیتاست نیز ماتریسی تصادفی در نظر گرفته شود که در طول زمان در حال تغییر است. مبه عنوان مثال، این ماتریس می‌تواند یک فرآیند تصادفی مارکوف را تشکیل دهد. مدل‌سازی مسئله در این شرایط و با اعمال محدودیت‌هایی از تغییرات ممکن در دیتاست، می‌تواند به نتایج بسیار با اهمیتی منجر شود. در این مدل جدید می‌توان به بررسی این موضوع پرداخت که استفاده داده‌های دیتاست در بخش از زمان‌ها برای انتخاب تحلیل و استفاده از زمان‌هایی دیگر برای تخمین آن، به چه میزان می‌تواند ایجاد سویدگی کند.

- در مدل تحلیل تطبیقی داده، به خصوص در مثال متغیرهای گاوسی، برای ارائه‌ی کران سویدگی، توزیع نویز به نحوی انتخاب شده که محاسبات ساده گردند. پیدا کردن توزیع‌های نویزی که منجر به کوچک‌ترین سویدگی شوند می‌تواند بسیار جالب و کاربردی باشد.

- در تئوری یادگیری ماشین، الگوریتم‌های آنلاین متعددی وجود دارند که به صورت «تطبیقی» به پردازش داده‌ها می‌پردازند. چهارچوب فعلی می‌تواند به عنوان ابزاری نیرومند برای تحلیل بایاس‌های چنین الگوریتم‌های به کار رود.

مراجع

- [1] D. Russo and J. Zou, "How much does your data exploration overfit? controlling bias via information usage," *IEEE Transactions on Information Theory*, vol.66, no.1, pp.302–323, 2019.
- [2] O. Bousquet and A. Elisseeff, "Stability and generalization," *Journal of Machine Learning Research*, vol.2, no.Mar, pp.499–526, 2002.
- [3] T. Poggio, R. Rifkin, S. Mukherjee, and P. Niyogi, "General conditions for predictivity in learning theory," *Nature*, vol.428, no.6981, pp.419–422, 2004.
- [4] S. Shalev-Shwartz, O. Shamir, N. Srebro, and K. Sridharan, "Learnability, stability and uniform convergence," *Journal of Machine Learning Research*, vol.11, no.Oct, pp.2635–2670, 2010.
- [5] D. McAllester, "A pac-Bayesian tutorial with a dropout bound," *arXiv preprint arXiv:1307.2118*, 2013.

- [6] M. Mohri, A. Rostamizadeh, and A. Talwalkar, “Foundations of machine learning,” 2018.

آ پیش‌نیازهای موردنیاز جهت فهم اثبات‌های قضایا

۱.آ متغیرهای تصادفی زیر-گوسی

تعریف ۱.آ. متغیر تصادفی X با میانگین $\mu = \mathbb{E}[X]$ را زیر-گوسی می‌نامیم هرگاه عدد مثبتی مانند σ وجود داشته باشد، به قسمی که برای هر $\lambda \in \mathbb{R}$ داشته باشیم:

$$\mathbb{E} \left[e^{\lambda(X-\mu)} \right] \leq e^{\frac{\sigma^2 \lambda^2}{2}}$$

ثابت σ را پارامتر این متغیر تصادفی می‌نامیم. به عنوان مثال، یک متغیر تصادفی گوسی با واریانس σ^2 ، خود یک متغیر تصادفی زیر-گوسی با پارامتر σ است. همچنین تعداد زیادی از متغیرهای تصادفی غیر گوسی، زیر-گوسی هستند.

قضیه ۱.آ. برای هر متغیر تصادفی زیر-گوسی X با متوسط $\mu = \mathbb{E}[X]$ و پارامتر σ داریم:

$$\mathbb{P}[|X - \mu| \geq t] \leq 2e^{-\frac{t^2}{2\sigma^2}} \quad (۱۶)$$

اثبات. از نامساوی مارکف می‌دانیم:

$$\mathbb{P}[X - \mu \geq t] = \mathbb{P}[e^{\lambda(X-\mu)} \geq e^{\lambda t}] \leq \frac{\mathbb{E} \left[e^{\lambda(X-\mu)} \right]}{e^{\lambda t}}$$

حال با توجه به تعریف متغیرهای تصادفی زیر-گوسی داریم:

$$\mathbb{P}[X - \mu \geq t] \leq \frac{\mathbb{E} \left[e^{\lambda(X-\mu)} \right]}{e^{\lambda t}} \leq \exp\left(\frac{\sigma^2 \lambda^2}{2} - \lambda t\right)$$

نامساوی بالا به ازای هر $\lambda \in \mathbb{R}$ برقرار است، من جمله λ ای که طرف راست را کمینه کند. در نتیجه:

$$\mathbb{P}[X - \mu \geq t] \leq \inf_{\lambda} \left\{ \exp\left(\frac{\sigma^2 \lambda^2}{2} - \lambda t\right) \right\} = e^{-\frac{t^2}{2\sigma^2}} \quad (۱۷)$$

همچنین اگر متغیر تصادفی X ، زیر-گوسی باشد، $-X$ هم زیر-گوسی است و به طور مشابه، داریم:

$$\mathbb{P}[-X + \mu \geq t] = \mathbb{P}[X - \mu \leq -t] \leq e^{-\frac{t^2}{2\sigma^2}}$$

و می‌توان نوشت:

$$\mathbb{P}[|X - \mu| \geq t] \leq \mathbb{P}[X - \mu \geq t] + \mathbb{P}[X - \mu \leq -t] \leq 2e^{-\frac{t^2}{2\sigma^2}}$$

□

متغیرهای تصادفی زیر-گوسی خواص گوناگونی دارند، تعدادی از این خواص را در قضایای بعدی مشاهده می‌کنیم.

قضیه ۲.آ. فرض کنید X یک متغیر تصادفی زیر-گوسی با امید ریاضی $\mathbb{E}[X] = 0$ باشد، در این صورت اگر متغیر تصادفی Z را به صورت $Z \sim \mathcal{N}(0, 2\sigma^2)$ در نظر بگیریم، داریم:

$$\mathbb{P}[|X| \geq s] \leq \sqrt{8e} \mathbb{P}[|Z| \geq s] \quad \forall s \geq 0 \quad (۱۸)$$

اثبات. از قضیه ۱.۲ داریم:

$$\mathbb{P}[X \geq t] \leq e^{-\frac{t^2}{2\sigma^2}} \quad \forall t \geq 0$$

و از طرف دیگر، با توجه به کران Mills ratio برای توزیع‌های گاوسی داریم:

$$\mathbb{P}[Z \geq t] \geq \left(\frac{\sqrt{2}\sigma}{t} - \frac{(\sqrt{2}\sigma)^3}{t^3} \right) e^{-\frac{t^2}{4\sigma^2}}$$

حال دو حالت زیر را در نظر می‌گیریم: حالتی که $t \in [0, 2\sigma]$ باشد. در این حالت، داریم:

$$\mathbb{P}[Z \geq t] \geq \mathbb{P}[Z \geq 2\sigma] = \left(\frac{1}{\sqrt{2}} - \frac{1}{2\sqrt{2}} \right) e^{-1} = \frac{1}{\sqrt{8}e}$$

و از آنجا که $\mathbb{P}[X \geq t] \leq 1$ داریم:

$$\frac{\mathbb{P}[X \geq t]}{\mathbb{P}[Z \geq t]} \leq \sqrt{8}e$$

حالتی که $t > 2\sigma$ باشد. در این حالت اگر کران Mills ratio را با کران به دست آمده در قضیه ۱.۲ ترکیب کنیم و تعریف کنیم $s = \frac{t}{\sigma}$ ، داریم:

$$\begin{aligned} \sup_{t > 2\sigma} \frac{\mathbb{P}[X \geq t]}{\mathbb{P}[Z \geq t]} &\leq \sup_{s > 2} \frac{e^{-\frac{s^2}{4}}}{\frac{\sqrt{2}}{s} - \frac{2\sqrt{2}}{s^3}} \\ &\leq \sup_{s > 2} s^3 e^{-\frac{s^2}{4}} \\ &\leq \sqrt{8}e \end{aligned}$$

□

پس در هر دو حالت نامساوی (۱۸) برقرار است.

قضیه ۳.۲. فرض کنید X یک متغیر تصادفی با امید ریاضی $\mathbb{E}[X] = 0$ باشد و بتوانیم عدد ثابتی مانند c و یک متغیر تصادفی $Z \sim \mathcal{N}(0, \tau^2)$ بیابیم، به قسمی که

$$\mathbb{P}[|X| \geq s] \leq c\mathbb{P}[|Z| \geq s] \quad \forall s \geq 0 \quad (۱۹)$$

در این صورت ثابتی مانند $\theta \geq 0$ وجود دارد، به گونه‌ای که:

$$\mathbb{E}[X^{2k}] \leq \frac{(2k)!}{2^k k!} \theta^{2k} \quad \forall k \in \mathbb{N} \quad (۲۰)$$

اثبات. اگر $Z \sim \mathcal{N}(0, \tau^2)$ و X یک متغیر تصادفی باشد که در رابطه‌ی $\mathbb{P}[|X| \geq s] \leq c\mathbb{P}[|Z| \geq s]$ به ازای هر $s \geq 0$ صدق کند، با توجه به این که X^{2k} یک متغیر تصادفی نامنفی است، می‌توان نوشت:

$$\begin{aligned} \mathbb{E}[X^{2k}] &= \int_0^{+\infty} \mathbb{P}[X^{2k} > s] ds \\ &= \int_0^{+\infty} \mathbb{P}[|X| > s^{\frac{1}{2k}}] ds \\ &\leq c \int_0^{+\infty} \mathbb{P}[|Z| > s^{\frac{1}{2k}}] ds \\ &= c\mathbb{E}[Z^{2k}] \end{aligned}$$

و از آن جا که متغیر تصادفی Z ، گاوسی است، می‌دانیم:

$$\mathbb{E} [Z^{2k}] = \frac{(2k)!}{2^k k!} \tau^{2k} \quad \forall k \in \mathcal{N}$$

در نتیجه:

$$\begin{aligned} \mathbb{E} [X^{2k}] &\leq c \mathbb{E} [Z^{2k}] \\ &= c \frac{(2k)!}{2^k k!} \tau^{2k} \\ &\leq \frac{(2k)!}{2^k k!} (c\tau)^{2k} \end{aligned}$$

□

در نتیجه اگر قرار دهیم $\theta = c\tau$ به حکم می‌رسیم.

قضیه ۴.۵. فرض کنید X یک متغیر تصادفی با امید ریاضی $\mathbb{E}[X] = 0$ باشد و ثابتی مانند $\theta \geq 0$ وجود داشته باشد، به قسمی که:

$$\mathbb{E} [X^{2k}] \leq \frac{(2k)!}{2^k k!} \theta^{2k} \quad \forall k \in \mathbb{N} \quad (21)$$

در این صورت متغیر تصادفی X ، زیر گاوسی با پارامتر $\sigma = \theta\sqrt{2}$ است.

