

کتاب راهنمای یادگیری ماشین

پردازش رادیو و ژاک و مارتا وایت

فهرست

۶	مرجع علامت گذاری
۱۰	پیشگفتار: شروع با یک مثال رگرسیون خطی
۱۲	مقدمه‌ای بر مدل سازی احتمالی
۱۴	۱-۱ نظریه احتمال و متغیرهای تصادفی
۱۶	۱-۲ تعریف توزیع
۱۷	۱-۲-۱ توابع جرم احتمال
۱۹	۱-۲-۲ توابع چگالی احتمال
۲۴	۱-۳ متغیرهای تصادفی چند متغیره
۲۶	۱-۳-۱ توزیع‌های مشروط
۲۷	۱-۳-۲ متغیرهای تصادفی مستقل
۲۹	۱-۴ امیدهای ریاضی و گشتاور
۳۳	۱-۵ چند متغیره PMF و PDF
۳۵	مقدمه‌ای بر بهینه سازی
۳۵	۲-۱ مسئله بهینه سازی اساسی و نقاط ثابت
۳۷	۲-۲ گرادیان کاهشی
۳۹	۲-۳ انتخاب اندازه گام
۴۰	۲-۴ خواص بهینه سازی
۴۲	اصول اولیه تخمین پارامتر
۴۲	۳-۱ نقشه و برآورد ماکسیم احتمال
۴۸	۳-۲ ماکسیم احتمال برای توزیع‌های شرطی
۴۹	۳-۳ [پیشرفته] رابطه بین به ماکسیم رساندن احتمال و واگرایی $Kullback - Leibler$
۵۱	مقدمه‌ای بر مسائل پیش‌بینی
۵۲	۴-۱ مسائل یادگیری تحت نظارت
۵۲	۴-۱-۱ رگرسیون و طبقه بندی
۵۴	۴-۱-۲ تصمیم گیری در مورد نحوه فرمول بندی کردن مسئله
۵۵	۴-۲ یادگیری بدون نظارت و یادگیری نیمه نظارت
۵۵	۴-۳ طبقه بندی بهینه و مدل های رگرسیون
۵۶	۴-۳-۱ نمونه هایی از هزینه ها

- ۵۷ _____ ۲-۳-۴ استخراج پیش‌بینی‌کننده‌های بهینه
- ۵۹ _____ ۳-۳-۴ خطای قابل کاهش و کاهش ناپذیر
- ۶۰ _____ ۴-۴ [پیشرفته] مدل‌های بهینه بیز
- ۶۲ _____ رگرسیون خطی
- ۶۲ _____ ۱-۵ فرمول ماکسیمم احتمال
- ۶۵ _____ ۲-۵ رگرسیون مینیمم مربعات معمولی (OLS).
- ۶۷ _____ ۱-۲-۵ تابع خطای وزنی
- ۶۸ _____ ۲-۲-۵ پیش‌بینی چندین خروجی به طور همزمان
- ۶۸ _____ ۳-۵ رگرسیون خطی برای مسائل غیر خطی
- ۶۹ _____ ۱-۳-۵ برازش منحنی چند جمله‌ای
- ۷۱ _____ ۴-۵ ثبات و مبادله بایاس واریانس
- ۷۱ _____ ۱-۴-۵ حساسیت راه حل OLS
- ۷۴ _____ ۲-۴-۵ منظم سازی
- ۷۵ _____ ۳-۴-۵ انتظار و واریانس برای راه حل منظم
- ۷۸ _____ ۵-۵ مبادله بایاس واریانس
- ۸۰ _____ اصول بهینه‌سازی پیشرفته‌تر
- ۸۰ _____ ۱-۶ گرادینت کاهشی در توابع چند متغیره
- ۸۲ _____ ۲-۶ خواص هسین
- ۸۳ _____ ۳-۶ مدیریت مجموعه داده‌های بزرگ
- ۸۵ _____ ۴-۶ بهینه‌سازی غیر هموار اما همچنان مستمر
- ۸۶ _____ ۵-۶ روش‌های بیشتر برای انتخاب اندازه گام‌ها
- ۸۷ _____ مدل‌های خطی تعمیم یافته
- ۸۸ _____ ۱-۷ انتقال نمایی و توزیع پواسون
- ۹۰ _____ ۲-۷ توزیع‌های خانوادگی نمایی
- ۹۱ _____ ۳-۷ فرمول‌بندی کردن مدل‌های خطی تعمیم یافته
- ۹۴ _____ طبقه‌بندی‌های خطی
- ۹۵ _____ ۱-۸ رگرسیون لجستیک
- ۹۶ _____ ۱-۱-۸ پیش‌بینی برچسب‌های کلاس
- ۹۷ _____ ۲-۱-۸ برآورد حداکثر احتمال برای رگرسیون لجستیک

- ۸-۱-۳ مسائل مربوط به به حداقل رساندن فاصله اقلیدسی _____ ۹۹
- ۸-۲ طبقه‌بندی کننده بیز ساده _____ ۱۰۱
- ۸-۲-۱ ویژگی‌های باینری و طبقه‌بندی خطی _____ ۱۰۲
- ۸-۲-۲ بیز ساده پیوسته _____ ۱۰۴
- ۸-۳ رگرسیون لجستیک چند جمله ای _____ ۱۰۵
- نمایش‌هایی برای یادگیری ماشینی _____ ۱۰۷
- ۹-۱ شبکه‌های تابع پایه شعاعی و نمایش هسته _____ ۱۰۷
- ۹-۲ بازنمایی‌های یادگیری _____ ۱۰۹
- ۹-۲-۱ شبکه‌های عصبی _____ ۱۰۹
- ۹-۲-۲ یادگیری بدون نظارت و فاکتورسازی ماتریس _____ ۱۱۵
- ارزیابی الگوریتم‌های یادگیری _____ ۱۱۸
- ۱۰-۱ مقدمه‌ای کوتاه بر مرزهای تعمیم _____ ۱۱۹
- ۱۰-۱-۱ نابرابری‌های تمرکز _____ ۱۲۰
- ۱۰-۱-۲ پیچیدگی یک کلاس تابع _____ ۱۲۱
- ۱۰-۱-۳ مرزهای تعمیم _____ ۱۲۲
- ۱۰-۲ مقایسه الگوریتم‌های یادگیری _____ ۱۲۲
- ۱۰-۳ به دست آوردن نمونه‌های خطا _____ ۱۲۴
- ۱۰-۴ معیارهای عملکرد برای مدل‌های طبقه بندی _____ ۱۲۵
- مطالب اضافی برای نظریه احتمال _____ ۱۲۸
- A.1 بدیهیات احتمال _____ ۱۲۸
- A.2 چند pmf مفید دیگر _____ ۱۲۹
- A.3 چند pdf مفید دیگر _____ ۱۳۰
- A.4 متغیرهای تصادفی _____ ۱۳۱
- A.4.1 تعریف رسمی متغیر تصادفی _____ ۱۳۳
- A.4.2 مثالی از استقلال مشروط _____ ۱۳۶
- A.4.3 اطلاعات اضافی برای انتظارات و لحظات _____ ۱۳۶
- A.5 مخلوط‌های توزیع _____ ۱۳۹
- A.6 نمایش گرافیکی توزیع‌های احتمال _____ ۱۴۱
- پیوست B _____ ۱۴۶

۱۴۶	B.1 قوانین اساسی برای گرادیان
۱۴۷	پیوست C
۱۴۷	C.1 دیدگاه جبری
۱۴۷	C.1.1 چهار زیرفضای اساسی
۱۴۹	C.1.2 به حداقل رساندن $\ Ax - b\ _2^2$
۱۵۲	پیوست D
۱۵۵	پیوست E

مرجع علامت گذاری

مجموعه نمادگذاری‌ها

\mathcal{X} مجموعه‌ای عمومی از مقادیر. به عنوان مثال، $\mathcal{X} = \{0, 1\}$ مجموعه‌ای است که فقط شامل 0 و 1 است، $\mathcal{X} = [0, 1]$ فاصله بین 0 تا 1 و $\mathcal{X} = \mathbb{R}$ مجموعه اعداد حقیقی است. بسته به مکان استفاده آن، نمادهایی مانند A ، B ، Ω و موارد دیگر نیز به عنوان مجموعه استفاده خواهند شد.

$P(\mathcal{X})$ مجموعه توان \mathcal{X} ، مجموعه‌ای شامل تمام زیر مجموعه‌های ممکن \mathcal{X} .

$[a, b]$ بازه بسته با شرط $a < b$ ، شامل a و b .

(a, b) بازه باز با شرط $a < b$ ، بدون a و b در مجموعه.

$(a, b]$ بازه نیم‌باز با شرط $a < b$ ، شامل b اما نه a .

$[a, b)$ بازه نیم‌باز با شرط $a < b$ ، شامل a اما نه b .

نماد برداری و ماتریس

x متغیرهای کوچک کمرنگ معمولاً اسکالر هستند. با این حال، هنگامی که $x \in \mathcal{X}$ ، جایی که \mathcal{X} مشخص نشده است، x ممکن است بردار، یک شی ساختار یافته مانند گراف و غیره را نشان دهد.

\mathbf{x} متغیرهای با حروف کوچک پررنگ بردار هستند. به طور پیش فرض، بردارها بردارهای ستونی هستند.

\mathbf{X} متغیرهای با حروف بزرگ پررنگ ماتریس هستند. به نظر می‌رسد یک متغیر تصادفی چند متغیره، \mathbf{X} ، اما متغیر تصادفی مورب است. اغلب از زمینه مشخص می‌شود که چه زمانی این یک متغیر تصادفی چند متغیره و چه زمانی یک ماتریس است.

\mathbf{X}^T جابجایی ماتریس. برای دو ماتریس \mathbf{A} و \mathbf{B} ، آن را برقرار می‌کند

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$$

یک ماتریس $n \times d$ متشکل از n بردار هر یک از ابعاد d را می‌توان به صورت بیان کرد

$$\mathbf{X} = [\mathbf{x}_1 \mathbf{x}_2 \dots \mathbf{x}_n]^T$$

X_i : ردیف i -ام ماتریس. یک بردار سطری.

$X_{:j}$: ستون j ماتریس. یک بردار ستونی

تاپل‌ها، بردارها و دنباله‌ها

(x_1, x_2, \dots, x_d) یک تاپل؛ به عنوان مثال، یک لیست مرتب از عناصر d . وقتی $(x_1, x_2, \dots, x_d) \in \mathbb{R}^d$ ، تاپل به عنوان بردار ستون $\mathbf{x} = [x_1 x_2 \dots x_d]^T$ تلقی می‌شود.

a_1, \dots, a_m دنباله‌ای از m آیت‌ها. متغیرهای شاخص روی این دنباله‌ها معمولاً متغیرهای i, j یا k هستند. به عنوان مثال، $\sum_{i=1}^m a_i$ یا اگر هر a_i بردار بعد d باشد، شاخص دوگانه $\sum_{i=1}^m \sum_{j=1}^d a_{ij}$ است.

نماد توابع

$f: \mathcal{X} \rightarrow \mathcal{Y}$ تابع در دامنه \mathcal{X} به هم دامنه \mathcal{Y} تعریف می‌شود و مقادیر $x \in \mathcal{X}$ را می‌گیرد و آنها را به $f(x) \in \mathcal{Y}$ ارسال می‌کند.

$\frac{dx}{dx}(x)$ مشتق تابع در $x \in \mathcal{X}$ که در آن $f: \mathcal{X} \rightarrow \mathbb{R}$ برای $\mathcal{X} \subset \mathbb{R}$

$\nabla f(x)$ گرادیان یک تابع در $x \in \mathcal{X}$ که در آن $f: \mathcal{X} \rightarrow \mathbb{R}$ برای $\mathcal{X} \subset \mathbb{R}^d$ آن را نگه می‌دارد

$$\nabla f(x) = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_d} \right).$$

$H_{f(x)}$ ماتریس هسین¹ یک تابع در $x \in \mathcal{X}$ که در آن $f: \mathcal{X} \rightarrow \mathbb{R}$ برای $\mathcal{X} \subset \mathbb{R}^d$ آن را نگه می‌دارد

$$H_{f(x)} = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_d} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & & \\ \vdots & & \ddots & \\ \frac{\partial^2 f}{\partial x_d \partial x_1} & & \dots & \frac{\partial^2 f}{\partial x_d^2} \end{bmatrix}$$

$\ell: \mathbb{R}^d \rightarrow \mathbb{R}$ یک تابع هزینه که نشان دهنده خطای پیش‌بینی رخ داده توسط وزن‌های داده شده، $\ell(\mathbf{w})$. اگر مشترک باشد، ℓ_i معمولاً هزینه را در نمونه i -ام با $\frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w})$ برای n نمونه نشان می‌دهد،

¹ Hessian matrix

$c: \mathbb{R}^d \rightarrow \mathbb{R}$ یک تابع هدف عمومی که می‌خواهیم آن را برای متغیر آموخته شده w کمینه کنیم. این می‌تواند، برای مثال، یک هزینه به علاوه یک تنظیم کننده باشد.