اثبات. به ازای هر $\lambda \in \mathbb{R}$ داریم:

$$\mathbb{E} [e^{\lambda X}] \leq 1 + \sum_{k=2}^{\infty} \frac{|\lambda|^k \mathbb{E} [|X|^k]}{k!}$$

(جمله‌ی متناظر با $k = 1$ به دلیل این که $\mathbb{E}[X] = 0$ حذف شده است) اگر X حول صفر تقارن داشته باشد، جملات دارای k فرد از بین می‌روند و داریم:

$$\mathbb{E} [e^{\lambda X}] \leq 1 + \sum_{k=1}^{\infty} \frac{\lambda(2k) \mathbb{E} [|X|^{2k}]}{(2k)!} \frac{(2k)! \theta^{2k}}{2^k k!} = e^{\frac{\lambda^2 \theta^2}{2}}$$

که نتیجه می‌دهد X یک متغیر تصادفی زیر-گاوسی با پارامتر θ است.

اگر X متقارن نباشد، می‌توانیم یک کران بالا برای جملات مربوط به k های فرد بیابیم:

$$\begin{aligned} \mathbb{E} [|\lambda X|^{2k+1}] &\leq \sqrt{\mathbb{E} [|\lambda X|^{2k}] \mathbb{E} [|\lambda X|^{2k+2}]} \\ &\leq \frac{1}{2} \left(\lambda^{2k} \mathbb{E} [X^{2k}] + \lambda^{2k+2} \mathbb{E} [X^{2k+2}] \right) \end{aligned}$$

(نامساوی اول از نامساوی کوشی-شوارتز نتیجه می‌شود و نامساوی دوم از نامساوی میانگین حسابی-هندسی.) در نتیجه داریم:

$$\begin{aligned} \mathbb{E} [|\lambda X|^{2k+1}] &\leq 1 + \left(\frac{1}{2} + \frac{1}{2 \times 3!} \right) \lambda^2 \mathbb{E} [X^2] + \\ &\quad + \sum_{k=2}^{\infty} \left(\frac{1}{(2k)!} + \frac{1}{2} \left[\frac{1}{(2k-1)!} + \frac{1}{(2k+1)!} \right] \right) \lambda^{2k} \mathbb{E} [X^{2k}] \\ &\leq \sum_{k=0}^{\infty} 2^k \frac{\lambda^{2k} \mathbb{E} [X^{2k}]}{(2k)!} \\ &\leq \sum_{k=0}^{\infty} 2^k \frac{\lambda^{2k}}{(2k)!} \frac{(2k)! \theta^{2k}}{2^k k!} \\ &= \exp\left(\frac{(\sqrt{2}\lambda\theta)^2}{2}\right) \end{aligned}$$

□ که نتیجه می‌دهد X یک متغیر تصادفی زیر-گوسی با پارامتر $\theta\sqrt{2}$ است.

قضیه ۵.۲. فرض کنید X یک متغیر تصادفی با امید ریاضی $\mathbb{E}[X] = 0$ باشد و بتوانیم عدد ثابتی مانند c و یک متغیر تصادفی $Z \sim \mathcal{N}(0, \tau^2)$ بیابیم، به قسمی که

$$\mathbb{P}[|X| \geq s] \leq c\mathbb{P}[|Z| \geq s] \quad \forall s \geq 0 \quad (22)$$

در این صورت X یک متغیر تصادفی زیر-گوسی با پارامتر $\sqrt{2}c\tau$ است.

اثبات. با توجه به اینکه

$$\mathbb{P}[|X| \geq s] \leq c\mathbb{P}[|Z| \geq s] \quad \forall s \geq 0$$

با استفاده از قضیه ۳.۲ داریم:

$$\mathbb{E}[X^{2k}] \leq \frac{(2k)!}{2^k k!} \theta^{2k} \quad \forall k \in \mathbb{N}$$

که $\theta = c\tau$. حال با توجه به قضیه ۴.۲ می‌توانیم نتیجه بگیریم که متغیر تصادفی X ، یک متغیر تصادفی زیر-گوسی با پارامتر $\sigma = \sqrt{2}\theta = \sqrt{2}c\tau$ است. □

قضیه ۶.۲. فرض کنید X یک متغیر تصادفی زیر-گوسی با امید ریاضی $\mathbb{E}[X] = 0$ و پارامتر σ باشد، در این صورت:

$$\mathbb{E}\left[e^{\frac{sX^2}{2\sigma^2}}\right] \leq \frac{1}{\sqrt{1-s}} \quad \forall s \in [0, 1) \quad (23)$$

اثبات. درستی حکم برای $s = 0$ واضح است. برای $s \in (0, 1)$ ، از آنجا که X یک متغیر تصادفی زیر-گوسی است، داریم:

$$\mathbb{E}\left[e^{\lambda X}\right] \leq e^{\frac{\lambda^2 \sigma^2}{2}}$$

دو طرف نامساوی فوق را در $e^{-\frac{\lambda^2 \sigma^2}{2s}}$ ضرب می‌کنیم:

$$\mathbb{E}\left[e^{\lambda X - \frac{\lambda^2 \sigma^2}{2s}}\right] \leq e^{\frac{\lambda^2 \sigma^2 (s-1)}{2s}}$$

از آنجا که رابطه‌ی اخیر به ازای هر $\lambda \in \mathbb{R}$ برقرار است، می‌توانیم از دوطرف آن روی λ انتگرال بگیریم. انتگرال سمت راست نامساوی عبارت است از:

$$\int_{-\infty}^{\infty} \exp\left(\frac{\lambda^2 \sigma^2 (s-1)}{2s}\right) d\lambda = \frac{1}{\sigma} \sqrt{\frac{2\pi s}{1-s}}$$

و انتگرال سمت چپ نامساوی:

$$\begin{aligned} \int_{-\infty}^{\infty} \mathbb{E}\left[e^{\lambda X - \frac{\lambda^2 \sigma^2}{2s}}\right] d\lambda &= \mathbb{E}\left[\int_{-\infty}^{\infty} e^{\lambda X - \frac{\lambda^2 \sigma^2}{2s}} d\lambda\right] \\ &= \mathbb{E}\left[\frac{\sqrt{2\pi s}}{\sigma} e^{\frac{sX^2}{2\sigma^2}}\right] \\ &= \frac{\sqrt{2\pi s}}{\sigma} \mathbb{E}\left[e^{\frac{sX^2}{2\sigma^2}}\right] \end{aligned}$$

در نتیجه:

$$\begin{aligned}\mathbb{E}\left[e^{\frac{sX^2}{2\sigma^2}}\right] &= \frac{\sigma}{\sqrt{2\pi s}} \int_{-\infty}^{\infty} \mathbb{E}\left[e^{\lambda X - \frac{\lambda^2 \sigma^2}{2s}}\right] d\lambda \\ &\leq \frac{\sigma}{\sqrt{2\pi s}} \int_{-\infty}^{\infty} \exp\left(-\frac{\lambda^2 \sigma^2 (s-1)}{2s}\right) d\lambda \\ &= \frac{\sigma}{\sqrt{2\pi s}} \frac{1}{\sigma} \sqrt{\frac{2\pi s}{1-s}} \\ &= \frac{1}{\sqrt{1-s}}\end{aligned}$$

□

۲.۲ متغیرهای تصادفی زیر-نمایی

تعریف ۲.۲. متغیر تصادفی X با امید ریاضی $\mu = \mathbb{E}[X]$ را زیر-نمایی می‌نامیم اگر پارامترهای نامنفی (ν, α) وجود داشته باشند، به قسمی که برای هر λ که $|\lambda| < \frac{1}{\alpha}$ داشته باشیم:

$$\mathbb{E}\left[e^{\lambda(X-\mu)}\right] \leq e^{\frac{\nu^2 \lambda^2}{2}}$$

از تعریف متغیرهای تصادفی زیر-گاوسی و زیر-نمایی می‌توان نتیجه گرفت که هر متغیر تصادفی زیر-گاوسی با پارامتر σ زیر-نمایی هم هست با پارامترهای $\nu = \sigma$ و $\alpha = 0$ ، ولی عکس آن الزاماً درست نیست. می‌توان قضیه‌ای مشابه قضیه ۱.۲ برای متغیرهای تصادفی زیر-نمایی بیان کرد:

قضیه ۲.۲. برای هر متغیر تصادفی زیر-نمایی با متوسط $\mu = \mathbb{E}[X]$ و پارامترهای ν, α داریم:

$$\mathbb{P}[X - \mu \geq t] \leq \begin{cases} e^{-\frac{t^2}{2\nu^2}} & 0 \leq t \leq \frac{\nu^2}{\alpha} \\ e^{-\frac{t}{2\alpha}} & t > \frac{\nu^2}{\alpha} \end{cases} \quad (24)$$

اثبات. از نامساوی مارکف می‌دانیم:

$$\mathbb{P}[X - \mu \geq t] = \mathbb{P}[e^{\lambda(X-\mu)} \geq e^{\lambda t}] \leq \frac{\mathbb{E}[e^{\lambda(X-\mu)}]}{e^{\lambda t}}$$

و در نتیجه با استفاده از تعریف متغیرهای تصادفی زیر-نمایی می‌توان نوشت:

$$\mathbb{P}[X - \mu \geq t] \leq \frac{\mathbb{E}[e^{\lambda(X-\mu)}]}{e^{\lambda t}} \leq \exp\left(\frac{\nu^2 \lambda^2}{2} - \lambda t\right) \quad \forall \lambda \in [0, \frac{1}{\alpha})$$

حال می‌خواهیم مقدار کمینه‌ی عبارت $g(\lambda, t) = \frac{\nu^2 \lambda^2}{2} - \lambda t$ را محاسبه کنیم. مقدار کمینه‌ی این عبارت در $\lambda^* = \frac{t}{\nu^2}$ رخ می‌دهد. اگر داشته باشیم $\frac{t}{\nu^2} \in [0, \frac{1}{\alpha})$ و به تبع آن $0 \leq t \leq \frac{\nu^2}{\alpha}$ ، می‌توانیم λ^* را به جای λ قرار دهیم و در این صورت داریم:

$$\mathbb{P}[X - \mu \geq t] \leq e^{-\frac{t^2}{2\nu^2}}$$

اگر $t > \frac{\nu^2}{\alpha}$ باشد، از آن‌جا که تابع $g(\cdot, t)$ در بازه‌ی $[0, \lambda^*]$ به طور یکنوا کاهشی است، مقدار کمینه‌ی آن در مرز $\lambda = \frac{1}{\alpha}$ رخ می‌دهد، در نتیجه:

$$\mathbb{P}[X - \mu \geq t] \leq e^{-\frac{t}{\alpha} + \frac{1}{2\alpha} \frac{\nu^2}{\alpha}} \leq e^{-\frac{t}{\alpha} + \frac{1}{2\alpha} t} = e^{-\frac{t}{2\alpha}}$$

□

۳.آ بیان دیگری از فاصله‌ی KL

می‌دانیم فاصله‌ی KL به صورت زیر تعریف می‌شود:

$$D_{KL}(p||q) = \sum_{x \in \mathcal{X}} p(x) \ln\left(\frac{p(x)}{q(x)}\right) \quad (۲۵)$$

تعریف معادلی برای فاصله‌ی KL در قضیه‌ی زیر بیان شده است:

قضیه ۸.آ. فرض کنید $p(x)$ و $q(x)$ دو توزیع روی الفبای یکسان \mathcal{X} باشند. در این صورت می‌توان نوشت:

$$D_{KL}(p||q) = \sup_f \left\{ \mathbb{E}_p[f(X)] - \ln \mathbb{E}_q[e^{f(X)}] \right\} \quad (۲۶)$$

که سوپریمم روی همه‌ی توابع f گرفته می‌شود که در آن‌ها $\mathbb{E}_p[f(X)]$ و $\mathbb{E}_q[e^{f(X)}]$ خوش‌تعریف باشند.