متغیرهای تصادفی و احتمالات

- X یک متغیر تصادفی تک متغیره با حروف بزرگ نوشته می‌شود.
- \mathcal{X} فضای مقادیر برای متغیر تصادفی.
- x متغیر کوچک یک نمونه یا نتیجه است، $x \in \mathcal{X}$.
- \mathbf{X} یک متغیر تصادفی چند متغیره با حروف بزرگ پررنگ نوشته می‌شود.
- \mathbf{x} متغیر پررنگ با حروف کوچک یک نمونه چند متغیره است. در موارد خاص، زمانی که مقدار متغیر به عنوان یک بردار در نظر گرفته می‌شود، از \mathbf{x} استفاده می‌کنیم.
- $\mathcal{N}(\mu, \sigma^2)$ یک توزیع گاوسی تک متغیره، با پارامترهای μ, σ_2 .
- \sim نشان می‌دهد که یک متغیر به عنوان یک توزیع بیان می‌شود. مثال، $X \sim \mathcal{N}(\mu, \sigma_2)$.

پارامترها و برآورد

- \mathcal{D} یک مجموعه داده، معمولاً از n عنصر ورودی‌های چند متغیره $\mathbf{X} \in \mathbb{R}^{n \times d}$ و خروجی‌های تک متغیره $y \in \mathbb{R}^n$ یا خروجی‌های چند متغیره $Y \in \mathbb{R}^{n \times m}$ تشکیل شده است. مجموعه داده همچنین به عنوان مجموعه‌ای از تاپل‌های نمایه شده نامیده می‌شود. به عنوان مثال، $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$.
- M یک مدل عمومی را برای بحث در مورد تخمین پارامتر کلی نشان می‌دهد. به عنوان مثال، $M = \theta$ برای برخی از پارامترهای θ ، مانند میانگین توزیع گاوسی.
- ω پارامترهای حقیقی برای مدل‌های رگرسیون خطی (تعمیم یافته) و طبقه بندی، معمولاً با $\omega \in \mathbb{R}^d$.
- w پارامترهای تقریبی برای مدل‌های رگرسیون خطی (تعمیم یافته) و طبقه بندی، معمولاً با $w \in \mathbb{R}^d$. هنگامی که w را به عنوان راه حل ماکسیمم احتمال در برخی از داده‌ها مورد بحث قرار می‌دهیم، $w_{ML}(\mathcal{D})$ را می‌نویسیم تا نشان دهیم که تغییرپذیری از \mathcal{D} ناشی می‌شود.

$$\begin{aligned} \max_{a \in B} f(a) & \text{ ماکسیمم مقدار یک تابع } f \text{ در بین مقادیر } a \text{ در مجموعه } B. \\ \operatorname{argmax}_{a \in B} f(a) & \text{ مورد } a \text{ در مجموعه } B \text{ که ماکسیمم مقدار } f(a) \text{ را تولید می‌کند.} \end{aligned}$$

نرم‌ها

$\|x\|$ یک نرم روی X .

$\|x\|_2$ نرم ℓ_2 در یک بردار، $\|x\|_2 = \sqrt{\sum_{i=1}^d x_i^2}$. این نرم فاصله اقلیدسی را از مبدأ دستگاه مختصات به X می‌دهد. یعنی طول بردار x است.

$\|x\|_2^2 = \sum_{i=1}^d x_i^2$ نرم مجذور^۱ ℓ_2 در یک بردار،

$\|x\|_p$ نرم کلی^۲ ℓ_p در یک بردار، $\|x\|_p = \left(\sum_{i=1}^d |x_i|^p\right)^{1/p}$

$\|X\|_F$ نرم فروبنیوس^۳ یک ماتریس n در d است. یعنی

$$\|X\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^d X_{ij}^2} = \sqrt{\sum_{i=1}^n \|X_{i:}\|_2^2} = \sqrt{\sum_{j=1}^d \|X_{:j}\|_2^2}$$

فرمول‌ها و قوانین مفید

$$\log\left(\frac{x}{y}\right) = \log(x) - \log(y)$$

$$\log(x^y) = y \log(x)$$

$$\sum_{i=1}^m a_i \int_x f_i(x) p(x) dx = \int_x \sum_{i=1}^m a_i f_i(x) p(x) dx \quad \triangleright \text{Can bring sum into integral}$$

$$\frac{d}{dx} \int_x f(x) p(x) dx = \int_x \frac{d}{dx} f(x) p(x) dx \quad \triangleright \text{Can (almost always) bring derivative into integral}$$

¹ squared

² general

³ Frobenius

پیشگفتار: شروع با یک مثال رگرسیون خطی

مفهوم یادگیری ماشین شامل طیف وسیعی از تکنیک‌ها برای عمل یادگیری از داده‌ها است. یک هدف اصلی (و هدف ما عمدتاً در این کتاب) پیش‌بینی است. بسیاری از تکنیک‌ها یک تابع $(f: \mathbb{R}^d \rightarrow \mathbb{R})$ را یاد می‌گیرند که یک ورودی‌های آن ویژگی‌ها یا نشانه‌هایی در مورد یک آیتم را وارد میکند و خروجی یک پیش‌بینی در مورد آن آیتم است. برای مثال، در نظر بگیرید که میخواهید قیمت یک خانه را بر اساس اطلاعات مربوط به آن خانه حدس بزنید یا پیش‌بینی کنید. این ویژگی‌ها ممکن است عبارت باشد از قدمت خانه، مساحت آن و فاصله آن تا نزدیک‌ترین بازار. بدون هیچ نمونه قبلی از هزینه‌های خانه، یعنی بدون هیچ داده‌ای، ممکن است حدس زدن این قیمت دشوار باشد. با این حال، تصور کنید مجموعه‌ای از ویژگی‌های چند خانه و هزینه‌های فروش مربوطه را برای خانه‌هایی که امسال فروخته شده‌اند به شما داده می‌شود. فرض کنید $\mathbf{x} \in \mathbb{R}^d$ بردار ویژگی‌های یک خانه باشد و در این مورد (قدمت خانه، مساحت و فاصله خانه تا نزدیک‌ترین نانوايي) $\mathbf{x} = [x_1 x_2 x_3] =$ [age, size, distancetobakery] باشد و هدف ما پیدا کردن قیمت ($y = price$) باشد. اگر ما از قبل ۱۰ نمونه از قیمت خانه‌ها را داشته باشیم، یک مجموعه داده: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{10}, y_{10})$ ، که در آن (\mathbf{x}_i, y_i) جفت (ویژگی، قیمت) برای خانه i در مجموعه نمونه‌های شما است. یک هدف طبیعی یافتن تابعی است که داده‌ها را به دقت بازآفرینی کند، به‌عنوان مثال با تلاش برای یافتن تابع f که منجر به تفاوت کوچکی بین پیش‌بینی، $f(\mathbf{x}_i)$ و قیمت واقعی، y_i ، برای هر خانه شود. ما می‌توانیم این را به عنوان یک مشکل بهینه‌سازی در نظر گرفته و فرمول‌بندی کنیم. تصور کنید فضایی از توابع ممکن \mathcal{H} داریم که می‌توانیم تابع f را از بین آن‌ها انتخاب کنیم. برای یک مورد ساده، اجازه دهید تصور کنیم که تابع خطی است: $f(\mathbf{x}) = w_0 + x_1 w_1 + x_2 w_2 + x_3 w_3$ که برای هر $\mathbf{w} = [w_0, w_1, w_2, w_3] \in \mathbb{R}^d$ که در آن w_0 همان نقطه قطع تابع خطی، یا همان عرض از مبدا است. می‌توانیم سعی کنیم تابعی از کلاس توابع خطی پیدا کنیم که مجموع مربع اختلاف‌ها را به مقدار مینیمم برساند یعنی:

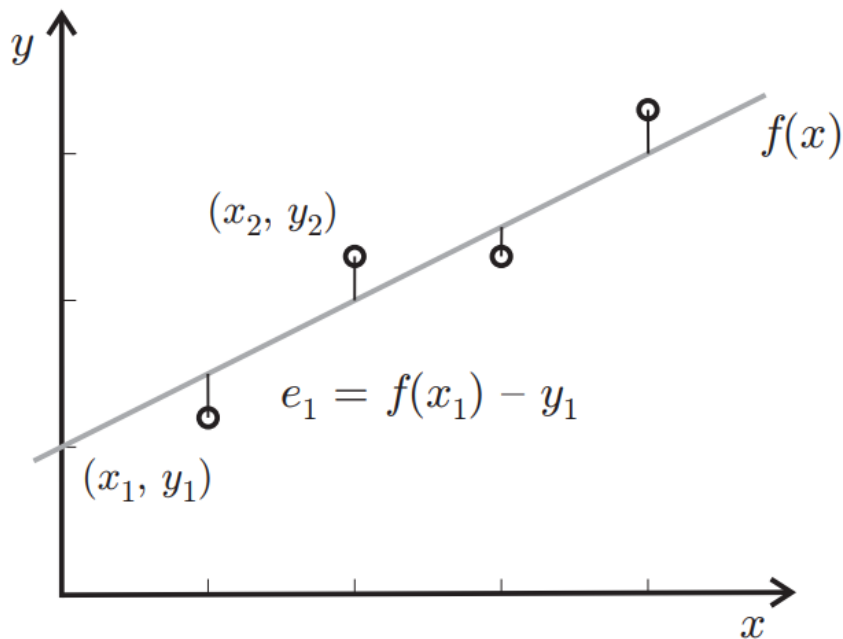
$$\min_{f \in H} \sum_{i=1}^{10} (f(\mathbf{x}_i) - y_i)^2$$

همانطور که بعداً در فصل ۵ خواهیم دید، حل این مسئله بهینه‌سازی برای توابع خطی ساده است. این روش شامل نوشتن صریح بهینه‌سازی بر حسب پارامترهای \mathbf{w} و حل \mathbf{w} بهینه است تا اندازه اختلاف‌ها را تا حد امکان کوچک می‌کند.

$$\text{Err} \stackrel{\text{def}}{=} \sum_{i=1}^{10} (f(x_i) - y_i)^2 = \sum_{i=1}^{10} \left(w_0 + \sum_{j=1}^d w_j x_{ij} - y_i \right)^2$$

راه حل یک خط مستقیم است که سعی می‌کند به بهترین نحو با اهداف مشاهده شده \mathcal{Y} مطابقت داشته باشد. یک تصویر ساده از چنین تابعی، تنها برای یک ویژگی، در شکل ۱ نشان داده شده است. هنگامی که این تابع را داریم، وقتی خانه جدیدی را می‌بینیم، امیدواریم که به اندازه کافی شبیه خانه‌های قبلی باشد تا این تابع به اندازه کافی آن را پیش‌بینی کند. قیمت خانه

تابع آموخته شده f بین این 10 نقطه درونیابی می‌شود تا نقاط نادیده را پیش‌بینی کند. اما، یک سوال طبیعی این است که آیا ما به خوبی درونیابی کردیم و آیا f آموخته شده می‌تواند پیش‌بینی دقیقی در مورد خانه‌های جدید ایجاد کند؟ اگر می‌خواهید از این تابع آموخته شده f در عمل استفاده کنید، می‌خواهید چنین خصوصیتی داشته باشید.



شکل ۱: نمونه‌ای از یک رگرسیون خطی مناسب در مجموعه داده $D = (1,1.2), (2,2.3), (3,2.3), (4,3.3)$ که وظیفه فرایند بهینه‌سازی برای یافتن بهترین تابع خطی $f(x) = w_0 + w_1x$ است به طوری که مجموع مجذور خطاها $(e_1^2 + e_2^2 + e_3^2 + e_4^2)$ به مقدار مینیمم می‌رسد.