۱.۳.آ اثبات قضیه‌ی (۸.آ)

اثبات. ابتدا فرض کنید تابع f را به این صورت تعریف کنیم:

$$f(x) = \ln\left(\frac{p(x)}{q(x)}\right).$$

در این حالت مشاهده می‌کنیم که:

$$\begin{aligned} \mathbb{E}_p[f(X)] - \ln \mathbb{E}_q[e^{f(X)}] &= \sum_{x \in \mathcal{X}} p(x) \ln\left(\frac{p(x)}{q(x)}\right) - \ln\left(\sum_{x \in \mathcal{X}} q(x) \frac{p(x)}{q(x)}\right) \\ &= \sum_{x \in \mathcal{X}} p(x) \ln\left(\frac{p(x)}{q(x)}\right) - \ln(1) \\ &= \sum_{x \in \mathcal{X}} p(x) \ln\left(\frac{p(x)}{q(x)}\right) = D_{KL}(p||q) \end{aligned}$$

در نتیجه:

$$\begin{aligned} \sup_f \left\{ \mathbb{E}_p[f(X)] - \ln \mathbb{E}_q[e^{f(X)}] \right\} &\geq \left[\mathbb{E}_p[f(X)] - \ln \mathbb{E}_q[e^{f(X)}] \right]_{f(x)=\ln(\frac{p(x)}{q(x)})} \\ &= D_{KL}(p||q) \end{aligned} \quad (۲۷)$$

از طرف دیگر، داریم:

$$\begin{aligned} \mathbb{E}_p[f(X)] - \ln \mathbb{E}_q[e^{f(X)}] &= \mathbb{E}_p[f(X)] - \mathbb{E}_p[\ln \mathbb{E}_q[e^{f(X)}]] \\ &= \mathbb{E}_p \left[f(X) - \ln \mathbb{E}_q[e^{f(X)}] \right] \\ &= \mathbb{E}_p \left[\ln \frac{e^{f(X)}}{\mathbb{E}_q[e^{f(X)}]} \right] \\ &= \sum_{x \in \mathcal{X}} p(x) \ln \frac{e^{f(x)}}{\sum_{x' \in \mathcal{X}} q(x') e^{f(x')}} \end{aligned}$$

اگر توزیع $q^{(f)}(x)$ را به صورت زیر تعریف کنیم:

$$q^{(f)}(x) = \frac{q(x)e^{f(x)}}{\sum_{x' \in \mathcal{X}} q(x')e^{f(x')}}.$$

می‌توانیم بنویسیم:

$$\begin{aligned} \mathbb{E}_p[f(X)] - \ln \mathbb{E}_q[e^{f(X)}] &= \sum_{x \in \mathcal{X}} p(x) \ln \frac{e^{f(x)}}{\sum_{x' \in \mathcal{X}} q(x')e^{f(x')}} \\ &= \sum_{x \in \mathcal{X}} p(x) \ln \frac{q^{(f)}(x)}{q(x)} \end{aligned}$$

و در نتیجه:

$$\begin{aligned} D_{KL}(p||q) - \left(\mathbb{E}_p[f(X)] - \ln \mathbb{E}_q[e^{f(X)}] \right) &= \sum_{x \in \mathcal{X}} p(x) \ln \left(\frac{p(x)}{q(x)} \right) - \sum_{x \in \mathcal{X}} p(x) \ln \left(\frac{q^{(f)}(x)}{q(x)} \right) \\ &= \sum_{x \in \mathcal{X}} p(x) \ln \left(\frac{p(x)}{q^{(f)}(x)} \right) \\ &= D_{KL}(p||q^{(f)}) \geq 0 \end{aligned}$$

در نتیجه:

$$D_{KL}(p||q) \geq \sup_f \left\{ \mathbb{E}_p[f(X)] - \ln \mathbb{E}_q[e^{f(X)}] \right\} \quad (28)$$

و از مقایسه‌ی (27) و (28) داریم:

$$D_{KL}(p||q) = \sup_f \left\{ \mathbb{E}_p[f(X)] - \ln \mathbb{E}_q[e^{f(X)}] \right\}$$

□

ب اثبات قضایای بیان‌شده

ب.۱ قضایای فصل ۲

ب.۱.۱ اثبات قضیه‌ی (۱.۲)

اثبات. با توجه به این‌که $\phi = (\phi_1, \dots, \phi_m)$ یک متغیر تصادفی برداری پیوسته است و T یک متغیر تصادفی اسکالر گسسته، اگر تحقق‌های ϕ را با $\varphi = (\varphi_1, \dots, \varphi_m)$ و تحقق‌های T را با i نشان دهیم، داریم:

$$I(T; \phi) = \sum_{i=1}^m \int_{\mathbb{R}^m} f_{\phi, T}(\varphi, i) \log \frac{f_{\phi, T}(\varphi, i)}{p_T(i) f_{\phi}(\varphi)} d\varphi \quad (29)$$

$$= \sum_{i=1}^m \int_{\mathbb{R}^m} f_{\phi|T}(\varphi|i) p_T(i) \log \frac{f_{\phi|T}(\varphi|i)}{f_{\phi}(\varphi)} d\varphi \quad (30)$$

$$= \sum_{i=1}^m p_T(i) \int_{\mathbb{R}^m} f_{\phi|T}(\varphi|i) \log \frac{f_{\phi|T}(\varphi|i)}{f_{\phi}(\varphi)} d\varphi \quad (31)$$

$$= \sum_{i=1}^m \mathbb{P}[T = i] D_{KL}(f_{\phi|T}||f_{\phi}) \quad (32)$$

و اگر تعریف کنیم:

$$\phi_{/i} = (\phi_1, \dots, \phi_{i-1}, \phi_{i+1}, \dots, \phi_m), \quad \varphi_{/i} = (\varphi_1, \dots, \varphi_{i-1}, \varphi_{i+1}, \dots, \varphi_m)$$

می توان نوشت:

$$I(T; \phi) = \sum_{i=1}^m p_T(i) \int_{\mathbb{R}} \int_{\mathbb{R}^{m-1}} f_{\phi_i|T}(\varphi_i|i) f_{\phi_{/i}|T, \phi_i}(\varphi_{/i}|i, \varphi_i) \log \frac{f_{\phi_i|T}(\varphi_i|i) f_{\phi_{/i}|T, \phi_i}(\varphi_{/i}|i, \varphi_i)}{f_{\phi_i}(\varphi_i) f_{\phi_{/i}}(\varphi_{/i})} d\varphi_{/i} d\varphi_i \quad (33)$$

$$= \sum_{i=1}^m p_T(i) \int_{\mathbb{R}} \int_{\mathbb{R}^{m-1}} f_{\phi_i|T}(\varphi_i|i) f_{\phi_{/i}|T, \phi_i}(\varphi_{/i}|i, \varphi_i) \log \frac{f_{\phi_i|T}(\varphi_i|i)}{f_{\phi_i}(\varphi_i)} d\varphi_{/i} d\varphi_i \\ + \sum_{i=1}^m p_T(i) \int_{\mathbb{R}} \int_{\mathbb{R}^{m-1}} f_{\phi_i|T}(\varphi_i|i) f_{\phi_{/i}|T, \phi_i}(\varphi_{/i}|i, \varphi_i) \log \frac{f_{\phi_{/i}|T, \phi_i}(\varphi_{/i}|i, \varphi_i)}{f_{\phi_{/i}}(\varphi_{/i})} d\varphi_{/i} d\varphi_i \quad (34)$$

$$\geq \sum_{i=1}^m p_T(i) \int_{\mathbb{R}} f_{\phi_i|T}(\varphi_i|i) \log \frac{f_{\phi_i|T}(\varphi_i|i)}{f_{\phi_i}(\varphi_i)} d\varphi_i \quad (35)$$

$$= \sum_{i=1}^m \mathbb{P}[T = i] D_{KL}(f_{\phi_i|T} || f_{\phi_i}) \quad (36)$$

حال اگر تعریف کنیم $g(\phi_i) = \lambda(\phi_i - \mu_i)$ و از قضیه ی (؟؟) با $p = f_{\phi_i|T}$ و $q = f_{\phi_i}$ استفاده کنیم، داریم:

$$D_{KL}(f_{\phi_i|T} || f_{\phi_i}) = \sup_f \left\{ \mathbb{E}_p[f(\phi_i)] - \ln \mathbb{E}_q[e^{f(\phi_i)}] \right\} \quad (37)$$

$$= \sup_f \left\{ \mathbb{E}[f(\phi_i)|T = i] - \ln \mathbb{E}[e^{f(\phi_i)}] \right\} \quad (38)$$

$$\geq \mathbb{E}[g(\phi_i)|T = i] - \ln \mathbb{E}[e^{g(\phi_i)}] \quad (39)$$

$$= \mathbb{E}[\lambda(\phi_i - \mu_i)|T = i] - \ln \mathbb{E}[e^{\lambda(\phi_i - \mu_i)}] \quad (40)$$

و از آن جا که طبق فرض قضیه، متغیر تصادفی $\phi_i - \mu_i$ یک متغیر تصادفی زیر-گوسی با پارامتر σ فرض شده است، با توجه به اینکه $\mathbb{E}[\phi_i - \mu_i] = \mathbb{E}[\phi_i] - \mu_i = \mu_i - \mu_i = 0$ از تعریف (؟؟) داریم:

$$\mathbb{E}[e^{\lambda(\phi_i - \mu_i)}] \leq e^{\frac{\lambda^2 \sigma^2}{2}} \quad (41)$$

و در نتیجه:

$$D_{KL}(f_{\phi_i|T} || f_{\phi_i}) \geq \mathbb{E}[\lambda(\phi_i - \mu_i)|T = i] - \ln \mathbb{E}[e^{\lambda(\phi_i - \mu_i)}] \quad (42)$$

$$\geq \lambda(\mathbb{E}[\phi_i|T = i] - \mu_i) - \frac{\lambda^2 \sigma^2}{2} \quad (43)$$

اگر تعریف کنیم $\Delta_i = \mathbb{E}[\phi_i|T = i] - \mu_i$ ، با توجه به اینکه (۴۲) برای همه ی مقادیر λ برقرار است، می توان نوشت:

$$D_{KL}(f_{\phi_i|T} || f_{\phi_i}) \geq \sup_{\lambda} \left\{ \lambda \Delta_i - \frac{\lambda^2 \sigma^2}{2} \right\} \quad (44)$$

$$= \left[\lambda \Delta_i - \frac{\lambda^2 \sigma^2}{2} \right]_{\lambda = \frac{\Delta_i}{\sigma^2}} = \frac{\Delta_i^2}{2\sigma^2} \quad (45)$$

و در نتیجه با ترکیب روابط (۲۶) و (۴۵) داریم:

$$2\sigma^2 I(T; \phi) \geq \sum_{i=1}^m \mathbb{P}[T = i] \Delta_i^2 = \mathbb{E}[\Delta_T^2] \quad (46)$$

و در نتیجه، با توجه به خواص امید ریاضی شرطی و نامساوی Jensen می توان نوشت:

$$|\mathbb{E}[\phi_T - \mu_T]| \leq \mathbb{E}[|\phi_T - \mu_T|] \quad (47)$$

$$= \mathbb{E}_T[\mathbb{E}[|\phi_i - \mu_i| \mid T = i]] \quad (48)$$

$$= \mathbb{E}[|\Delta_T|] \quad (49)$$

$$= \mathbb{E} \left[\sqrt{\Delta_T^2} \right] \quad (50)$$

$$\leq \sqrt{\mathbb{E}[\Delta_T^2]} \quad (51)$$

$$\leq \sigma \sqrt{2I(T; \phi)} \quad (52)$$

□

ب. ۲.۱. اثبات قضیه ی (۲.۲)

اثبات. مانند اثبات قضیه ی (۱.۲) می توان نوشت:

$$I(T; \phi) \geq \sum_{i=1}^m \mathbb{P}[T = i] D_{KL}(f_{\phi_i|T} \| f_{\phi_i}) \quad (53)$$

همچنین، از آن جا که $\phi_i - \mu_i$ یک متغیر تصادفی زیر-گوسی با پارامتر σ_i است، مشابه اثبات قضیه ی (۱.۲) اگر تعریف کنیم $\Delta_i = \mathbb{E}[\phi_i | T = i] - \mu_i$ ، داریم:

$$D_{KL}(f_{\phi_i|T} \| f_{\phi_i}) \geq \sup_{\lambda} \left\{ \lambda \Delta_i - \frac{\lambda^2 \sigma_i^2}{2} \right\} \quad (54)$$

$$= \left[\lambda \Delta_i - \frac{\lambda^2 \sigma_i^2}{2} \right]_{\lambda = \frac{\Delta_i}{\sigma_i^2}} = \frac{\Delta_i^2}{2\sigma_i^2} \quad (55)$$

در نتیجه داریم:

$$|\Delta_i| \leq \sigma_i \sqrt{2D_{KL}(f_{\phi_i|T} \| f_{\phi_i})} \quad (56)$$

حال می‌توان نوشت:

$$|\mathbb{E}[\phi_T - \mu_T]| \leq \mathbb{E}[|\phi_T - \mu_T|] \quad (57)$$

$$= \mathbb{E}_T[\mathbb{E}[|\phi_i - \mu_i| | T = i]] \quad (58)$$

$$= \mathbb{E}[|\Delta_T|] \quad (59)$$

$$= \sum_{i=1}^m |\Delta_i| \mathbb{P}[T = i] \quad (60)$$

$$\leq \sum_{i=1}^m \sigma_i \mathbb{P}[T = i] \sqrt{2D_{KL}(f_{\phi_i|T} || f_{\phi_i})} \quad (61)$$

$$\leq \sqrt{\sum_{i=1}^m \mathbb{P}[T = i] \sigma_i^2} \sqrt{2 \sum_{i=1}^m \mathbb{P}[T = i] D_{KL}(f_{\phi_i|T} || f_{\phi_i})} \quad (62)$$

$$\leq \sqrt{\mathbb{E}[\sigma_T^2]} \sqrt{2I(T; \phi)} \quad (63)$$

□

ب. ۳.۱. اثبات قضیه (۳.۲)

اثبات. مشابه اثبات قضیه (۱.۲) می‌توان نوشت:

$$I(T; \phi) \geq \sum_{i=1}^m \mathbb{P}[T = i] D_{KL}(f_{\phi_i|T} || f_{\phi_i}) \quad (64)$$

همچنین، از آن‌جا که $\phi_i - \mu_i$ یک متغیر تصادفی زیر-نمایی با پارامترهای (σ, b) است، مشابه اثبات قضیه (۱.۲) اگر تعریف کنیم $\Delta_i = \mathbb{E}[\phi_i | T = i] - \mu_i$ داریم:

$$D_{KL}(f_{\phi_i|T} || f_{\phi_i}) \geq \sup_{\lambda < \frac{1}{b}} \left\{ \lambda \Delta_i - \frac{\lambda^2 \sigma_i^2}{2} \right\} \quad (65)$$

$$\geq \left[\lambda \Delta_i - \frac{\lambda^2 \sigma_i^2}{2} \right]_{\lambda = \frac{1}{b}} = \frac{\Delta_i}{b} - \frac{\sigma^2}{2b^2} \quad (66)$$

در نتیجه:

$$I(T; \phi) \geq \sum_{i=1}^m \mathbb{P}[T = i] D_{KL}(f_{\phi_i|T} || f_{\phi_i}) \quad (67)$$

$$\geq \sum_{i=1}^m \mathbb{P}[T = i] \left(\frac{\Delta_i}{b} - \frac{\sigma^2}{2b^2} \right) \quad (68)$$

$$= \frac{\mathbb{E}_T[\mathbb{E}[\phi_i | T = i] - \mu_i]}{b} - \frac{\sigma^2}{2b^2} \quad (69)$$

$$= \frac{\mathbb{E}[\phi_T - \mu_T]}{b} - \frac{\sigma^2}{2b^2} \quad (70)$$

و در نتیجه:

$$\mathbb{E}[\phi_T - \mu_T] \leq bI(T; \phi) + \frac{\sigma^2}{2b} \quad (71)$$

در حالتی که $b < 1$ باشد، می‌توان $\lambda = \frac{1}{\sqrt{b}} < \frac{1}{b}$ قرار داد، در نتیجه:

$$D_{KL}(f_{\phi_i|T}||f_{\phi_i}) \geq \sup_{\lambda < \frac{1}{b}} \left\{ \lambda \Delta_i - \frac{\lambda^2 \sigma_i^2}{2} \right\} \quad (72)$$

$$\geq \left[\lambda \Delta_i - \frac{\lambda^2 \sigma_i^2}{2} \right]_{\lambda = \frac{1}{\sqrt{b}}} = \frac{\Delta_i}{\sqrt{b}} - \frac{\sigma^2}{2b} \quad (73)$$

و داریم:

$$I(T; \phi) \geq \sum_{i=1}^m \mathbb{P}[T = i] D_{KL}(f_{\phi_i|T}||f_{\phi_i}) \quad (74)$$

$$\geq \sum_{i=1}^m \mathbb{P}[T = i] \left(\frac{\Delta_i}{\sqrt{b}} - \frac{\sigma^2}{2b} \right) \quad (75)$$

$$= \frac{\mathbb{E}[\phi_T - \mu_T]}{\sqrt{b}} - \frac{\sigma^2}{2b} \quad (76)$$

و:

$$\mathbb{E}[\phi_T - \mu_T] \leq \sqrt{b} I(T; \phi) + \frac{\sigma^2}{2\sqrt{b}} \quad (77)$$

□

ب. ۴.۱. اثبات قضیه (۴.۲)

اثبات. تعریف می‌کنیم $U_i = \phi_i - \mu_i$. بنا به فرض می‌دانیم U_i زیر-گوسی با پارامتر σ است. همچنین تعریف می‌کنیم $Y_i = |U_i| - \gamma_i$ و $\gamma_i = \mathbb{E}[|U_i|] = \mathbb{E}[|\phi_i - \mu_i|]$ داریم:

$$\mathbb{P}[|Y_i| \geq s] = \mathbb{P}[Y_i \geq s] + \mathbb{P}[Y_i \leq -s] \quad (78)$$

$$= \mathbb{P}[|U_i| \geq s + \gamma_i] + \mathbb{P}[|U_i| \leq \gamma_i - s] \quad (79)$$

از قضیه (۲.۴) می‌دانیم:

$$\mathbb{P}[|U_i| \geq s + \gamma_i] \leq \sqrt{8e} \mathbb{P}[|Z| \geq s + \gamma_i] \leq \sqrt{8e} \mathbb{P}[|Z| \geq s] \quad (80)$$

که در آن، $Z \sim \mathcal{N}(0, 2\sigma^2)$ از طرف دیگر، داریم:

$$\mathbb{P}[|U_i| \leq \gamma_i - s] \leq \frac{\mathbb{P}[|Z| \geq s]}{\mathbb{P}[|Z| \geq \gamma_i]} \quad (81)$$

در نتیجه:

$$\mathbb{P}[|Y_i| \geq s] \leq \left(\sqrt{8e} + \frac{1}{\mathbb{P}[|Z| \geq \gamma_i]} \right) \mathbb{P}[|Z| \geq s] \quad (82)$$

در نتیجه، با توجه به قضیه (۵.۴)، Y_i یک متغیر تصادفی زیر-گوسی با پارامتر $\left(\sqrt{8e} + \frac{1}{\mathbb{P}[|Z| \geq \gamma_i]} \right) 2\sigma$ است. از آنجا که U_i زیر-گوسی با پارامتر σ است، واریانس آن کمتر از σ^2 خواهد بود و در نتیجه داریم:

$$\gamma_i = \mathbb{E}[|U_i|] = \mathbb{E}[\sqrt{U_i^2}] \leq \sqrt{\mathbb{E}[U_i^2]} \leq \sigma$$

که نتیجه می‌دهد:

$$\mathbb{P}[|Z| \geq \gamma_i] > \mathbb{P}[|Z| \geq \sigma] > 0.1 \quad (۸۳)$$

در نتیجه Y_i یک متغیر تصادفی زیر-گوسی با پارامتر $c\sigma$ است که $c < 36$ حال داریم:

$$\mathbb{E}[|\phi_T - \mu_T| - \gamma_T] = \mathbb{E}[Y_T] \leq c\sigma\sqrt{2I(T; \phi)} \quad (۸۴)$$

و از آنجا که برای هر i ، $\gamma_i \leq \sigma$ ، نتیجه می‌گیریم که $\gamma_T < \sigma$ و در نتیجه:

$$\mathbb{E}[|\phi_T - \mu_T|] \leq \sigma + c\sigma\sqrt{2I(T; \phi)} \leq \sigma + 36\sigma\sqrt{2I(T; \phi)} \quad (۸۵)$$

□

ب. ۵.۱. اثبات قضیه‌ی (۵.۲)

اثبات. تعریف می‌کنیم $Y_i = \phi_i - \mu_i$ و $\gamma_i = \mathbb{E}[(\phi_i - \mu_i)^2]$ از آنجا که Y_i زیر-گوسی با پارامتر σ است، واریانس آن کمتر از σ^2 است و در نتیجه $\gamma_i < \sigma^2$. از طرف دیگر، از آنجا که Y_i یک متغیر تصادفی زیر-گوسی با پارامتر σ است، با توجه به قضیه‌ی (۶.۲) داریم:

$$\mathbb{E}\left[e^{\frac{\lambda Y_i^2}{2\sigma^2}}\right] \leq \frac{1}{\sqrt{1-\lambda}} \quad \forall \lambda \in [0, 1) \quad (۸۶)$$

و می‌توان نوشت:

$$\mathbb{E}\left[e^{\frac{\lambda(Y_i^2 - \gamma_i)}{2\sigma^2}}\right] \leq \mathbb{E}\left[e^{\frac{\lambda Y_i^2}{2\sigma^2}}\right] \leq \frac{1}{\sqrt{1-\lambda}} \leq e^{10\lambda^2} \quad \forall \lambda \in [0, 0.1) \quad (۸۷)$$

اگر تعریف کنیم $t = \frac{\lambda}{\sigma^2}$ ، می‌توان نوشت:

$$\mathbb{E}\left[e^{t(Y_i^2 - \gamma_i)}\right] \leq e^{10\sigma^4 t^2} \quad \forall t \in [0, \frac{0.1}{\sigma^2}) \quad (۸۸)$$

و در نتیجه، $Y_i^2 - \gamma_i$ یک متغیر تصادفی زیر-نمایی با پارامترهای $(\sqrt{5}\sigma^2, 10\sigma^2)$ است. حال از قضیه‌ی (۳.۲) داریم:

$$\mathbb{E}[Y_T^2] \leq 10\sigma^2 I(T; \mathbf{Y}^2) + \frac{5\sigma^4}{20\sigma^2}$$

که $\mathbf{Y}^2 = (Y_1^2, \dots, Y_m^2)$ ، در نتیجه:

$$\mathbb{E}[(\phi_T - \mu_T)^2] \leq \gamma_T + 10\sigma^2 I(T; \mathbf{Y}^2) + \frac{5\sigma^4}{20\sigma^2} \quad (۸۹)$$

$$\leq \sigma^2 + 10\sigma^2 I(T; \mathbf{Y}^2) + \frac{5\sigma^4}{20\sigma^2} \quad (۹۰)$$

$$= \sigma^2(1.25 + 10I(T; \mathbf{Y}^2)) \quad (۹۱)$$

$$\leq \sigma^2(1.25 + 10I(T; \phi)) \quad (۹۲)$$

□

و نامساوی آخر از نامساوی پردازش داده‌ها نتیجه شده است.