در مورد درستی ادعای بالا، پاسخ به چنین سوالاتی بسیار دشوار است. ما می‌توانیم اصلاحات بصری انجام دهیم تا امیدوار باشیم پیش‌بینی‌های دقیق‌تری را ارائه دهیم، مانند گسترش کلاس توابع خطی به توابع غیر خطی پیچیده. اما، این تغییرات عملکردی هنوز به مشخص کردن صحت پیش‌بینی در خانه‌های جدید کمک نمی‌کند. در عوض، چیزی که از دست رفته مفهوم اعتماد به پیش‌بینی است. چقدر به پیش‌بینی‌ها اطمینان داریم؟ آیا خانه‌های قبلی را به اندازه کافی دیدیم که از این پیش‌بینی مطمئن باشیم؟ اعتبارسنجی منابع این داده‌ها چیست؟ همه این نوع سوالات نیاز به محاسبات احتمالی دارند. در این کتاب، ما با ارائه مقدمه‌ای بر احتمال شروع می‌کنیم تا مبنایی برای مقابله با عدم قطعیت در یادگیری ماشین فراهم کنیم. پس از آن که ابزارهای احتمالی را برای درک بهتر نحوه نزدیک شدن به پاسخ به این سوالات در اختیار داشتیم، به یادگیری این توابع باز می‌گردیم. بسیاری از پیشینه ریاضی مورد نیاز برای این کتاب، شامل درک اولیه احتمال و بهینه‌سازی است. البته این کتاب تلاش خواهد کرد تا بیشتر پیشینه مورد نیاز را در سراسر آن ارائه کند.

فصل ۱

مقدمه‌ای بر مدل‌سازی احتمالی

مدل‌سازی جهان پیرامون و پیش‌بینی وقوع رویدادها، یک تلاش چند رشته‌ای است که بر پایه‌های محکم نظریه احتمال، آمار و علوم رایانه قرار دارد. اگرچه این زمینه‌ها در فرآیند مدل‌سازی در هم تنیده شده‌اند، اما نقش‌های نسبتاً قابل تشخیصی دارند و تا حدی می‌توان آنها را به صورت جداگانه مورد مطالعه قرار داد. تئوری احتمالات زیرساخت‌های ریاضی را برای دستکاری احتمالات به ارمغان می‌آورد و ما را با طیف وسیعی از مدل‌ها با ویژگی‌های نظری کاملاً درک شده مجهز می‌کند. آمار چارچوب‌هایی را برای فرمول‌بندی استنتاج و فرآیند محدود کردن فضای مدل بر اساس داده‌های مشاهده شده و تجربه ما به منظور یافتن و سپس تحلیل راه‌حل‌ها فراهم می‌کند. علوم کامپیوتر نظریه‌ها، الگوریتم‌ها و نرم‌افزارهایی را برای مدیریت داده‌ها، محاسبه راه‌حل‌ها و مطالعه رابطه بین راه‌حل‌ها و منابع موجود (زمان، مکان، معماری کامپیوتر و غیره) در اختیار ما قرار می‌دهد. به این ترتیب، این سه رشته چارچوب کمی اصلی را برای همه علوم تجربی و فراتر از آن تشکیل می‌دهند. تئوری احتمالات و آمار سابقه نسبتاً طولانی دارند. ریشه شکل‌گیری هر دو را به طور می‌توان در قرن هفدهم دنبال کرد. نظریه احتمال از تلاش برای درک بازی‌های شانسی و قمار توسعه یافته است. مکاتبات بین بلز پاسکال و پیر دو فرما در سال ۱۶۵۴ به عنوان قدیمی‌ترین سابقه نظریه احتمالات مدرن است. از سوی دیگر، آمار از ابتکارات جمع‌آوری داده‌ها و تلاش‌ها برای درک روندها در جامعه (به عنوان مثال، تولید، علل مرگ و میر، قیمت زمین) و امور سیاسی (مانند درآمدهای عمومی، مالیات، ارتش) سرچشمه می‌گیرد. این دو رشته در قرن هجدهم با استفاده از داده‌ها برای اهداف استنتاجی در نجوم، جغرافیا و علوم اجتماعی شروع به ادغام کردند. افزایش پیچیدگی مدل‌ها و در دسترس بودن داده‌ها در قرن نوزدهم بر اهمیت ماشین‌های محاسباتی تاکید کرد. این به ایجاد پایه‌های حوزه علوم کامپیوتر در قرن بیستم کمک کرد، که عموماً به معرفی معماری فون نویمان و رسمی‌سازی مفهوم یک الگوریتم نسبت داده می‌شود. همگرایی این سه رشته در حال حاضر به جایگاه یک نظریه اصولی استنتاج احتمالی با کاربردهای گسترده در علم، تجارت، پزشکی، نظامی، مبارزات سیاسی و غیره رسیده است. مفاهیمی مانند توزیع بولتزمن، الگوریتم ژنتیک یا شبکه عصبی تأثیر فیزیک، زیست‌شناسی، روان‌شناسی و مهندسی را نشان می‌دهند. ما به فرآیند مدل‌سازی، استنتاج و تصمیم‌گیری بر اساس مدل‌های احتمالی به عنوان استدلال احتمالی یا استدلال تحت عدم قطعیت اشاره خواهیم کرد. نوعی استدلال در شرایط عدم قطعیت جزء ضروری زندگی روزمره است. به عنوان مثال، هنگام رانندگی، ما اغلب بر اساس انتظارات خود در مورد بهترین راه تصمیم می‌گیریم. در حالی که این موقعیت‌ها معمولاً مستلزم استفاده صریح از احتمالات و مدل‌های احتمالی نیستند، یک خودروی بدون راننده مانند شوfer گوگل^۱ باید از آنها استفاده کند. و همینطور یک نرم افزار تشخیص هرزنامه در یک اکانت ایمیل، یک سیستم تشخیص تقلب در کارت اعتباری، یا الگوریتمی که استنباط می‌کند که آیا یک جهش ژنتیکی خاص منجر به بیماری می‌شود یا خیر. بنابراین، ابتدا باید مفهوم احتمال را درک

¹ Google Chauffeur

کنیم و سپس یک نظریه رسمی برای ترکیب شواهد (مثلاً داده‌های جمع‌آوری‌شده از ابزارها) به منظور اتخاذ تصمیم‌های خوب در طیف وسیعی از موقعیت‌ها معرفی کنیم.

در سطح پایه، از احتمالات برای تعیین کمیت احتمال وقوع رویدادها استفاده می‌شود. همانطور که ژاکوب برنولی در کار خود به نام هنر حدس زدن (۱۷۱۳) به طرز درخشانی بیان می‌کند: «حدس زدن [پیش‌بینی] درباره چیزی مانند اندازه‌گیری احتمال آن است. بنابراین، ما هنر حدس زدن [علم پیش‌بینی] اتفاقی را به عنوان هنر اندازه‌گیری احتمالات اشیاء با دقت هر چه تمام‌تر تعریف می‌کنیم تا در قضاوت‌ها و اعمال، همیشه آنچه را که پیدا شده است انتخاب کنیم یا دنبال کنیم. و نتایج بهتر، رضایت بخش‌تر، ایمن‌تر یا با دقت بیشتری مورد توجه قرار گیرد.» تکنیک‌های مدل‌سازی احتمالی بسیاری از مفاهیم شهودی را رسمیت می‌بخشد. به طور خلاصه، آنها ابزارهایی را برای تجزیه و تحلیل دقیق ریاضی و استنتاج، اغلب در حضور شواهد، در مورد رویدادهایی که تحت تأثیر عواملی هستند که ما یا به طور کامل درک نمی‌کنیم یا کنترلی بر آنها نداریم، ارائه می‌دهند. برای ارائه بینشی سریع به مفهوم عدم قطعیت و مدل‌سازی، پرتاب یک تاس شش وجهی مناسب را در نظر بگیرید. اگر موقعیت اولیه، نیرو، اصطکاک، تو رفتگی‌های شکل و سایر عوامل فیزیکی را با دقت در نظر بگیریم و محاسبه کنیم، سپس آزمایش را اجرا کنیم، می‌توانیم به طور دقیق (یا اینطور فکر می‌کنیم) نتیجه را پیش‌بینی کنیم. اما قوانین فیزیکی ممکن است شناخته شده نباشند، ممکن است ترکیب آنها دشوار باشد یا حتی ممکن است چنین اقداماتی توسط قوانین آزمایش مجاز نباشد. بنابراین، عملاً مفید است که به سادگی فرض کنیم که هر نتیجه به یک اندازه محتمل است. در واقع، اگر تاس را بارها پرتاب کنیم، در واقع مشاهده می‌کنیم که هر عدد تقریباً به یک اندازه مشاهده می‌شود. تخصیص یک شانس (احتمال) برابر برای هر نتیجه از چرخش تاس، روشی کارآمد و ظریف برای مدل‌سازی عدم قطعیت‌های ذاتی آزمایش فراهم می‌کند.

مثال واقعی‌تر دیگری که در آن جمع‌آوری داده‌ها مبنایی را برای مدل‌سازی احتمالی ساده فراهم می‌کند، وضعیت رانندگی هر روز به محل کار و پیش‌بینی مدت زمانی است که فردا به مقصد می‌رسیم. اگر «زمان برای کار» را برای چند ماه ثبت می‌کردیم، مشاهده می‌کردیم که سفرها معمولاً بسته به عوامل داخلی (مثلاً سرعت ترجیحی ما برای روز) و همچنین عوامل خارجی (مانند آب‌وهوا، شلوغی جاده‌ها، تصادفات، یک راننده کند) زمان‌های متفاوتی طول می‌کشد. در حالی که این رویدادها، در صورت شناخته شدن، می‌توانند برای پیش‌بینی مدت زمان دقیق رفت و آمد مورد استفاده قرار گیرند، انتظار داشتن اطلاعات کامل نه تنها غیرواقعی است بلکه ما قابلیت مشاهده جزئی نداریم. ارائه راه‌هایی برای جمع‌آوری عوامل خارجی از طریق جمع‌آوری داده‌ها در یک دوره زمانی و ارائه توزیع زمان رفت و آمد مفید است. چنین توزیعی، در غیاب هر گونه اطلاعات دیگری، پس از آن استدلال در مورد رویدادهایی مانند رسیدن به موقع به یک جلسه مهم در ساعت ۹ صبح را تسهیل می‌کند.

تکنیک‌های مدل‌سازی احتمالی، یک چارچوب را برای برخورد با چنین آزمایش‌های تکراری تحت تأثیر تعدادی از عوامل خارجی که کنترل یا دانش کمی روی آنها داریم، ارائه می‌کند. با چنین چارچوب‌هایی، می‌توانیم نحوه پیش‌بینی‌های خود را بهتر درک کرده و بهبود ببخشیم، زیرا می‌توانیم پیش‌فرض‌های خود را در مورد عدم قطعیت خود با وضوح بیشتری مشخص کنیم و به صراحت درباره نتایج احتمالی استدلال کنیم. در این فصل به معرفی احتمالات و نظریه احتمالات، از ابتدا می‌پردازیم. از آنجایی که احتمال یک مفهوم اساسی در یادگیری ماشینی است، ارزش آن را دارد که بفهمیم از کجا آمده است. با این وجود، با پیروی از محتوای این یادداشت‌ها، یک درمان مختصر خواهد بود و بیشتر بر آنچه برای درک مطالب در فصل‌های بعدی نیاز است تمرکز خواهد کرد.