ب. ۶.۱. اثبات قضیه (۶.۲)

اثبات. کران بالا از قضیه (۵.۲) قابل اثبات است، در اینجا به کران پایین می‌پردازیم. تعریف می‌کنیم $M = \max_i \phi_i$ و $M_{-i} = \max_{j \neq i} \phi_j$. در نتیجه $T = i$ خواهد شد، اگر و تنها اگر $M_{-i} \leq \phi_i$ باشد. تعریف می‌کنیم $I = \{i : \mathbb{E}[M_{-i}] \geq \mu_i + 1\}$ ، در نتیجه می‌توان آنتروپی را بسط داد:

$$H(T) = \sum_{i \notin I} \mathbb{P}[T = i] \log \left(\frac{1}{\mathbb{P}[T = i]} \right) + \sum_{i \in I} \mathbb{P}[T = i] \log \left(\frac{1}{\mathbb{P}[T = i]} \right) \quad (۹۳)$$

می‌دانیم که M و M_{-i} ، به ازای هر $i \in \{1, \dots, m\}$ ، زیر-گوسی با پارامتر ۱ هستند. در نتیجه می‌توان نوشت:

$$\mathbb{E}[e^{\lambda(M - \mathbb{E}[M])}] \leq e^{\frac{\lambda^2}{2}} \quad (۹۴)$$

$$\mathbb{E}[(M - \mathbb{E}[M])^2] \leq 1 \quad (۹۵)$$

$$\mathbb{P}[M \geq \mathbb{E}[M] + \lambda] \leq e^{\frac{-\lambda^2}{2}} \quad (۹۶)$$

همچنین اگر $X \sim (0, 1)$ ، برای هر $x > 0$ داریم:

$$\mathbb{P}[X > x] \geq \frac{1}{\sqrt{2\pi}} \left(\frac{x}{x^2 + 1} \right) e^{\frac{-x^2}{2}} \quad (۹۷)$$

حال سعی می‌کنیم در جمله اول مجموع (۹۳)، یک کران پایین برای $\mathbb{P}[T = i]$ بیابیم. چون در جمله اول هستیم، $\mathbb{E}[M_{-i}] < \mu_i + 1$ است.

$$\mathbb{P}[T = i] = \mathbb{P}[M_{-i} < \phi_i] \quad (۹۸)$$

$$\geq \mathbb{P}[M_{-i} < \mathbb{E}[M_{-i}] + \lambda] \cdot \mathbb{P}[\phi_i > \mathbb{E}[M_{-i}] + \lambda] \quad (۹۹)$$

$$\geq \mathbb{P}[M_{-i} < \mathbb{E}[M_{-i}] + \lambda] \cdot \mathbb{P}[\phi_i > \mu_i + 1 + \lambda] \quad (۱۰۰)$$

$$\geq \left(1 - e^{\frac{-\lambda^2}{2}}\right) \frac{1}{\sqrt{2\pi}} \left(\frac{1 + \lambda}{(1 + \lambda)^2 + 1} \right) e^{\frac{-(1+\lambda)^2}{2}} \quad (۱۰۱)$$

$$\triangleq p(\lambda) \quad (۱۰۲)$$

در نتیجه:

$$\sum_{i \notin I} \mathbb{P}[T = i] \log \left(\frac{1}{\mathbb{P}[T = i]} \right) \leq \mathbb{P}[T \notin I] \max_{i \notin I} \log \left(\frac{1}{\mathbb{P}[T = i]} \right) \quad (۱۰۳)$$

$$\leq \log \left(\frac{1}{p(1)} \right) \triangleq c_{-I} \quad (۱۰۴)$$

محاسبه مستقیم نشان می‌دهد که $c_{-I} < 5$. حال به سراغ جمله دوم می‌رویم، تعریف می‌کنیم $X = \phi_i - \mu_i \sim \mathcal{N}(0, 1)$ و $Y = M_{-i} - \mu_i$. حال می‌توان نوشت:

$$\mathbb{P}(T = i) = \int_{-\infty}^{\infty} \mathbb{P}(X > x) \mathbb{P}(Y = dx) \quad (۱۰۵)$$

$$\geq \int_1^{\infty} \mathbb{P}(X > x) \mathbb{P}(Y = dx) \quad (۱۰۶)$$

$$= \mathbb{P}(Y \geq 1) \int_1^{\infty} \mathbb{P}(X > x) \mathbb{P}(Y = dx | Y \geq 1) \quad (۱۰۷)$$

$$\geq \frac{\mathbb{P}(Y \geq 1)}{\sqrt{2\pi}} \int_1^{\infty} \left(\frac{x}{x^2 + 1} \right) e^{-x^2/2} \mathbb{P}(Y = dx | Y \geq 1). \quad (۱۰۸)$$

و با استفاده از نامساوی Jensen داریم:

$$\log \mathbb{P}(T = i) \geq \log(1/\sqrt{2\pi}) + \log(\mathbb{P}(Y \geq 1)) + \int_1^\infty \left(\log \left(\frac{x}{x^2 + 1} \right) - x^2/2 \right) \times \mathbb{P}(Y = dx | Y \geq 1), \quad (109)$$

که می‌تواند به صورت زیر نوشته شود:

$$\begin{aligned} \log \left(\frac{1}{\mathbb{P}(T = i)} \right) &\leq \log(\sqrt{2\pi}) + \log \left(\frac{1}{\mathbb{P}(Y \geq 1)} \right) \\ &\quad + \mathbb{E} \left[\log \left(\frac{Y^2 + 1}{Y} \right) | Y > 1 \right] + \frac{\mathbb{E}[Y^2 | Y > 1]}{2}. \end{aligned} \quad (110)$$

اگر $Y \geq 1$ باشد، داریم $\log((Y^2 + 1)/Y) \leq \log(1 + Y) \leq Y \leq Y^2$ در نتیجه:

$$\log \left(\frac{1}{\mathbb{P}(T = i)} \right) \leq \log(\sqrt{2\pi}) + \log \left(\frac{1}{\mathbb{P}(Y \geq 1)} \right) + 1.5\mathbb{E}[Y^2 | Y > 1]. \quad (111)$$

حال،

$$\mathbb{E}[Y^2 | Y > 1] \leq \frac{\mathbb{E}[Y^2]}{\mathbb{P}(Y > 1)} \quad (112)$$

$$= \frac{(\mathbb{E}[(Y - \mathbb{E}[Y])^2] + \mathbb{E}[Y]^2)}{\mathbb{P}(Y > 1)}. \quad (113)$$

از آنجا که $Y = M_{-i} - \mu_i$ ، واریانس Y کمتر از ۱ است. در نتیجه $\mathbb{P}(Y > 1) \geq 1 - 1/\sqrt{e}$ و

$$\log \left(\frac{1}{\mathbb{P}(T = i)} \right) \leq \log(\sqrt{2\pi}) + \log \left(\frac{1}{\mathbb{P}(Y \geq 1)} \right) + \frac{1.5(1 + \mathbb{E}[Y]^2)}{\mathbb{P}(Y \geq 1)} \quad (114)$$

$$< 5 + 4\mathbb{E}[Y]^2. \quad (115)$$

حال، اگر کران‌های به دست آمده را در کنار هم قرار دهیم، داریم:

$$H(T) = \sum_i \mathbb{P}(T = i) \log(1/\mathbb{P}(T = i)) \quad (116)$$

$$\leq c_{-I} + 5 + 4 \sum_{i \in I} \mathbb{P}(T = i) (\mathbb{E}[M_{-i}] - \mu_i)^2 \quad (117)$$

$$\leq c_{-I} + 5 + 4 \sum_{i \in I} \mathbb{P}(T = i) (\mathbb{E}[M] - \mu_i)^2 \quad (118)$$

$$\leq c_{-I} + 5 + 4 \|\mathbb{E}[M] - \mu_T\|^2 \quad (119)$$

که در آن، $\|X\| \equiv \sqrt{\mathbb{E}[X^2]}$ حال داریم:

$$\|\mathbb{E}[M] - \phi_T\| = \mathbb{E}[(\phi_T - \mathbb{E}[\phi_T])^2] \leq 1. \quad (120)$$

و در نتیجه:

$$\|\mathbb{E}[M] - \mu_T\| = \|\mathbb{E}[M] - \phi_T + \phi_T - \mu_T\| \leq 1 + \|\phi_T - \mu_T\|. \quad (121)$$

و می‌توانیم نتیجه بگیریم:

$$\|\mathbb{E}[M] - \mu_T\|^2 \leq (1 + \|\phi_T - \mu_T\|)^2 \leq 2 + 2\|\phi_T - \mu_T\|^2 \quad (122)$$

از اینجا نتیجه می‌شود که:

$$H(T) \leq c_{-I} + 5 + 8 + 8\|\phi_T - \mu_T\|^2 \quad (123)$$

یا:

$$\|\phi_T - \mu_T\|^2 \geq c_1 H(T) - c_2 \quad (124)$$

□

که در آن $c_1 = \frac{1}{8}$ و $c_2 = \frac{c_{-I}+13}{8} < 2.5$.

ب. ۲. اثبات قضایای فصل ۳

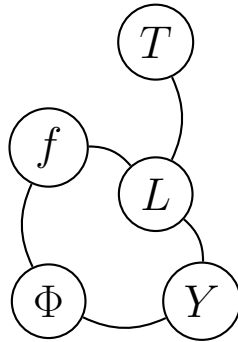
ب. ۱.۲. اثبات قضیه (۲.۰)

اثبات. بگیریم $\phi_i = \hat{L}(f_i)$ و $\mu_i = L(f_i)$. متغیر تصادفی T را نیز اندیس تصادفی‌ای در نظر بگیرید که $\hat{f}(\mathbf{x}) = f_T(\mathbf{x})$. توابع f_1, f_2, \dots, f_m هستند که طبقه‌بندی‌های مختلف را بر روی نمونه‌های $\mathbf{x} = (x_1, x_2, \dots, x_n)$ انجام می‌دهد (به تعریف تابع رشد مراجعه نمایید). با این تعریف‌ها، می‌توان از قضیه (۱.۲) مقاله برای بررسی مسئله طبقه‌بندی استفاده کرد. هدف، کران زدن بر روی $\mathbb{E}[\phi_T - \mu_T]$ است.

حال به بیان چند گزاره ساده در زمینه متغیرهای تصادفی زیرگوسی می‌پردازیم. اگر $X \sim \text{Bern}(p)$ متغیری تصادفی با پارامتر p باشد، آن‌گاه $X - p$ متغیری زیرگوسی با پارامتر $\frac{1}{4}$ است. به طور مشابه، اگر X_1, \dots, X_n متغیرهای تصادفی برنولی با پارامترهای p_1, \dots, p_n باشند، $\sum_{i=1}^n (X_i - p_i)$ یک متغیر تصادفی زیرگوسی با پارامتری کمتر از $1/4\sqrt{n}$ است. با استفاده از قضیه (۱.۲) داریم:

$$\mathbb{E}[\phi_T - \mu_T] \leq \sqrt{\frac{I(T; \phi)}{2n}}$$

شبکه‌ی مارکوفی زیر را می‌توان برای متغیرهای مسئله در نظر گرفت:



شکل ۳: شبکه‌ی مارکوفی متغیرها در مسئله طبقه‌بندی

با توجه به این شبکه‌ی مارکوفی و با استفاده از نامساوی پردازش اطلاعات، می‌توان نوشت:

$$I(T; \phi) \leq I(T; \mathbf{Y}) = I(\hat{f}(\mathbf{x}); \mathbf{Y})$$

از طرفی، با توجه به تعریف تابع رشد مجموعه‌ی توابع \mathcal{F} داریم:

$$I(T; \phi) \leq H(T) \leq \log(\Pi_{\mathcal{F}}(n)) \quad (125)$$

با استفاده از لم ساور-شلاح، این عبارت را می‌توان با بُعد VC مجموعه‌ی \mathcal{F} که محدود و برابر با d است، کران زد:

$$\Pi_{\mathcal{F}}(n) \leq \begin{cases} 2^n & n < d \\ (\frac{en}{d})^d & n \geq d \end{cases} \quad (۱۲۶)$$

با ترکیب (۱۲۵) و (۱۲۶) داریم: $I(\hat{f}(\mathbf{x}), \mathbf{Y}) \leq d \log_+(\frac{en}{d})$.

□

ب.۳ اثبات قضایای فصل ۴

ب.۱.۳ اثبات لم (۱.۴)

اثبات. از آن جایی که $H_k \perp\!\!\!\perp T_{k+1} | \phi$ ، با استفاده از نامساوی پردازش اطلاعات می‌توان نوشت:

$$I(T_{k+1}; \phi) \leq I(H_k; \phi).$$

در نتیجه داریم:

$$I(H_k; \phi) = \sum_{i=1}^k I((T_i, Y_{T_i}); \phi | H_{i-1}).$$

لذا $I((T_i, Y_{T_i}); \phi | H_{i-1})$ Let $\phi_{(-i)} = (\phi_j : j \neq i)$ با این وجود می‌توان نوشت:

$$\begin{aligned} I((T_i, Y_{T_i}); \phi | H_{i-1}) &= I(T_i; \phi | H_{i-1}) \\ &\quad + I(Y_{T_i}; \phi | H_{i-1}, T_i) \\ &= I(Y_{T_i}; \phi | H_{i-1}, T_i) \\ &= I(Y_{T_i}; \phi_{T_i} | H_{i-1}, T_i) \\ &\quad + I(Y_{T_i}; \phi_{(-T_i)} | H_{i-1}, T_i, \phi_{T_i}) \\ &= I(Y_{T_i}; \phi_{T_i} | H_{i-1}, T_i), \end{aligned}$$

عبارت آخر از استقلال Y_{T_i} و $\phi_{(-T_i)}$ به شرط ϕ_{T_i} ، Y_{T_i} به دست آمده است. به کمک این نامساوی داریم:

$$I(T_{k+1}; \phi) \leq I(H_k; \phi) = \sum_{i=1}^k I(Y_{T_i}; \phi_{T_i} | H_{i-1}, T_i)$$

□

ب.۲.۳ اثبات قضیه‌ی (۱.۴)

اثبات. فرض کنید $\phi_i \sim \mathcal{N}(\mu_i, \frac{\sigma^2}{n})$ و (ϕ_1, \dots, ϕ_k) متغیرهای مشترکاً گاوسی. اگر در مرحله‌ی j ام داشته باشیم $Y_{T_j} = \phi_{T_j} + W_j$ که $W_j \sim \mathcal{N}(0, \frac{\omega_j^2}{n})$ و متغیرهای W_j مستقل از ϕ باشند، آنگاه

$$\mathbb{E}[|Y_{T_{k+1}} - \mu_{T_{k+1}}|] \leq \frac{\sigma}{\sqrt{n}} + c_1 \left(\frac{\omega_{k+1}}{\sqrt{n}} + \sigma^2 \sqrt{\frac{\sum_{j=1}^k \omega_j^{-2}}{n}} \right).$$

با فرض $\omega_j = \sigma j^{\frac{1}{4}}$:

$$\begin{aligned}
\mathbb{E}[|Y_{T_{k+1}} - \mu_{T_{k+1}}|] &\leq \mathbb{E}[|Y_{T_{k+1}} - \phi_{T_{k+1}}|] + \mathbb{E}[|\phi_{T_{k+1}} - \mu_{T_{k+1}}|] \\
&\leq \sqrt{\frac{2\omega_{k+1}}{\pi n}} + \mathbb{E}[|\phi_{T_{k+1}} - \mu_{T_{k+1}}|] \\
&\leq \sqrt{\frac{2\omega_{k+1}}{\pi n}} + \frac{\sigma}{\sqrt{n}} + c \cdot \sigma \sqrt{\frac{2I(T_{k+1}; \phi)}{n}}
\end{aligned}$$

در نامساوی دوم، از امیدریاضی متغیر تصادفی نیم-گوسی استفاده شده است. با اعمال لم (۱.۴)، داریم:

$$I(T_{k+1}; \phi) \leq \sum_{i=1}^k I(Y_{T_i}; \phi_{T_i} | H_{i-1}, T_i)$$

از آنجایی که ϕ_i ها مشترکاً نرمال هستند، توزیع $\mathbb{P}(\phi_j | H_{i-1})$ نیز متغیری گاوسی و واریانسی کمتر از $\frac{\sigma^2}{n}$ است. همچنین به شرط H_{i+1} ، متغیر T_i مستقل از (ϕ_1, ϕ_2, \dots) و (Y_1, Y_2, \dots) می باشد. نتیجه این بحث این است که $\phi_{T_i} | H_{i-1}, T_i$ نیز دارای توزیع نرمال و واریانسی کمتر از σ^2/n است. با توجه به اطلاعات متقابل در متغیرهای گاوسی:

$$I(Y_{T_i}; \phi_{T_i} | H_{i-1}, T_i) \leq \frac{\sigma^2/n}{2\omega_i^2/n} = \frac{\sigma^2}{2\omega_i^2}$$

و بنابراین:

$$I(T_{k+1}; \phi) \leq \left(\frac{\sigma^2}{2}\right) \sum_{i=1}^k \omega_i^{-2}.$$

با جای گذاری این کران در عبارت مربوط به سویدگی، می توان نوشت:

$$\mathbb{E}[|Y_{T_{k+1}} - \mu_{T_{k+1}}|] \leq \sqrt{\frac{2\omega_{k+1}}{\pi n}} + \frac{\sigma}{\sqrt{n}} + c\sigma^2 \sqrt{\frac{\sum_{i=1}^k \omega_i^{-2}}{n}},$$

□

پ برخی دیگر از کاربردهای مسئله‌ی کاهش سویدگی از طریق کنترل اطلاعات استفاده شده

پ.۱ فیلتر کردن به کمک آماره‌های حاشیه‌ای

فرض کنید بعد از مشاهده‌ی دیتاست D ، T انتخاب شده باشد. دیتاست D مقادیر ϕ_1, \dots, ϕ_m را مشخص می کند ولی شامل اطلاعات دیگری نیز هست. داریم:

$$I(T; \phi) = H(T) - H(T|\phi) \quad (۱۲۷)$$

$$\leq H(T) - I(T; D|\phi) \quad (۱۲۸)$$

$$= (1 - \alpha)H(T) \quad (۱۲۹)$$

در این معادلات، $\alpha = I(T; D|\phi)/H(T)$ است. این پارامتر، بیانگر کسر عدم قطعیتی است که علاوه بر توابع ϕ در D وجود دارد. در بسیاری از مواقع، $I(T; \phi)$ خیلی کوچک تر از $H(T)$ است که خود کمتر از $\log(m)$ می باشد. یک مثال از این سناریو، انتخاب T بر اساس آماره‌های D می باشد، یک مثال از این مسئله، انتخاب ویژگی مبتنی بر واریانس است. فرض کنید n نمونه از

m ویژگی زیستی در اختیار داریم. مقدار ویژگی i ام در نمونه‌ی j ام را $X_{i,j}$ بنامید. دیتاست مورد استفاده $D = \{X_{i,j}\}$ است. تابع ϕ_i را میانگین نمونه‌ی ویژگی i ام بگیرد:

$$\phi_i = \frac{1}{n} \sum_{j=1}^n X_{i,j} \quad (۱۳۰)$$

ما علاقه‌مند به ویژگی‌هایی هستیم که میانگینی به طرز معنادار متفاوت با صفر داشته باشند. برخی از روش‌ها، در ابتدا یک مرحله‌ی فیلترینگ انجام می‌دهند و تنها ویژگی‌های با واریانس بزرگ را نگه می‌دارند و بقیه را حذف می‌کنند. این فیلترینگ با این استدلال انجام می‌شود که ویژگی‌های با واریانس کم، اطلاعاتی در بر ندارند و احتمالاً خطای سیستماتیک هستند. سوال طبیعی این است که آیا این فیلترینگ منجر به بایاس می‌شود؟ مسئله را به این صورت فرمول‌بندی می‌کنیم که

$$T = \operatorname{argmax}_i \sum_{j=1}^n (X_{i,j} - \phi_i)^2 \quad (۱۳۱)$$

یعنی تنها یک ویژگی انتخاب شده و این ویژگی دارای بزرگترین واریانس می‌باشد. تمام نتایج امکان تعمیم به حالتی که k ویژگی با واریانس بزرگ انتخاب شده‌اند نیز می‌باشد. قضیه‌ی (۱۰۲) بیان می‌کند که سویدگی $\mathbb{E}[\phi_T - \mu_T]$ کوچک است، اگر $I(T; \phi)$ کوچک باشد.

اگر میانگین نمونه و واریانس نمونه خیلی وابسته نباشند، این شرط برقرار می‌شود. به عنوان مثال می‌دانیم که اگر نمونه‌های گاوسی و i.i.d باشند، ϕ_1, \dots, ϕ_m از V_1, \dots, V_m مستقل و $I(T; \phi) = 0$ است، در نتیجه روش سویدگی ایجاد نمی‌کند. توجه داشته باشید که در این حالت $I(T; D)$ می‌تواند بسیار بزرگ باشد ولی قضیه‌ی (۱۰۲) تنها به $I(T; \phi)$ وابسته است. این بدین معناست که انتخاب T تنها وابسته به بخشی از اطلاعات دیتاست است که در توابع ϕ خود را نشان نداده است.