۱-۱ نظریه احتمال و متغیرهای تصادفی

نظریه احتمال به عنوان شاخه ای از ریاضیات است که به اندازه گیری احتمال وقوع رویدادها می پردازد. در قلب نظریه احتمال، مفهوم آزمایش وجود دارد. یک آزمایش می تواند فرآیند پرتاب کردن یک سکه، چرخاندن تاس، بررسی دمای فردا یا تعیین محل کلیدها باشد. وقتی انجام می شود، هر آزمایش یک نتیجه دارد، که عنصری است که از مجموعه ای از گزینه های از پیش تعریف شده، به طور بالقوه بی نهایت اندازه گرفته شده است. نتیجه یک چرخش تاس عددی بین یک تا شش است. دمای فردا ممکن است یک عدد حقیقی باشد. نتیجه مکان کلیدها می تواند مجموعه ای مجزا از مکان ها مانند میز آشپزخانه، زیر کاناپه، دفتر کار و غیره باشد. از بسیاری جهات، هدف اصلی مدل سازی احتمالی، فرمول بندی یک سوال خاص یا یک فرضیه مربوط به دنیای فیزیکی به عنوان یک آزمایش این است که داده ها را جمع آوری کرده و سپس یک مدل بسازید. هنگامی که یک مدل ایجاد شد، می توانیم معیارهای کمی مجموعه ای از نتایج مورد علاقه خود را محاسبه کنیم و اعتمادی را که باید به این معیارها داشته باشیم، ارزیابی کنیم. ما می توانیم قواعد احتمال را بر اساس مجموعه ای ساده از اصول به نام اصول احتمال بسازیم. اجازه دهید فضای نمونه (Ω) مجموعه ای ناتمام از نتایج و فضای رویداد (\mathcal{E}) مجموعه ای ناتمامی از زیر مجموعه های Ω باشد. به عنوان مثال، $\Omega = \{1, 2, 3\}$ و یک رویداد ممکن است $A = \{1, 3\} \in \mathcal{E}$ ، که در آن فضای رویداد این است که 1 یا 3 مشاهده می شود. فضای رویداد \mathcal{E} باید ویژگی های زیر را برآورده کند

$$1. \quad A \in \mathcal{E} \Rightarrow A^c \in \mathcal{E} \quad \text{به طوری که } A^c \text{ مکمل مجموعه } A \text{ است } A^c = \Omega - A$$

$$2. \quad A_1, A_2, \dots \in \mathcal{E} \Rightarrow \bigcup_{i=1}^{\infty} A_i \in \mathcal{E}$$

$$3. \quad \mathcal{E} \text{ یک مجموعه ناتمامی است} \quad \{\emptyset \in \mathcal{E} \text{ و } \Omega \in \mathcal{E}\}$$

اگر \mathcal{E} این سه ویژگی را برآورده کند، آنگاه به (Ω, \mathcal{E}) فضای قابل اندازه گیری گفته می شود. اکنون می توانیم اصول احتمال را تعریف کنیم، که روشن تر می کند که چرا این دو شرط برای فضای رویداد ما برای تعریف احتمالات معنادار بر روی رویدادها مورد نیاز است. تابع $P : \mathcal{E} \rightarrow [0, 1]$ اصول احتمال را برآورده می کند اگر

$$1. \quad P(\Omega) = 1$$

$$2. \quad A_1, A_2, \dots \in \mathcal{E}, A_i \cap A_j = \emptyset \forall i, j \Rightarrow P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i)$$

که به آن اندازه گیری احتمال یا توزیع احتمال گفته می شود. چندتایی (Ω, \mathcal{A}, P) فضای احتمال نامیده می شود. زیبایی این اصول در فشردگی و ظرافت آنها نهفته است. بسیاری از عبارات مورد استفاده را می توان از بدیهیات احتمال استخراج کرد. برای مثال، واضح است که $P(A^c) = 1 - P(A)$. این امر روشن تر می کند که چرا ما نیاز داریم که اگر یک رویداد در فضای رویداد باشد، مکمل آن نیز باید در فضای رویداد باشد: اگر بتوانیم احتمال یک رویداد را اندازه گیری کنیم، می دانیم که احتمال رخ ندادن آن رویداد به طور مشابه 1 منهای این احتمال است. ما مستلزم این هستیم که اگر دو رویداد A_1, A_2 بدون اشتراک باشند، آنگاه: $P(A_1 \cup A_2) = P(A_1) + P(A_2)$ احتمال وقوع هر یک از این رویدادها برابر مجموع احتمالات آنها است، زیرا هیچ همپوشانی وجود ندارد. در نتیجه در رویدادها ویژگی دیگری که می توانیم استنباط کنیم این است که ما همیشه $\emptyset, \Omega \in \mathcal{E}$ داریم، جایی که \emptyset مربوط به رویدادی است که در آن هیچ اتفاقی نمی افتد - که باید احتمال آن صفر باشد. اگر به قوانین دیگری که می توان استخراج کرد علاقه مند هستید، به پیوست A.1 مراجعه کنید.

مثال ۱: [متغیرهای گسسته (قابل شمارش)] احتمالات ریختن یک تاس را مدلسازی کنید.

فضای نتیجه مجموعه محدود $\Omega = \{1, 2, 3, 4, 5, 6\}$ و فضای رویداد \mathcal{E} مجموعه همه زیر مجموعه‌ها است، یعنی $P(\Omega) = \{\emptyset, \{1\}, \{2\}, \dots, \{2, 3, 4, 5, 6\}, \Omega\}$ که از همه زیر مجموعه‌های ممکن Ω تشکیل شده است. یک توزیع احتمال طبیعی روی (Ω, \mathcal{E}) به هر تاس $\frac{1}{6}$ شانس وقوع می‌دهد، که به صورت $P(\{x\}) = \frac{1}{6}$ برای $x \in \Omega$ ، $P(\{1, 2\}) = \frac{2}{6}$ تعریف می‌شود. و غیره.

مثال ۲: متغیرهای پیوسته (غیر قابل شمارش) احتمال‌های زمان توقف خودرو را در محدوده ۳ تا ۶ ثانیه مدل کنید.

فضای نتیجه بازه پیوسته $\Omega = [3, 6]$ است. یک رویداد می‌تواند این باشد که خودرو در عرض ۳ تا ۳٫۱ ثانیه متوقف شود و به صورت $A = [3, 3.1] \in \mathcal{E}$ در نظر گرفته شود. احتمال $P(A)$ چنین رویدادی کم است، زیرا زمان توقف بسیار سریع خواهد بود. سپس می‌توانیم تمام فواصل زمانی ممکن برای فضای رویداد و احتمالات مربوطه را در نظر بگیریم. ما قبلاً دیدیم که این برای متغیرهای پیوسته کمی پیچیده‌تر خواهد بود، و بنابراین ما با دقت بیشتری نحوه تعریف \mathcal{E} و P را در زیر در بخش ۱٫۲٫۲ نشان می‌دهیم.

این دو مثال دو مورد متداول را که با آن مواجه خواهیم شد نشان می‌دهد: متغیرهای گسسته و متغیرهای پیوسته. عبارات فوق - قابل شمارش و غیرقابل شمارش - نشان می‌دهد که آیا یک مجموعه قابل شمارش است یا خیر. به عنوان مثال، مجموعه اعداد طبیعی را می‌توان شمارش کرد، و بنابراین قابل شمارش است، در حالی که مجموعه اعداد حقیقی را نمی‌توان شمارش کرد - همیشه یک عدد حقیقی دیگر بین هر دو عدد حقیقی وجود دارد - و بنابراین غیرقابل شمارش است. اگرچه این تمایز منجر به تفاوت‌های واقعی می‌شود - مانند استفاده از مجموع برای مجموعه‌های قابل شمارش و انتگرال‌ها برای مجموعه‌های غیرقابل شمارش - قالب و شهود تا حد زیادی بین این دو دسته منتقل می‌شوند. ما بیشتر بر روی متغیرهای گسسته و پیوسته تمرکز خواهیم کرد. بسیاری از ایده‌های مشابه به متغیرهای مختلط نیز منتقل می‌شوند، جایی که فضاهای نتیجه از هر دو مجموعه گسسته و پیوسته مانند $\Omega = [0, 1] \cup \{2\}$ تشکیل شده‌اند. علاوه بر این، برای تنظیم غیرقابل شمارش، ما به طور خاص مجموعه‌های پیوسته را مورد بحث قرار می‌دهیم، به عنوان مثال، آن‌ها اتحادیه‌های بازه‌های پیوسته مانند $\Omega = [0, 1] \cup [5, 10]$ هستند. از آنجا که تقریباً تمام مجموعه‌های غیرقابل شمارش که می‌خواهیم در نظر بگیریم پیوسته هستند، برای تعیین چنین فضاهایی از اصطلاحات پیوسته و غیرقابل شمارش استفاده می‌کنیم. در نهایت، مجموعه‌های گسسته می‌توانند متناهی باشند، مانند $\{1, 2, 3\}$ ، یا قابل شمارش بی نهایت، مانند اعداد طبیعی. مجموعه‌های پیوسته به وضوح نامتناهی هستند و گفته می‌شود که به طور غیرقابل شمارش نامتناهی هستند.

قبل از توضیح بیشتر در مورد چگونگی تعریف توزیع احتمال، ابتدا متغیرهای تصادفی را معرفی می‌کنیم و از اینجا به بعد به طور دقیق به متغیرهای تصادفی می‌پردازیم. یک متغیر تصادفی به ما اجازه می‌دهد تا تبدیل‌های فضاهای احتمال را با دقت بیشتری تعریف کنیم. هنگامی که آن تبدیل را اجرا می‌کنیم، می‌توانیم فضای احتمال زیربنایی را فراموش کنیم و می‌توانیم روی رویدادها و توزیع فقط بر روی متغیر تصادفی تمرکز کنیم. این در واقع همان کاری است که شما به طور طبیعی هنگام تعریف احتمالات بر روی متغیرها انجام می‌دهید، بدون اینکه نیازی به فرمول‌بندی کردن آن به صورت ریاضی باشد. البته در اینجا آن را به فرمول‌های ریاضی تبدیل می‌کنیم.

دوباره مثال تاس را در نظر بگیرید، جایی که اکنون در عوض ممکن است بخواهید بدانید: احتمال دیدن یک عدد کوچک در بازه $(1 - 3)$ یا یک عدد بزرگ در بازه $(4 - 6)$ چقدر است؟ می‌توانیم فضای احتمال جدیدی را با $\Omega_x = \{\text{high}, \text{low}\}$

تعریف میشود: $P_x(\{\text{low}\}) = P_x(\{\text{high}\}) = \frac{1}{2}$ و $\mathcal{E}_x = \{\text{high}, \text{low}\}$ به این صورت $X = \Omega \rightarrow \Omega_x$ تابع تبدیل کنیم. تعریف کنیم.

$$X(\omega) \stackrel{\text{def}}{=} \begin{cases} \text{low} & \text{if } \omega \in \{1, 2, 3\} \\ \text{high} & \text{if } \omega \in \{4, 5, 6\} \end{cases}$$

$$X(\omega) \stackrel{\text{def}}{=} \begin{cases} \text{low} & \text{if } \omega \in \{1, 2, 3\} \\ \text{high} & \text{if } \omega \in \{4, 5, 6\} \end{cases}$$

توزیع P_x بلافاصله از این تبدیل تعیین می شود. به عنوان مثال $P_x(\{\text{low}\}) = P(\{\omega : X(\omega) = \text{low}\})$ زیرا فضای احتمال low احتمال دیدن 1، 2 یا 3 را نشان می دهد. اکنون می توانیم به سؤالاتی در مورد احتمال پاسخ دهیم. دیدن یک عدد کوچک از یک عدد عضو مجموع high.

این تابع X یک متغیر تصادفی نامیده می شود. این کمی گیج کننده است که نه تصادفی است و نه یک متغیر، زیرا X یک تابع است. با این حال، از اینجا به بعد، برای تابع X به جای نوشتن $P(\{\omega : X(\omega)\})$ با نوشتن عباراتی مانند $P_x(X = x)$ یا $P_x(X \in A)$ مانند یک متغیر تصادفی رفتار می کنیم $P(\{\omega : X(\omega) = x\})$ یا $P(\{\omega : X(\omega) \in A\})$. برای صحت موضوع، می توانیم به یاد داشته باشیم که تابعی است که در فضای احتمالی زیربنایی پیچیده تر تعریف شده است. اما، در عمل، می توانیم مستقیماً بر حسب متغیر تصادفی X و احتمالات مرتبط با آن فکر کنیم. به طور مشابه، حتی برای نقش تاس، می توانیم تصدیق کنیم که فضای احتمالی پیچیده تری وجود دارد که توسط پویایی تاس تعریف می شود. هنگامی که فقط احتمالات نتایج گسسته از 1-6 را در نظر می گیریم، ما قبلاً به طور ضمنی تغییری را در بالای احتمالات سیستم فیزیکی اعمال کرده ایم. هنگامی که یک متغیر تصادفی داریم، یک فضای احتمال معتبر $(\Omega_x, \mathcal{E}_x, P_x)$ را تعریف می کند. بنابراین، همه قواعد احتمال یکسان اعمال می شوند، درک یکسانی از نحوه تعریف توزیع ها، و غیره. در واقع، ما همیشه می توانیم یک متغیر تصادفی X را که مطابق با هیچ تبدیلی نیست، تعریف کنیم تا فضای احتمال اصلی را بدست آوریم. به همین دلیل، با فرض اینکه همیشه با متغیرهای تصادفی سروکار داریم، بدون از دست دادن کلیت، می توانیم جلو برویم. ما اندیس ها را حذف می کنیم و (Ω, \mathcal{E}, P) را برای X تعریف می کنیم. برای بحث عمیق تر در مورد متغیرهای تصادفی، به پیوست A.4 مراجعه کنید.