در حالت کلی‌ترین مسئله فرض کنید برای هر ویژگی i دو آماره‌ی حاشیه‌ای ϕ_i و ψ_i محاسبه شده‌اند. بنا بر نامساوی پردازش اطلاعات داریم $I(T; D) \leq I(\psi; \phi)$. در این جا $T = f(\psi_i)$ است و در تحلیل داده، علاقه‌مند به مقادیر ϕ_i هستیم و نتیجه‌ی تحلیل ما ϕ_T است. اگر ψ_i ها اطلاعات زیادی از ϕ_i ها در بر نداشته باشند، فیلترینگ داده بر اساس ψ_i ها منجر به سویدگی بزرگ نمی‌شود. رابطه‌ی متغیرها در این مسئله به فرم زنجیره‌ی مارکوف $T - \psi - \phi$ مدل می‌شود. با نامساوی پردازش اطلاعات داریم: $I(T; \phi) \leq I(\phi, \psi)$ که منتج به همان نتیجه‌ای که اگر ψ ها و ϕ ها اطلاعات متقابل کمی داشته باشند، سویدگی کم است می‌شود. این کران کران تیزی نیست و برای ارائه‌ی یک کران تیز، به نامساوی پردازش اطلاعات قوی روی می‌آوریم.

تعریف ۱۰. متغیرهای تصادفی X و Y در نامساوی قوی پردازش اطلاعات با پارامتر $\eta \in [0, 1]$ صدق می‌کنند، اگر برای هر متغیر تصادفی U با زنجیره‌ی $U - X - Y$ داشته باشیم:

$$I(U; Y) \leq \eta I(U; X)$$

η_{XY} را کوچکترین ضریبی بگیرد که نامساوی فوق برای تمام U ها برقرار باشد.

یکی از ویژگی‌های جالب ضریب η_{XY} این است که اگر $(X_1, Y_1), \dots, (X_n, Y_n)$ دنباله‌ای مستقل باشد، $\eta_{XY} = \max_i \eta_{X_i Y_i}$ است. همچنین اگر برای یک متغیر تصادفی Z زنجیره‌ی $X - Y - Z$ برقرار باشد، $\eta_{XZ} \leq \eta_{YZ}$ است.

مثال ۱۰. فرض کنید $D = (X_1, \dots, X_n)$ شامل n متغیر تصادفی i.i.d باشد و $\psi = (X_1, \dots, X_k)$ که $k \leq n$ است، یک زیرمجموعه‌ی n تایی از آن باشد. در این صورت $\eta_{\psi\phi} \leq \eta_{\psi D} \leq \frac{k}{n}$ است.

توجه داشته باشید که در مثال ما، $I_{\psi, \phi}$ تنها تابع توزیع مشترک ψ و ϕ است و ارتباطی با تابع $T = f(\psi)$ ندارد. این بدین معناست که از این نامساوی، می‌توان برای ارائه‌ی کرانی تیز بر روی سویدگی استفاده کرد.

قضیه ۱۰. اگر $\phi_i - \mu_i$ ها متغیرهای تصادفی σ -زیرگاوسی بوده و $T - \psi - \phi$ برقرار باشد، داریم:

$$\mathbb{E}[\phi_T - \mu_T] \leq \sigma \sqrt{2\eta_{\psi\phi} I(T; \psi)} \quad (۱۳۲)$$

□

اثبات. با توجه به قضیه‌ی (۱۰۲) و تعریف نامساوی پردازش اطلاعات قوی، بدیهی است.

پ.۲ استفاده کردن از اطلاعات و مسئله طبقه بندی

در این بخش به کاربرد قضیه مطرح شده در این مقاله در کران زدن خطای جواب مسئله طبقه بندی در یادگیری ماشین می پردازیم. فرض کنید n نمونه آموزشی در اختیار باشد. هر نمونه، شامل یک بردار $X_1, X_2, \dots, X_n \in \mathcal{X}$ است که هر یک به صورت i.i.d. از توزیع \mathcal{D} انتخاب شده اند که توزیع \mathcal{D} در اختیار نمی باشد. برای هر یک از این n نمونه، متغیر $Y_i \in \{0, 1\}$ داده شده است. فرض کنید که برجسب داده ی $X_i = x_i$ از توزیع $\mathbb{P}(Y_i | X_i = x_i)$ گرفته شده باشد.

تعریف پ.۲. یک طبقه بند، تابعی مانند f از \mathcal{X} به $\{0, 1\}$ است. خطای یک طبقه بند f بدین صورت تعریف می شود:

$$L(f) = \mathbb{E}_{X \sim \mathcal{D}} [\mathbf{1}(f(X) \neq Y)]$$

خطای طبقه بند بر روی داده های آموزشی نیز تعریف می شود:

$$\hat{L}(f) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(f(x_i) \neq y_i)$$

با این تعریف داریم $\mathbb{E}[\hat{L}(f)] = L(f)$.

الگوریتم یادگیری ماشین به دنبال تابعی مناسب است که رابطه ی بین ویژگی ها و برجسب ها را مدل کند. هدف این است که تابع \hat{f} از یک مجموعه ی داده شده ی توابع مثل \mathcal{F} است به نحوی پیدا شود که $L(\hat{f})$ کمینه شود. به دلیل اینکه الگوریتم به توزیع \mathcal{D} دسترسی ندارد، حل این مسئله بهینه سازی به صورت مستقیم ممکن نیست در نتیجه، به کمک روشی دیگر، تابع \hat{f} را انتخاب می کنیم. یکی از این روش ها، استفاده از

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \hat{L}(f)$$

است. به این الگوریتم Empirical Risk Minimization یا ERM می گویند.

در حالت کلی، با چنین روش هایی ممکن است دچار مشکلی به نام over-fitting شویم. این به این معنای این است که با وجود کوچک بودن $\hat{L}(\hat{f})$ مقدار $L(\hat{f})$ از $\min_{f \in \mathcal{F}} L(f)$ بسیار بزرگ تر باشد. تئوری یادگیری ماشین به دنبال بررسی تئوری فاصله ی $|L(\hat{f}) - \min_{f \in \mathcal{F}} L(f)|$ است. به طور خاص در تئوری یادگیری ماشین، برای الگوریتم ها و مجموعه های توابع مختلف، به دنبال تابع $m_{\mathcal{F}}$ در تعریف زیر هستند.

تعریف پ.۳. فرض کنید m نمونه ی i.i.d. $S = \{(X_1, Y_1), (X_2, Y_2), \dots, (X_m, Y_m)\}$ از توزیع \mathcal{D} در اختیار باشد. توزیع \mathcal{D} دلخواه است. الگوریتم A با دریافت این m نمونه، تابع $\hat{f} = A_S$ را انتخاب می کند. می گوئیم مجموعه ی توابع \mathcal{F} توسط الگوریتم A قابل یادگیری به مفهوم PAC هستند، اگر تابع چند جمله ای $m_{\mathcal{F}}(\cdot, \cdot)$ وجود داشته باشد که به ازای هر $(\epsilon, \delta) \in [0, 1]^2$ ، اگر $m \geq m_{\mathcal{F}}(\epsilon, \delta)$ باشد، با احتمال حداقل $1 - \delta$ داشته باشیم:

$$\mathbb{E}_{X \sim \mathcal{D}} [\mathbf{1}(A_S(X) \neq Y)] \leq \min_{f \in \mathcal{F}} \mathbb{E}_{X \sim \mathcal{D}} [\mathbf{1}(f(X) \neq Y)] + \epsilon \quad (۱۳۳)$$

توجه داشته باشید که $m_{\mathcal{F}}(\cdot, \cdot)$ تنها به مجموعه ی \mathcal{F} ربط دارد و شرط فوق باید به ازای هر توزیع \mathcal{D} با یک تابع مشخص و ثابت $m_{\mathcal{F}}(\cdot, \cdot)$ برقرار باشد.

مثال پ.۲. مثالی از مسئله طبقه بندی این است که $X_i \in \mathbb{R}^d$ و $\mathcal{F} = \{f_{\theta} : \theta \in \mathbb{R}^d\}$ باشد که در آن $f_{\theta}(x) = 1(x^T \theta \geq 0)$ باشد. الگوریتم A نیز θ ای را می یابد که در آن $\hat{L}(f_{\theta})$ کمینه شود. به این مسئله، طبقه بندی خطی می گوئیم. در این مثال، به وضوح می دانیم به افزایش d و با ثابت نگاه داشتن تعداد نمونه، ریسک over-fitting زیاد می شود.

در تئوری یادگیری ماشین، نتایج گوناگونی درباره ی کران فاصله ی

$$\left| \mathbb{E}_{X \sim \mathcal{D}} [\mathbf{1}(A_S(X) \neq Y)] - \min_{f \in \mathcal{F}} \mathbb{E}_{X \sim \mathcal{D}} [\mathbf{1}(f(X) \neq Y)] \right|$$

بر اساس معیارهای پیچیدگی کلاس \mathcal{F} وجود دارد. انتظار داریم با تعداد مشخص نمونه، با افزایش پیچیدگی کلاس \mathcal{F} ، فاصله ی مذکور بزرگ تر باشد. یکی از این معیارهای پیچیدگی، بُعد Vapnik-Chervonenkis یا VC Dimension است.

تعریف پ. ۴. بُد VC مجموعه‌ی توابع \mathcal{F} برابر بزرگترین کاردینالیتی مجموعه‌ی $S \subseteq \mathcal{X}$ است به نحوی که توابع \mathcal{F} قادر باشند تمام $2^{|S|}$ برچسب‌گذاری ممکن را بر روی نمونه‌های آن انجام دهند.

تعریف پ. ۵. تعداد برچسب‌گذاری‌هایی که توابع مجموعه‌ی \mathcal{F} می‌توانند بر روی یک مجموعه‌ی $S \subseteq \mathcal{X}$ انجام دهند را با $\Pi_{\mathcal{F}}(S)$ نمایش می‌دهیم. تابع رشد مجموعه‌ی توابع \mathcal{F} به این صورت تعریف می‌شود:

$$\Pi_{\mathcal{F}}(m) = \max_{S: |S|=m} \Pi_{\mathcal{F}}(S)$$

لم پ. ۱۰. ساور-شلاح: فرض کنید \mathcal{F} مجموعه‌ای از توابع با بعد VC محدود d باشد. در این صورت، به ازای هر $m \in \mathcal{N}$ داریم:

$$\Pi_{\mathcal{F}}(m) \leq \left(\frac{em}{d}\right)^d \quad (۱۳۴)$$

اثبات. این قضیه در درس تئوری یادگیری ماشین برای ما اثبات شده است و قضیه‌ی (۳۰۵) کتاب [۶] است. \square

قضیه پ. ۲. فرض کنید $\mathbf{x} = (x_1, \dots, x_n)$ ، $\mathbf{Y} = (Y_1, \dots, Y_n)$ ، $\hat{f}(\mathbf{x}) = (\hat{f}(x_1), \dots, \hat{f}(x_n))$ و $\log_+(z) = \max\{1, \log(z)\}$ باشند. در این صورت:

$$\mathbb{E}[L(\hat{f}) - \hat{L}(\hat{f})] \leq \sqrt{\frac{I(\hat{f}(\mathbf{x}); Y)}{2n}} \quad (۱۳۵)$$

به طور خاص، اگر \mathcal{F} دارای بعد VC محدود d باشد، آنگاه

$$I(\hat{f}(\mathbf{x}); Y) \leq d \log_+\left(\frac{ne}{d}\right) \quad (۱۳۶)$$

نتیجه پ. ۱۰. در تعریف یادگیری PAC دیدیم که برای این نوع یادگیری، نیاز است کرانی ثابت برای تمام توزیع‌های ممکن ورودی به دست آوریم. کران‌های معمول در تئوری یادگیری ماشین نیز از این جنس هستند که مستقل از توزیع ورودی، برای هر توزیع دلخواهی خطا را کران می‌زنند. از طرفی، کران (۱۳۵) وابسته به توزیع داده‌هاست. این نوع کران‌ها می‌توانند منجر به فهم بهتر این مسئله شوند که چه نوع توزیع‌هایی به نحوی بهتر قابل یادگیری هستند. همچنین برای برخی توزیع‌ها، کران به دست آمده، می‌تواند بسیار تیز از کران‌های نظریه‌ی PAC باشد.