۱-۲ تعریف توزیع

اکنون می خواهیم بدانیم که چگونه P را برای اصول احتمال مشخص کنیم، تا احتمال X را برای یک رویداد A ، $P(X \in A)$ با فضای نتیجه Ω مدل کنیم. این کار دلهره آور به نظر می رسد، زیرا به نظر می رسد ما باید احتمال هر رویداد ممکن را تعریف کنیم - مجموعه ای از نتایج - و به گونه ای که اصول احتمال را برآورده کند، نه! خوشبختانه، در عوض می توانیم توزیع را با استفاده از تابعی تعریف کنیم که مستقیماً روی نمونه های Ω تعریف شده است. در نظر گرفتن جداگانه فضاهای نمونه گسسته (قابل شمارش) و پیوسته (غیر قابل شمارش) راحت است. برای حالت گسسته، توابع جرم احتمال و برای حالت پیوسته، توابع چگالی احتمال را تعریف می کنیم.

۱-۲-۱ توابع جرم احتمال

فرض کنید Ω یک فضای نمونه گسسته و $\mathcal{E} = P(\Omega)$ ، مجموعه زیر مجموعه‌های Ω باشد. تابع $p: \Omega \rightarrow [0, 1]$ تابع جرم احتمال^۱ به اختصار (pmf) نامیده می‌شود اگر

$$\sum_{\omega \in \Omega} p(\omega) = 1$$

احتمال هر رویداد $A \in \mathcal{E}$ به صورت تعریف شده است

$$P(A) \stackrel{\text{def}}{=} \sum_{\omega \in \Omega} p(\omega) = 1$$

راستی آزمایی اینکه P اصول احتمال را برآورده می‌کند و بنابراین یک توزیع احتمال است ساده است. بنابراین، برای متغیرهای تصادفی گسسته، ما اغلب $P(X = x)$ را می‌نویسیم، به این معنی که $P(X = x) = p(x)$ برای هر نتیجه $x \in \Omega$ ما به ندرت توزیع را مستقیماً تعریف می‌کنیم، و به جای آن، p pmf را که توزیع P را القا می‌کند، تعریف می‌کنیم.

مثال ۳: یک پرتاب از تاس شش وجهی منصفانه را در نظر بگیرید. یعنی $\Omega = \{1, 2, 3, 4, 5, 6\}$ و فضای رویداد $\mathcal{E} = P(\Omega)$. احتمال اینکه نتیجه عددی بزرگتر از 4 باشد چقدر است؟

اول، چون تاس منصفانه است، می‌دانیم که $p(\omega) = \frac{1}{6}$ برای $\forall \omega \in \Omega$. حال، اجازه دهید A یک رویداد در \mathcal{E} باشد که نتیجه آن بزرگتر از 4 باشد. یعنی $A = \{5, 6\}$ بدین ترتیب،

$$P(A) \stackrel{\text{def}}{=} \sum_{\omega \in A} p(\omega) = \frac{1}{3}$$

توجه داشته باشید که توزیع P بر روی عناصر \mathcal{E} تعریف شده است، در حالی که p بر روی عناصر Ω تعریف شده است. یعنی $P(\{1\}) = p(1)$, $P(\{2\}) = p(2)$, $P(\{1, 2\}) = p(1) + p(2)$, ...

بنابراین، برای مشخص کردن P ، باید نحوه تعیین pmf ، یعنی احتمال هر نتیجه گسسته را تعیین کنیم. pmf اغلب به عنوان جدولی از مقادیر احتمال مشخص می‌شود. برای مثال، برای مدل‌سازی احتمال تولد برای هر روز در سال، می‌توان جدولی از ۳۶۵ مقدار بین صفر و یک داشت، تا زمانی که مجموع احتمالات برابر با 1 باشد. این احتمالات را می‌توان از داده‌های مربوط به تولد افراد محاسبه کرد با استفاده از شمارش برای هر روز و عادی سازی تعداد کل افراد در جمعیت برای تخمین احتمال دیدن تولد در یک روز معین. چنین جدول مقادیر بسیار منعطف است و امکان تعیین مقادیر احتمال دقیق برای هر نتیجه را فراهم می‌کند. با این حال، چند pmf مفید وجود دارد که شکل عملکردی (محدودتر) دارند. ما سه pmf از این قبیل را در اینجا شرح می‌دهیم که در سراسر این کتاب از آنها استفاده خواهیم کرد. برای مثال‌های بیشتر از pmf ، به پیوست A.2 مراجعه کنید.

توزیع برنولی از مفهوم آزمایش برنولی ناشی می‌شود، آزمایشی که دو نتیجه ممکن دارد: موفقیت و شکست. در آزمایش برنولی، موفقیت با احتمال $\alpha \in [0, 1]$ و بنابراین، شکست با احتمال $(1 - \alpha)$ رخ می‌دهد. پرتاب یک سکه (سردم)، یک بازی

¹ Probability mass functions

بسکتبال (برد/باخت)، یا یک چرخش قالب (زوج/فرد) همگی به عنوان آزمایش برنولی دیده می‌شوند. ما این توزیع را با تنظیم فضای نمونه روی دو عنصر و تعریف احتمال یکی از آنها به عنوان α مدل می‌کنیم. به طور خاص، $\Omega = \{\text{success, failure}\}$ و داریم:

$$p(\omega) = \begin{cases} \alpha & \omega = \text{success} \\ 1 - \alpha & \omega = \text{failure} \end{cases}$$

که $\alpha \in (0, 1)$ یک پارامتر است. اگر به جای آن $\Omega = \{0, 1\}$ را در نظر بگیریم، می‌توانیم توزیع برنولی را به طور فشرده به صورت $p(\omega) = \alpha^\omega (1 - \alpha)^{1-\omega}$ برای $\omega \in \Omega$ بنویسیم. توزیع برنولی اغلب به صورت $\text{Bernoulli}(\alpha)$ نوشته می‌شود. همانطور که خواهیم دید، یک نگاشت رایج که در آن از برنولی استفاده می‌کنیم برای طبقه بندی دوتایی است، مثلاً جایی که سعی می‌کنیم پیش‌بینی کنیم که آیا بیمار آنفولانزا دارد (نتیجه 0) یا آنفولانزا ندارد (نتیجه 1).

توزیع یکنواخت برای فضاهای نمونه گسسته بر روی مجموعه محدودی از نتایج تعریف می‌شود که احتمال وقوع هر کدام به یک اندازه است. اجازه دهید $\Omega = \{1, \dots, n\}$; سپس برای $\forall \omega \in \Omega$ قرار دهیم

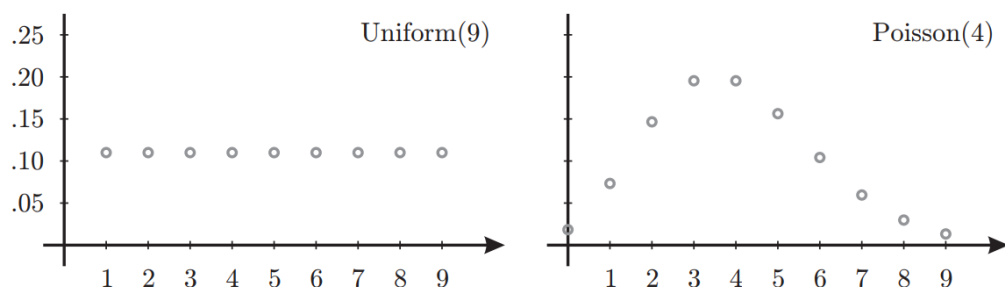
$$P(\omega) = \frac{1}{n}$$

توزیع یکنواخت شامل پارامترها نیست. با اندازه فضای نمونه تعریف می‌شود. ما به این توزیع به عنوان $\text{Uniform}(n)$ اشاره می‌کنیم. بعداً خواهیم دید که توزیع یکنواخت را می‌توان در فواصل محدود در فضاهای پیوسته نیز تعریف کرد.

توزیع پواسون منعکس کننده احتمال وقوع چند حادثه است (به طور ضمنی در یک بازه زمانی ثابت). به عنوان مثال، یک مرکز تماس که به احتمال زیاد ۵۰ تماس در ساعت دریافت می‌کند، با احتمال بسیار کمتری ۵ تماس یا ۱۰۰۰ تماس دریافت می‌کند. این را می‌توان با $\text{Poisson}(\lambda)$ مدل کرد، که در آن λ تعداد تماس‌های مورد انتظار را نشان می‌دهد. به طور رسمی تر، $\forall \omega \in \Omega$ و $\Omega = \{0, 1, \dots\}$

$$p(\omega) = \frac{\lambda^\omega e^{-\lambda}}{\omega!}$$

این تابع توده‌ای به شکل تپه است، جایی که بالای تپه بیشتر در مرکز λ قرار دارد و سمت چپ تپه کوتاه و پر شیب و یک دم سمت راست بلند و کم شیب وجود دارد. توزیع پواسون بر روی یک فضای نمونه بی‌نهایت تعریف شده است، اما همچنان قابل شمارش است. این در شکل ۱-۱ نشان داده شده است.



شکل ۱-۱: دو تابع جرم احتمالی، برای متغیرهای تصادفی گسسته. توزیع پواسون بیشتر روی محور x ادامه می‌یابد (برای متغیر $\omega \in N$)، با احتمال کاهش به صفر به صورت $\omega \rightarrow \infty$

تمرین ۱: ثابت کنید که $\sum_{\omega \in \Omega} p(\omega) = 1$ برای توزیع پواسون.

مثال ۴: به عنوان مقدمه‌ای برای تخمین پارامترهای توزیع‌ها، مثالی از نحوه استفاده از توزیع برنولی و تعیین پارامتر α برای برنولی را در نظر بگیرید. یک مثال متعارف برای توزیع‌های برنولی، پرتاب سکه است که در آن نتایج سر (H) یا دم (T) هستند. $P(X = H) = \alpha$ احتمال دیدن H و $P(X = T) = 1 - \alpha$ احتمال دیدن T است. ما معمولاً $\alpha = 0.5$ را فرض می‌کنیم. این سکه منصفانه (بی طرفانه) نامیده می‌شود. اگر سکه را بارها پرتاب کنیم، انتظار داریم تقریباً تعداد برابر H و T را ببینیم. با این حال، یک سکه ناهمگن ممکن است مقداری انحراف به سمت H یا T داشته باشد. اگر سکه را بارها پرتاب کنیم، اگر $\alpha > 0.5$ باشد، باید در نهایت ببینیم که H بیشتری مشاهده می‌شود و اگر $\alpha < 0.5$ باشد، باید T بیشتری را ببینیم.

چگونه ممکن است واقعاً مقدار α را تعیین کنیم؟ یک ایده شهودی این است که از آزمایش‌های مکرر (داده‌ها) استفاده کنید، درست همانطور که در بالا توضیح داده شد: سکه را بارها پرتاب کنید تا ببینید آیا می‌توانید انحراف را اندازه بگیرید. اگر تعداد 1000 تا H و 50 تا T را ببینید، یک حدس طبیعی برای انحراف $\alpha \approx \frac{1000}{1000+50} = 0.95$ است. چقدر به این راه حل اطمینان دارید؟ آیا قطعاً 0.95 است؟ و چگونه می‌توانیم به طور رسمی‌تر تعریف کنیم که چرا این باید راه حل باشد؟ این در واقع یک راه حل معقول است و مطابق با راه حل ماکسیمم احتمال است، همانطور که در فصل ۳ بحث خواهیم کرد.

۲-۱-۲ توابع چگالی احتمال

محاسبات فضاهای احتمال پیوسته مشابه فضاهای گسسته است، با توابع چگالی احتمال^۱ یا به اختصار pdf که جایگزین توابع جرم احتمال و انتگرال‌ها جایگزین مجموع می‌شوند. با این حال، در تعریف pdfها، ما نمی‌توانیم از جداول مقادیر استفاده کنیم و به فرم‌های تابعی محدود می‌شویم. دلیل اصلی این تفاوت از این واقعیت ناشی می‌شود که دیگر منطقی نیست که احتمال یک رویداد تکی را اندازه گیری کنیم. دوباره زمان توقف خودرو را در نظر بگیرید، که در مثال ۲ مورد بحث قرار گرفت. اینکه احتمال توقف خودرو را دقیقاً در 3.14159625 ثانیه بپرسید، چندان منطقی نیست. به طور واقع بینانه، احتمال چنین رویداد دقیقی بسیار کم است. در واقع، احتمال مشاهده دقیق زمان توقف صفر است، زیرا مجموعه {3.14159625} به عنوان زیر مجموعه $[3, 6]$ مجموعه‌ای با اندازه صفر است. اساساً، جرم صفر را در داخل بازه $[3, 6]$ می‌گیرد که در نهایت به طور غیرقابل شمارش بی‌نهایت است. در عوض، ما باید احتمالات فواصل، مانند $[4, 5]$ یا $[5.667, 5.668]$ را در نظر بگیریم.

برای فضاهای پیوسته، فرض می‌کنیم که مجموعه رویدادهای \mathcal{E} شامل تمام فواصل ممکن است که به آن میدان بول $B(\Omega)$ می‌گویند. به عنوان مثال، اگر $\Omega = R$ ، میدان بول $B(R)$ شامل تمام بازه‌های باز (به عنوان مثال $(0, 1)$)، بازه‌های بسته (مثلاً $[0, 1]$) و بازه‌های نیمه باز (مثلاً $(0, 1]$) در R ، و همچنین مجموعه‌هایی که می‌توان با تعداد قابل شمارشی از عملیات‌های مجموعه پایه روی آنها، مانند اجتماع، به دست آورد. این منجر به مجموعه‌ای محدودتر از رویدادها نسبت به مجموعه توان Ω می‌شود، که برای مثال شامل مجموعه‌هایی با یک رویداد تکی می‌شود. $B(R)$ هنوز یک مجموعه عظیم است - یک مجموعه نامتناهی غیرقابل شمارش - اما کوچکتر از $P(R)$ است. با این حال، به خوبی، $B(R)$ هنوز هم شامل تمام مجموعه‌هایی است که می‌خواهیم بتوانیم اندازه‌گیری کنیم است. میدان Boirel را می‌توان برای هر فضای قابل اندازه‌گیری، مانند فضاهای

¹ Probability density functions

با ابعاد بالاتر مانند $\Omega = \mathbb{R}^2$ ، با رویدادهایی مانند $\Omega \supset A = [0, 1] \times [0, 1]$ یا $A = [1, 2] \times [-1, 4]$ $[0, 0.1] \times [10, 1000]$ تعریف کرد.

اکنون Ω یک فضای نمونه پیوسته و $B(\Omega) = \mathcal{E}$ باشد. تابع $p : \Omega \rightarrow [0, \infty)$ تابع چگالی احتمال (pdf) نامیده می‌شود اگر

$$\int_{\Omega} p(\omega) d\omega = 1$$

احتمال یک رویداد $A \in B(\Omega)$ به صورت تعریف می‌شود

$$P(A) \stackrel{\text{def}}{=} \int_A p(\omega) d\omega$$

توجه داشته باشید که تعریف pdf تنها به داشتن محدوده $[0, 1]$ محدود نمی‌شود، بلکه به $(0, \infty]$ محدود می‌شود. برای pmf، احتمال یک رویداد تکی $\{\omega\}$ مقدار pmf در نقطه نمونه ω است. یعنی $P(\{\omega\}) = p(\omega)$. از آنجایی که توزیع‌های احتمال P به محدوده $[0, 1]$ محدود می‌شوند، این بدان معناست که pmf‌ها نیز باید به آن محدوده محدود شود. در مقابل، مقدار pdf در نقطه ω یک احتمال نیست. در واقع می‌تواند بزرگ‌تر از 1 باشد. اگر چه زمانی که به طور غیررسمی صحبت می‌شود معمولاً $p(x)$ را احتمال x می‌نامیم، به‌طور دقیق‌تر آن را چگالی در x می‌نامیم، زیرا قطعاً یک احتمال نیست. در واقع، همانطور که در بالا گفته شد، احتمال در هر نقطه منفرد 0 است (یعنی یک زیر مجموعه قابل شمارش از Ω مجموعه ای از اندازه‌گیری صفر است).

یک سردرگمی طبیعی این است که چگونه p می‌تواند با 1 ادغام شود، اما در واقع مقادیری بزرگتر از 1 داشته باشد. بازه کوچک $A = [x, x + \Delta x]$ را در نظر بگیرید با احتمال

$$P(A) = \int_x^{x+\Delta x} p(\omega) d\omega \approx p(x) \Delta x$$

مقدار بالقوه بزرگ تابع چگالی با بازه کوچک Δx جبران می‌شود تا عددی بین 0 و 1 به دست آید. بنابراین، حتی اگر $p(x)$ یک میلیون باشد، چگالی نقاط در یک بازه کوچک می‌باشد. احتمال یک رویداد باید هنوز ≤ 1 باشد. چگالی نشان می‌دهد که احتمال بالایی در اطراف آن نقطه وجود دارد. با داشتن یک چگالی بزرگ در اطراف x ، این نشان می‌دهد که چگالی برای نقاط دیگر صفر یا نزدیک به صفر است و pdf به شدت حول x به اوج خود رسیده است.

برخلاف pmf، ما نمی‌توانیم p pdf را به این راحتی تعریف کنیم تا به‌طور انعطاف‌پذیر احتمالات خاصی را برای هر نتیجه با جدولی از احتمالات ارائه کنیم. بلکه برای pdf معمولاً از pdf شناخته شده‌ای استفاده می‌کنیم که ویژگی‌های مورد نیاز را برآورده کند. علاوه بر این، بر خلاف حالت گسسته، ما هرگز $P(X = x)$ را نخواهیم نوشت، زیرا این عدد صفر خواهد بود. در عوض، ما معمولاً $P(X \in A)$ یا سوالات احتمالی صریح‌تر مانند $P(X \leq 5)$ را می‌نویسیم. ما در اینجا چهار pdf را نشان می‌دهیم که در سراسر این کتاب استفاده می‌شود. برای نمونه‌های بیشتر از PDFها، به پیوست A.3 مراجعه کنید.

توزیع یکنواخت با مقدار مساوی یک تابع چگالی احتمال در یک بازه محدود در \mathbb{R} تعریف می‌شود. بنابراین، برای $\Omega = [a, b]$ تابع چگالی احتمال یکنواخت $\forall \omega \in [a, b]$ به این صورت تعریف می‌شود.

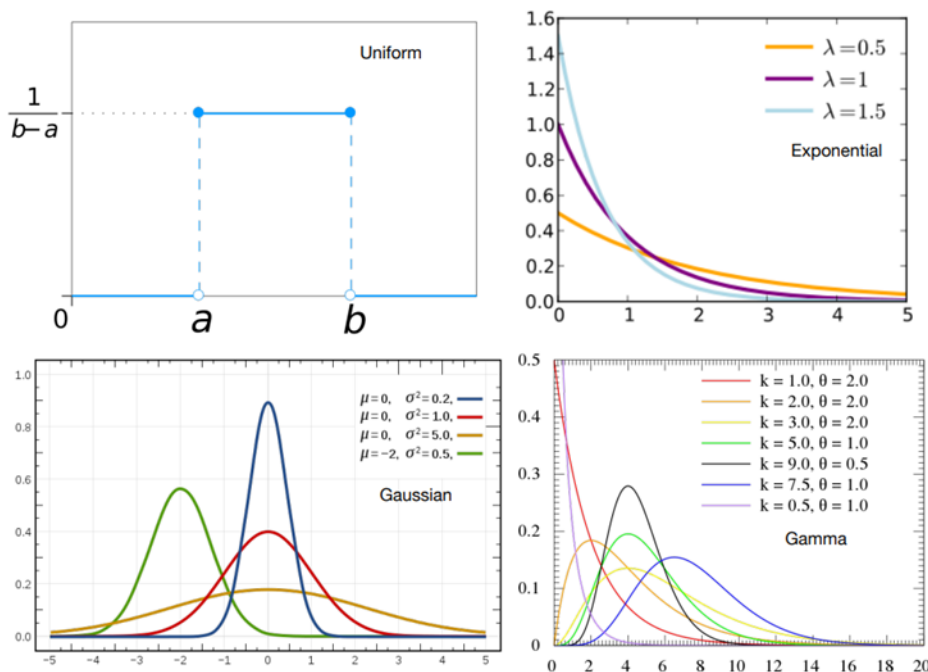
$$p(\omega) \stackrel{\text{def}}{=} \frac{1}{b - a}$$

همچنین می‌توان با در نظر گرفتن $\Omega = \mathbb{R}$ و تنظیم $p(\omega) = 0$ هر زمان که ω خارج از $[a, b]$ باشد، $\text{Uniform}(a, b)$ را تعریف کرد. این فرم مناسب است زیرا $\Omega = \mathbb{R}$ می‌تواند به طور مداوم برای همه توزیع‌های احتمال یک بعدی استفاده شود. وقتی این کار را انجام دادیم، به زیرمجموعه \mathbb{R} اشاره خواهیم کرد که در آن $p(\omega) > 0$ به عنوان پشتیبان تابع چگالی است.

توزیع نمایی بر روی مجموعه‌ای از اعداد غیر منفی تعریف می‌شود. یعنی $\Omega = [0, \infty)$. برای پارامتر $\lambda > 0$ pdf آن به صورت:

$$p(\omega) = \lambda e^{-\lambda\omega}$$

همانطور که از نام آن پیداست، این pdf شکل نمایی دارد و با افزایش بزرگی مقادیر x ، احتمال به شدت کاهش می‌یابد. مانند قبل، فضای نمونه را می‌توان به همه اعداد حقیقی گسترش داد، در این صورت برای $\omega < 0$ ، $p(\omega) = 0$ را قرار می‌دهیم.



شکل ۱-۲: چهار تابع چگالی احتمال، برای متغیرهای تصادفی پیوسته. تصاویر برگرفته از ویکی پدیا

توزیع گاوسی یا **توزیع نرمال** یکی از پرکاربردترین توزیع‌های احتمال است. بر روی $\Omega = \mathbb{R}$ ، با دو پارامتر $\mu \in \mathbb{R}$ و $\sigma > 0$ pdf تعریف شده است.

توزیع نرمال دارای خواص زیر است:

- میانگین = میانه = مد
- خط تقارن در وسط قرار می‌گیرد
- نیمی از داده‌ها کوچکتر از میانگین و نیمی دیگر بزرگتر از میانگین

$$p(\omega) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(\omega-\mu)^2}$$

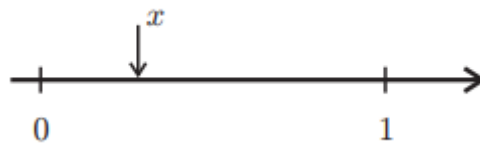
همانطور که در ادامه بحث خواهیم کرد، برای یک متغیر تصادفی که گاوسی توزیع شده است، پارامتر μ میانگین یا مقدار مورد انتظار و σ^2 واریانس است. ما به این توزیع به عنوان $\text{Gaussian}(\mu, \sigma^2)$ یا $N(\mu, \sigma^2)$ اشاره خواهیم کرد. وقتی میانگین صفر و واریانس یک باشد (واریانس واحد)، این گاوسی نرمال استاندارد نامیده می‌شود. نام این تابع گاوسی خاص به این دلیل نام دارد که بسیار مورد استفاده قرار می‌گیرد. هر دو توزیع گاوسی و نمایی اعضای خانواده وسیع‌تری از توزیع‌ها به نام خانواده نمایی طبیعی هستند. تعریف کلی این خانواده را بعداً در بخش ۷،۲ خواهیم دید.

توزیع لاپلاس شبیه به توزیع گاوسی است، اما در محدوده میانگین اوج بیشتری دارد. همچنین بر روی $\Omega = \mathbb{R}$ با دو پارامتر، $\mu \in \mathbb{R}$ و $b > 0$ pdf و تعریف شده است.