پ. ۳. مشاهده‌ی داده‌ها و تعداد کلاس‌ها در خوشه‌بندی

در بسیاری از مواقع، در مسائل خوشه‌بندی داده‌ها با روش‌هایی مثل K-Means، شخص پردازشگر داده، با مشاهده‌ی نمودارهای داده، تعداد خوشه‌ها، K ، را مشخص می‌کند. در این تحلیل و با توجه به قضیه‌ی (۱۰۲)، سویدگی ناشی از این کار، اگر $I(T; \phi)$ کوچک باشد، کم است. بر اساس نامساوی پردازش اطلاعات داریم: $I(T; \phi) \leq I(K; \phi)$ است. در مواردی که ساختار مشخصی در داده وجود دارد که الگوریتم خوشه‌بندی قابل درک آن است، معمولاً K حول مقادیر خاصی متمرکز است و $I(K; \phi) \leq H(K) \approx 0$ است. در این مثال، تعداد خوشه‌ها برای تحلیلگر داده مفید است ولی مصداق «استفاده‌ی بد از اطلاعات» نیست.

پ. ۴. کنترل سویدگی از طریق کنترل FDR

کنترل False Discovery Rate یکی از مسائل مطرح در هنگام انجام تعداد زیادی آزمون‌های فرضیه همزمان است. فرض کنید دیتاست بزرگ $D \in \mathbb{R}^{n \times m}$ شامل n نمونه از m ژن باشد. فرض کنید n_1 نمونه‌ی اول از بافتی سرطانی بوده و $n - n_1$ نمونه‌ی بعدی نیز از بافت سالم آمده باشند. اصولاً دانشمندان علاقه‌مند به یافتن ژن‌هایی هستند که در توده‌ی سرطانی و توده‌های سالم تفاوت معناداری داشته باشند. برای این کار، برای هر ژن، یک آزمون فرضیه انجام می‌شود که رد شدن فرضیه‌ی صفر آن، نشان دهنده‌ی این است که بین توزیع این ژن در بافت سرطانی و بافت سالم تفاوت وجود دارد. در این حالت، بسیار نامحتمل است

تفاوت بین ژن در بافت‌های سرطانی و سالم به صورت تصادفی بوده باشد. هدف تنظیم آستانه‌ی رد و قبول فرضیه در این آزمون فرض است.

در این‌جا به بررسی یک حالت بسیار کلی از این مسئله می‌پردازیم. فرض کنید $D \in \mathbb{R}^{n \times m}$ یک ماتریس تصادفی بوده و بردار $\phi \in \mathbb{R}^m$ تابعی از آن باشد. به عنوان مثال، این تابع می‌تواند آماره‌های ستون‌های ماتریس D باشد، به عنوان مثال، تفاوت میانگین ستون i ام در بافت‌های سرطانی و بافت‌های سالم. امید ریاضی این بردار را μ بگیرد، یعنی $\mathbb{E}\phi(D) = \mu$. اندیس‌های $\{1, \dots, m\}$ به دو بخش H_0 و H_1 تقسیم شده‌اند. یک فرآیند انتخاب، تابعی مثل $\psi : \mathbb{R}^{n \times m} \rightarrow \{0, 1\}^m$ است که $\psi(D)_i = 1$ نشان می‌دهد که ویژگی (یا ژن) i ام انتخاب شده است. مجموعه‌ی $S_1 \subseteq \{1, \dots, m\}$ را مجموعه‌ی ژن‌های انتخاب شده و مجموعه‌ی S_2 را مجموعه‌ی ژن‌های انتخاب نشده در نظر بگیرید. در این فرمول‌بندی، ستون‌های با اندیس عضو H_0 ، ستون‌هایی هستند که فرض صفر برای آن‌ها برقرار است و تفاوت ژن مربوطه در بافت سالم و سرطانی، معنادار نیست. تعریف‌های زیر، تعریف‌هایی مشابه تعریف‌های معمول خطای نوع اول و خطای نوع دوم هستند:

$$\hat{\alpha} = \frac{\#(H_0 \cap S_1)}{\#H_0} \quad \hat{\beta} = \frac{\#(H_1 \cap S_0)}{\#H_1}$$

ما علاقه‌مند هستیم که متغیرهایی به فرم زیر را بررسی کنیم:

$$\frac{1}{\#S_1} \sum_{i \in S_1} (\phi_i - \mu_i) \quad (137)$$

$$\frac{1}{\#S_1} \sum_{i \in S_1} |\phi_i - \mu_i| \quad (138)$$

$$\frac{1}{\#S_1} \sum_{i \in S_1} (\phi_i - \mu_i)^2 \quad (139)$$

این متغیرها، به معنای خطا (یا سویدگی) میانگین در متغیرهای انتخاب شده هستند. این کمیت‌ها را می‌توان به ترتیب به صورت $\mathbb{E}[\phi_T - \mu_T]$ ، $\mathbb{E}[|\phi_T - \mu_T|]$ ، $\mathbb{E}[(\phi_T - \mu_T)^2]$ نوشته‌شوند که در آن، متغیر تصادفی T به شرط D دارای توزیع یکنواخت بر روی متغیرهای انتخاب شده در S_1 دارد.

حال، FDR را به صورت $\text{FDR} = \mathbf{P}(T \in H_0)$ تعریف می‌کنیم. این متغیر، نشان می‌دهد که برای چه کسری از متغیرهای انتخاب شده، عضو H_0 بوده‌اند و اشتباه رخ داده است. قضیه‌ی فوق قضیه‌ای بسیار مهم است که نشان می‌دهد که در صورتی که در فاز آزمون فرضیه و انتخاب ویژگی خوب عمل کنیم و FDR کوچکی داشته باشیم، اطلاعات متقابل $I(T; \phi)$ نیز کوچک است و در نتیجه در فاز تخمین نیز دچار سویدگی بزرگی نخواهیم بود.

قضیه ۳.۰. در مسئله‌ی فوق، داریم:

$$I(T; \phi) \leq h(\text{FDR}) + (1 - \text{FDR}) \log\left(\frac{1}{1 - \beta}\right) + \text{FDR} \log\left(\frac{1}{\alpha}\right) + \xi \quad (140)$$

که در آن، $h(p)$ تابع آنتروپی باینری است. $\alpha = \mathbb{E}\hat{\alpha}$ و $\beta = \mathbb{E}\hat{\beta}$ خطاهای نوع اول و دوم هستند و داریم:

$$\xi = \mathbb{E}\left[\log_+\left(\frac{1 - \beta}{1 - \hat{\beta}}\right)\right] + \mathbb{E}\left[\log_+\left(\frac{\alpha}{\hat{\alpha}}\right)\right]$$

اثبات. تعریف می‌کنیم $\mathcal{X} = \mathbf{1}(T \in H_0)$. از آنجایی که \mathcal{X} تابعی یقینی از T است، می‌توان نوشت:

$$I(T; \phi) = I(T; \mathcal{X}, \phi) = I(T; \mathcal{X}) + I(T; \phi | \mathcal{X}) \leq H(\mathcal{X}) + I(T; \phi | \mathcal{X}) \quad (141)$$

با استفاده از توزیع \mathcal{X} و استفاده از نامساوی پردازش اطلاعات و قاعده‌ی زنجیر، داریم:

$$I(T; \phi | \mathcal{X}) \leq H(\mathcal{X}) + \mathbf{P}(\mathcal{X} = 0)I(T; \phi | \mathcal{X} = 0) + \mathbf{P}(\mathcal{X} = 1)I(T; \phi | \mathcal{X} = 1) \quad (142)$$

$$= h(\mathbf{P}(T \in H_1)) + \mathbf{P}(T \in H_1)I(T; \phi | T \in H_1) \quad (143)$$

$$+ \mathbf{P}(T \in H_0)I(T; \phi | T \in H_0) \quad (144)$$

$$\leq h(\mathbf{P}(T \in H_1)) + \mathbf{P}(T \in H_1)I(T; D | T \in H_1) \quad (145)$$

$$+ \mathbf{P}(T \in H_0)I(T; D | T \in H_0) \quad (146)$$

می‌توان به سادگی $I(T; D \mid T \in \mathcal{H}_1)$ را به فرم زیر کران زد:

$$I(T; D \mid T \in \mathcal{H}_1) = H(T \mid T \in \mathcal{H}_1) - H(T \mid T \in \mathcal{H}_1, D) \quad (۱۴۷)$$

$$\leq \log(\#\mathcal{H}_1) - \mathbb{E}[\log(\#(S_1 \cap \mathcal{H}_1)) \mid T \in \mathcal{H}_1] \quad (۱۴۸)$$

$$= \log(\#\mathcal{H}_1) - \mathbb{E}[\log((1 - \hat{\beta}) \cdot (\#\mathcal{H}_1)) \mid T \in \mathcal{H}_1] \quad (۱۴۹)$$

$$= -\mathbb{E}\left[\log(1 - \hat{\beta}) \mid T \in \mathcal{H}_1\right] \quad (۱۵۰)$$

$$= -\log(1 - \beta) + \mathbb{E}\left[\log\left(\frac{1 - \beta}{1 - \hat{\beta}}\right) \mid T \in \mathcal{H}_1\right] \quad (۱۵۱)$$

با استفاده از این کران، داریم:

$$\mathbb{P}(T \in \mathcal{H}_1)I(T; D \mid T \in \mathcal{H}_1) \leq -\mathbb{P}(T \in \mathcal{H}_1) \log(1 - \beta) \quad (۱۵۲)$$

$$+ \mathbb{E}\left[\log\left(\frac{1 - \beta}{1 - \hat{\beta}}\right) \mathbf{1}_{\{T \in \mathcal{H}_1\}}\right] \quad (۱۵۳)$$

$$\leq -\mathbb{P}(T \in \mathcal{H}_1) \log(1 - \beta) \quad (۱۵۴)$$

$$+ \mathbb{E}\left[\log_+\left(\frac{1 - \beta}{1 - \hat{\beta}}\right)\right] \quad (۱۵۵)$$

با تکرار محاسبات فوق، به نتیجه‌ی زیر می‌رسیم:

$$\mathbb{P}(T \in \mathcal{H}_0)I(T; X \mid T \in \mathcal{H}_0) \leq -\mathbb{P}(T \in \mathcal{H}_0) \log(\alpha) + \mathbb{E}\left[\log_+\left(\frac{\alpha}{\hat{\alpha}}\right)\right] \quad (۱۵۶)$$

□

با جای‌گذاری این نتایج در معادله‌ی (۱۴۱)، قضیه اثبات می‌شود.