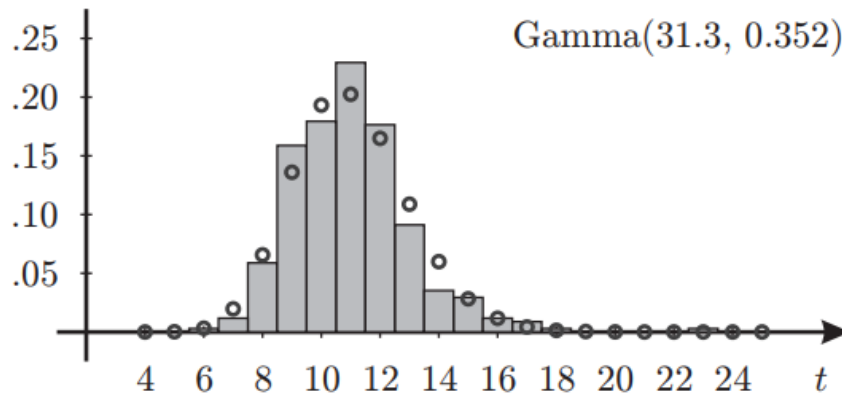
$$p(\omega) = \frac{1}{2b} e^{-\frac{1}{b}|\omega-\mu|}$$

توزیع گاما برای مدل‌سازی زمان‌های انتظار استفاده می‌شود و مشابه توزیع پواسون است اما برای متغیرهای پیوسته. بر روی $\Omega = (0, \infty)$ با پارامتر شکل $\alpha > 0$ و پارامتر نرخ $\beta > 0$ pdf تعریف شده است.

$$p(\omega) = \frac{\beta^\alpha}{\Gamma(\alpha)} \omega^{\alpha-1} e^{-\beta\omega}$$



شکل ۳-۱: انتخاب یک عدد تصادفی (x) از بازه واحد $[0, 1]$



شکل ۴-۱: هیستوگرام ضبط شده از زمان رفت و آمد (بر حسب دقیقه) تا محل کار. مجموعه داده شامل ۳۴۰ اندازه گیری است که در طی یک سال، برای مسافتی تقریباً ۳/۱ مایلی جمع آوری شده است. داده‌ها با استفاده از یک خانواده گاما از توزیع‌های احتمال، با پارامترهای مکان و

جایی که $\Gamma(\alpha)$ تابع گاما نامیده می‌شود. یک متغیر تصادفی که گاما توزیع شده است به صورت $X \sim \text{Gamma}(\alpha, \beta)$ نشان داده می‌شود.

مثال ۵: انتخاب یک عدد (x) بین ۰ و ۱ به طور یکنواخت و به طور تصادفی را در نظر بگیرید (شکل ۳، ۱). احتمال اینکه عدد بزرگتر یا مساوی $\frac{3}{4}$ باشد یا کمتر مساوی $\frac{1}{4}$ چقدر است؟

ما می‌دانیم که $\Omega = [0, 1]$ توزیع با pdf یکنواخت تعریف می‌شود، $p(\omega) = \frac{1}{b-a} = 1$ که $a = 0$ ، $b = 1$ بازه فضای نتیجه را تعریف می‌کند. ما رویداد مورد علاقه را به صورت $A = [0, \frac{1}{4}] \cup [\frac{3}{4}, 1]$ تعریف می‌کنیم و احتمال آن را به این صورت محاسبه می‌کنیم

$$P(A) = \int_0^{\frac{1}{4}} p(\omega) d\omega + \int_{\frac{3}{4}}^1 p(\omega) d\omega = (\frac{1}{4} - 0) + (1 - \frac{3}{4}) = \frac{1}{2}$$

اگر در عوض احتمال اینکه عدد بزرگتر از $\frac{3}{4}$ یا کمتر از $\frac{1}{4}$ باشد را بررسی کنیم چه؟ از آنجایی که احتمال هر رویداد فردی در حالت پیوسته ۰ است، اگر بازه‌های باز یا بسته را در نظر بگیریم، تفاوتی در ادغام وجود ندارد. بنابراین، احتمال همچنان $\frac{1}{2}$ خواهد بود

مثال ۶: بیایید تصور کنیم زمان رفت و آمد خود را برای سال جمع‌آوری کرده‌اید و می‌خواهید احتمال زمان رفت و آمد خود را مدل کنید تا به شما کمک کند بتوانید زمان رفت و آمد فردا را پیش‌بینی کنید. برای این تنظیم، متغیر تصادفی X شما مطابق با زمان رفت و آمد است و شما باید احتمالاتی را برای این متغیر تصادفی تعریف کنید. این داده‌ها را می‌توان گسسته در نظر گرفت و مقادیر را در دقیقه، $\{4, 5, 6, \dots, 26\}$ می‌گیرد. سپس می‌توانید هیستوگرام‌هایی از این داده‌ها ایجاد کنید (جدول مقادیر احتمال)، همانطور که در شکل ۱، ۴ نشان داده شده است، تا احتمال زمان رفت و آمد را نشان دهید.

با این حال، زمان رفت و آمد در واقع مجزا نیست، بنابراین شما می‌خواهید به عنوان یک مدل پیوسته مدل‌سازی کنید. یک انتخاب معقول توزیع گاما است. با این حال، چگونه می‌توان داده‌های ثبت شده را گرفته و پارامترهای α ، β را در توزیع گاما تعیین کرد؟ تخمین این پارامترها در واقع کاملاً ساده است، اگرچه به اندازه تخمین جداول مقادیر احتمال آشکار نیست. نحوه

انجام این کار را در فصل ۳ مورد بحث قرار می‌دهیم. توزیع گامای آموخته شده نیز در شکل ۱,۴ نشان داده شده است. با توجه به توزیع گاما، اکنون می‌توان این سوال را مطرح کرد: محتمل‌ترین زمان رفت و آمد امروز چقدر است؟ این مربوط به $\max_{\omega} p(\omega)$ است که به آن حالت توزیع می‌گویند. سوال طبیعی دیگر میانگین یا زمان مورد انتظار رفت و آمد است. برای به دست آوردن این، به مقدار مورد انتظار (میانگین) این توزیع گاما نیاز دارید که در زیر در بخش ۱,۴ تعریف می‌کنیم.

۱-۳ متغیرهای تصادفی چند متغیره

توسعه بسیاری از مفاهیم فوق به متغیرهای تصادفی چند متغیره - بردار متغیرهای تصادفی - گسترش می‌یابد، زیرا تعریف فضاهای نتیجه و احتمالات عمومی است. با این حال، مثال‌هایی که تاکنون ارائه شده‌اند، با متغیرهای تصادفی اسکالر سروکار داشته‌اند، زیرا برای متغیرهای تصادفی چند متغیره، باید نحوه تعامل متغیرها را درک کنیم. در این بخش، به چندین مفهوم جدید می‌پردازیم که تنها زمانی به وجود می‌آیند که متغیرهای تصادفی متعددی از جمله توزیع‌های مشترک، توزیع‌های شرطی، حاشیه‌ها و وابستگی بین متغیرها وجود داشته باشد.

اجازه دهید با یک مثال ساده‌تر شروع کنیم، با دو متغیر تصادفی گسسته X و Y با فضاهای نتیجه X و Y . یک تابع جرم احتمال مشترک p وجود دارد: $[0, 1] \rightarrow X \times Y$ و توزیع احتمال مشترک P ، مانند که

$$p(x, y) \stackrel{\text{def}}{=} P(X = x, Y = y)$$

جایی که pmf باید انجام شود

$$\sum_{x \in X} \sum_{y \in Y} p(x, y) = 1$$

به عنوان مثال، اگر $X = \{\text{young, old}\}$ و $Y = \{\text{no arthritis, arthritis}\}$ آنگاه pmf می‌تواند جدول احتمالات مشترک باشد که در تعریف فضاهای احتمال مطابقت دارد، زیرا

		Y	
		0	1
X	0	$P(X=0, Y=0) = 1/2$	$P(X=0, Y=1) = 1/100$
	1	$P(X=1, Y=0) = 1/10$	$P(X=1, Y=1) = 39/100$

جدول ۱-۱: جدول احتمال مشترک برای متغیرهای تصادفی X و Y .

$\Omega = X \times Y$ یک فضای معتبر است و $\sum_{\omega \in \Omega} p(\omega) = \sum_{(x,y) \in \Omega} p(x, y) = \sum_{x \in X} \sum_{y \in Y} p(x, y)$ متغیر تصادفی $Z = (X, Y)$ یک متغیر تصادفی چند متغیره است که دارای دو بعد است.

با نگاه کردن به احتمالات مشترک در جدول، می‌توانیم ببینیم که دو متغیر تصادفی با هم تعامل دارند. به عنوان مثال، احتمال مفاصل برای افراد جوان و مبتلا به آرتریت کم است. علاوه بر این، به نظر می‌رسد بزرگی بیشتری در ردیف‌های مربوط به جوان

وجود دارد، که نشان می‌دهد احتمالات تحت تأثیر نسبت افراد پیر یا جوان در جمعیت است. در واقع، می‌توان پرسید که آیا می‌توانیم این نسبت را فقط از این جدول دریابیم.

پاسخ کاملاً مثبت است، و ما را به توزیع‌های حاشیه‌ای و اینکه چرا ممکن است به توزیع‌های حاشیه‌ای اهمیت دهیم، هدایت می‌کند. با توجه به توزیع مشترک بر روی متغیرهای تصادفی، می‌توان امیدوار بود که بتوانیم احتمالات خاص‌تری را استخراج کنیم، مانند توزیع فقط روی یکی از آن متغیرها، که توزیع حاشیه‌ای نام دارد. حاشیه را می‌توان به سادگی با جمع کردن تمام مقادیر متغیر دیگر محاسبه کرد

$$P(X = \text{young}) = p(\text{young, no arthritis}) + p(\text{young, arthritis}) = \frac{51}{100}$$

یک فرد جوان یا بیماری آرتریت دارد یا ندارد، بنابراین جمع‌بندی این دو مورد احتمالی آن متغیر را مشخص می‌کند. بنابراین، با استفاده از داده‌های جمع‌آوری شده برای متغیر تصادفی $Z = (X, Y)$ ، می‌توان نسبت جمعیت جوان و نسبت پیر را تعیین کرد.

به طور کلی، می‌توانیم متغیر تصادفی d بعدی $\mathbf{X} = (X_1, X_2, \dots, X_d)$ را با نتایج بردار $\mathbf{x} = (x_1, x_2, \dots, x_d)$ در نظر بگیریم، به طوری که هر x_i از بین برخی از اشیاء انتخاب شود سپس، برای حالت گسسته، هر تابع $p: X_1 \times X_2 \times \dots \times X_d \rightarrow [0, 1]$ تابع جرم احتمال چند بعدی نامیده می‌شود اگر

$$\sum_{x_1 \in X_1} \sum_{x_2 \in X_2} \dots \sum_{x_d \in X_d} p(x_1, x_2, \dots, x_d) = 1$$

یا، برای حالت پیوسته، $p: X_1 \times X_2 \times \dots \times X_d \rightarrow [1, \infty]$ یک تابع چگالی احتمال چند بعدی است اگر

$$\int_{x_1} \int_{x_2} \dots \int_{x_d} p(x_1, x_2, \dots, x_d) dx_1 dx_2 \dots dx_d = 1$$

یک توزیع حاشیه‌ای برای زیرمجموعه‌ای از $\mathbf{X} = (X_1, X_2, \dots, X_d)$ با جمع یا ادغام بر روی متغیرهای باقی مانده تعریف می‌شود. برای حالت گسسته، توزیع حاشیه‌ای $p(x_i)$ به این صورت تعریف شده است

$$p(x_i) \stackrel{\text{def}}{=} \sum_{x_1 \in X_1} \dots \sum_{x_{i-1} \in X_{i-1}} \sum_{x_{i+1} \in X_{i+1}} \dots \sum_{x_d \in X_d} p(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_d)$$

که در آن متغیر x_i روی مقداری ثابت است و همه مقادیر ممکن متغیرهای دیگر را جمع می‌کنیم. به طور مشابه، برای حالت پیوسته، توزیع حاشیه‌ای $p(x_i)$ به این صورت تعریف شده است

$$p(x_i) \stackrel{\text{def}}{=} \int_{x_1} \dots \int_{x_{i-1}} \int_{x_{i+1}} \dots \int_{x_d} p(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_d) dx_1 \dots dx_{i-1} dx_{i+1} \dots dx_d$$

توجه داشته باشید که ما از p برای تعریف چگالی روی \mathbf{X} استفاده می‌کنیم، اما سپس این اصطلاح را بارگذاری می‌کنیم و همچنین از p برای چگالی فقط روی x_i استفاده می‌کنیم. برای نتایج دقیق‌تر، باید دو تابع مجزا (pdf) تعریف کنیم، مثلاً $p_{\mathbf{X}}$ برای چگالی روی متغیر تصادفی چند متغیره و p_{x_i} برای حاشیه. با این حال، استفاده ساده از p و استنتاج متغیر تصادفی از متن معمول است. در بیشتر موارد، واضح است؛ اگر اینطور نیست، ما به صراحت pdfها را با زیرنویس‌های اضافی برجسته می‌کنیم.

ما می‌توانیم pmf و pdf های چند متغیره رایج را تعریف کنیم که پسوند pmf و pdf اسکالر هستند. برخی از پسوندها - مانند جداول مقادیر احتمال و توزیع‌های یکنواخت واضح تر هستند، در حالی که برخی دیگر نیاز به اندیس‌های بیشتر دارند، مانند گاوسی. برای دیگران، مانند لاپلاس، پسوند ممکن است منحصر به فرد نباشد و چندین گزینه امکان پذیر است. ما در بخش پایانی این فصل، بخش ۱.۵، برای مرجع، افزونه‌هایی را که به آن نیاز خواهیم داشت، تعریف می‌کنیم. با این حال، ابتدا درک چگونگی تعامل چندین متغیر مفید خواهد بود، زیرا این امر بر گسترش از تک متغیره به چند متغیره تأثیر می‌گذارد. به ویژه، درک توزیع‌های شرطی و وابستگی مفید خواهد بود، که در ادامه به آن می‌پردازیم.

۱-۳-۱ توزیع‌های مشروط

احتمالات شرطی احتمالات یک متغیر تصادفی مانند X ، اطلاعاتی در مورد مقدار متغیر تصادفی دیگر Y به ما میدهد. به طور رسمی‌تر، احتمال شرطی $p(y|x)$ برای دو متغیر تصادفی X و Y به صورت تعریف می‌شود.

$$p(y|x) \stackrel{\text{def}}{=} \frac{p(x,y)}{p(x)} \quad (1.1)$$

که $p(x) > 0$

تمرین ۲: بررسی کنید که $p(y|x)$ بر روی تمام مقادیر $y \in Y$ برای یک $x \in X$ معین ثابت، با 1 ادغام می‌شود، و بنابراین شرایط یک تابع جرم احتمالی (چگالی) را برآورده می‌کند.

معادله (۱,۱) اکنون به ما امکان می‌دهد تا با توجه به مشاهدات x ، احتمال پیشین یک رویداد A را محاسبه کنیم.

$$P(Y \in A | X = x) = \begin{cases} \sum_{y \in A} p(y|x) & Y: \text{discrete} \\ \int_A p(y|x) dy & Y: \text{continuous} \end{cases}$$

نوشتن $p(x,y) = p(x|y)p(y) = p(y|x)p(x)$ قاعده ضرب نامیده می‌شود. گسترش بیش از دو متغیر ساده است. ما میتوانیم بنویسیم

$$p(x_1, \dots, x_d) = p(x_d | x_1, \dots, x_{d-1}) p(x_1, \dots, x_{d-1})$$

با استفاده تابع بازگشتی از قاعده ضرب، به دست می‌آوریم

$$\begin{aligned} p(x_1, \dots, x_d) &= p(x_d | x_1, \dots, x_{d-1}) p(x_1, \dots, x_{d-1}) \\ &= p(x_d | x_1, \dots, x_{d-1}) p(x_{d-1} | x_1, \dots, x_{d-2}) p(x_1, \dots, x_{d-2}) \\ &\quad \cdot \\ &\quad \cdot \\ &\quad \cdot \end{aligned}$$

$$p(x_d | x_1, \dots, x_{d-1}) p(x_{d-1} | x_1, \dots, x_{d-2}) \dots p(x_2 | x_1) p(x_1)$$

به طور فشرده تر

$$p(x_1, \dots, x_d) = p(x_1) \prod_{i=2}^d p(x_i | x_1, \dots, x_{i-1}) \quad (1.2)$$

که از آن به عنوان قانون زنجیره‌ای یا قاعده کلی ضرب یاد می‌شود. به عنوان مثال، برای سه متغیر، قاعده ضرب می‌دهد

$$p(x_1, x_2, x_3) = p(x_3 | x_2, x_1) p(x_2 | x_1) p(x_1)$$

این قاعده برای مجموعه‌ای از متغیرهای تصادفی نیز اعمال می‌شود، جایی که یک مجموعه را می‌توان به عنوان یک متغیر تصادفی در نظر گرفت. مثلاً،

$$p(x_1, x_2, x_3) = p(x_3 | x_2, x_1) p(x_2 | x_1) p(x_1)$$

این به این دلیل است که (x_2, x_3) یک فضای احتمال معتبر دارند، بنابراین می‌توانیم از قاعده ضرب برای دو متغیر استفاده کنیم: x_1 و (x_2, x_3) . با استفاده از قاعده ضرب، با دادن $p(x, y) = p(x|y)p(y) = p(y|x)p(x)$ می‌توانیم قانون بیز را نیز استخراج کنیم:

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)} \quad (1.3)$$

بنابراین، واقعاً فقط باید قاعده ضرب را به خاطر بسپارید تا به راحتی قانون بیز را به خاطر بیاورید

ممکن است متوجه شوید که ترتیب متغیرها در قاعده ضرب اهمیتی ندارد. در واقع تا حدودی جالب است که ما می‌توانیم توزیع شرطی $p(x|y)$ و $p(y)$ حاشیه‌ای را تعریف کنیم یا می‌توانیم $p(y|x)$ و $p(x)$ را تعریف کنیم و هر دو به طور معادل توزیع مشترک $p(x, y)$ این ویژگی به سادگی یک واقعیت از تعریف توزیع‌های شرطی است و در هنگام تخمین توزیع‌ها انعطاف‌پذیری را فراهم می‌کند. ما بیشتر از این هم‌ارزی در قالب قانون بیز، هنگام انجام تخمین پارامتر و ماکسیمم احتمال استفاده خواهیم کرد. برای کار در مدل‌های گرافیکی، که در اینجا مورد بحث قرار نمی‌گیرد، این انعطاف‌پذیری از اهمیت بیشتری برخوردار است.

۲-۳-۱ متغیرهای تصادفی مستقل

دو متغیر تصادفی مستقل هستند اگر عوامل توزیع احتمال مشترک آنها، حاصل ضرب حاشیه‌ها باشد

$$p(x, y) = p(x)p(y)$$

یکی از دلایل شهودی این تعریف را می‌توان با در نظر گرفتن X به شرط Y مشاهده کرد. اگر $p(x|y) = p(x)$ این بدان معناست که مقدار Y هیچ تأثیری بر توزیع روی X ندارد و بنابراین آنها مستقل هستند. از قاعده حاصلضرب، $p(x, y) = p(x|y)p(y)$ را می‌دانیم و از آنجایی که $p(x, y) = p(x)p(y)$ را بدست می‌آوریم. $p(y)$ همانطور که در بالا تعریف شد.

مفهوم استقلال را می‌توان به بیش از دو متغیر تصادفی تعمیم داد. به طور کلی‌تر، اگر بتوان توزیع احتمال مشترک هر زیرمجموعه‌ای از متغیرها را به عنوان حاصلضرب توزیع‌های احتمال حاشیه‌ای اجزای آن بیان کرد، به d متغیر تصادفی مستقل یا مشترکاً مستقل گفته می‌شود.

$$p(x_1, x_2, \dots, x_d) = p(x_1)p(x_2) \dots p(x_d)$$

شکل دیگری از استقلال، به نام استقلال شرطی، حتی بیشتر در یادگیری ماشین استفاده می‌شود. این نشان دهنده استقلال بین متغیرها در حضور برخی متغیرهای تصادفی دیگر (شواهد) است. به عنوان مثال،

$$p(x, y|z) = p(x|z)p(y|z)$$

جالب اینجاست که این دو شکل استقلال با هم ارتباطی ندارند: هیچ کدام متضمن دیگری نیست. X و Y می‌توانند مستقل باشند، اما نه به طور مشروط مستقل با توجه به Z . X و Y می‌توانند به صورت شرطی مستقل باشند، اما مستقل نیستند. ما این را در دو مثال ساده در شکل A.1 در پیوست نشان می‌دهیم.

در اینجا، مثالی ارائه می‌دهیم که مستقیماً با یادگیری ماشین مرتبط است، در مورد اینکه چرا به استقلال و استقلال مشروط اهمیت می‌دهیم. اگر دو متغیر مستقل باشند، این مفاهیم مدل‌سازی مهمی دارد. برای مثال، اگر ویژگی X و هدف Y مستقل باشند، X برای پیش‌بینی Y مفید نیست و بنابراین ویژگی مفیدی نیست. اگر دو متغیر با توجه به متغیر دیگری به صورت شرطی مستقل باشند، این نیز می‌تواند مفاهیم مدل‌سازی مهمی داشته باشد. به عنوان مثال، اگر ما دو ویژگی X_1 و X_2 با هدف Y داشته باشیم، که در آن X_2 و Y به طور مشروط مستقل از X_1 هستند، ویژگی X_2 اضافی است و می‌تواند به طور بالقوه کنار گذاشته شود.

به عنوان یک مثال عینی، اجازه دهید $X_1 = \text{temperature in Celcius}$ و $X_2 = \text{temperature in Fahrenheit}$ ، $Y = \text{plants need watering}$ قطعاً مستقل از X_2 نیست. با این حال، هنگامی که X_1 شناخته شد (یا داده شد)، دیگر اطلاعات اضافی از X_2 به دست نمی‌آید و بنابراین $p(y|x_1, x_2) = p(y|x_1) = p(y|x_2)$. به طور کلی، شناخت استقلال‌ها و استقلال‌های مشروط می‌تواند روند مدل‌سازی را اطلاع‌رسانی و ساده‌سازی کند و نمونه‌های متعددی را از نظر ساده‌سازی ماکسیمم احتمال برای متغیرهای تصادفی مستقل و با توزیع یکسان¹ یا به اختصار i.i.d. را برای داده‌ها و در بیز ساده برای طبقه بندی خواهیم دید. ما این بخش را با یک مثال دیگر، با استفاده از یک سکه مغرضانه، به پایان می‌بریم تا تمایز بین استقلال و استقلال مشروط را برجسته کنیم.

مثال ۷: [سکه مغرضانه و استقلال مشروط] فرض کنید یک تولیدکننده یک سکه مغرضانه تولید کرده است، جایی که به طور تصادفی سر (H) یا دم (T) را نمی‌دهد. در عوض، در واقع مقداری احتمال ناشناخته α برای دیدن H در هنگام چرخاندن سکه را دارد. از آنجایی که این سوگیری ناشناخته است، ما عدم قطعیت خود را با تعریف یک تصادفی (بایاس سکه) رمزگذاری می‌کنیم. به طور کلی، این متغیر تصادفی می‌تواند مقادیر $[0, 1]$ را بگیرد. برای اهداف این مثال، اجازه دهید این را کمی ساده‌تر کنیم و فرض کنیم که می‌دانیم بایاس یکی از $Z = \{0.1, 0.5, 0.8\}$ است. اگر سوگیری 0.5 باشد، به این معنی است که این یک سکه بی طرفانه (منصفانه) است. اجازه دهید فرض کنیم که احتمال هر سوگیری به یک اندازه محتمل است، یعنی $P(Z = z) = \frac{1}{3}$ ، زیرا سازنده هیچ دلیلی به ما نداده است که هر یک از 0.1، 0.5 یا 0.8 را محتمل‌تر بدانیم. حال تصور کنید که سکه را دو بار برگردانید و دو نتیجه X_1 و X_2 را ثبت کنید. این دو تلنگر جداگانه با دو متغیر تصادفی X_1 و X_2 مطابقت دارند. فضای نتیجه برای X_1 و X_2 ، $\{H, T\}$ است. با توجه به سوگیری واقعی سکه، α ، توزیع واقعی برنولی $P(X_i = H|Z = \alpha) = \alpha$ است. با این حال، ما سوگیری α را نمی‌دانیم. در عوض، ما آن را با یک متغیر تصادفی Z مدل می‌کنیم، به طوری که برای هر $Z = z$ داده شده، می‌دانیم که $P(X_i|Z = z) = z$ ، یعنی X_i یک متغیر تصادفی

¹ Independent and identically distributed random variables

برنولی با پارامتر Z است. از آنجایی که ما سوگیری واقعی را نمی‌دانیم، باید روی Z به حاشیه برسیم تا توزیع حاشیه ای را روی X_1 بدست آوریم.

$$\begin{aligned} P(X_1 = x) &= \sum_{z \in \mathcal{Z}} P(X_1 = x, Z = z) \\ &= \sum_{z \in \mathcal{Z}} P(X_1 = x | Z = z) P(Z = z) \\ &= P(X_1 = x | Z = 0.1) P(Z = 0.1) + P(X_1 = x | Z = 0.5) P(Z = 0.5) + P(X_1 = x | Z = 0.8) P(Z = 0.8) \end{aligned}$$

آیا X_1 و X_2 به صورت مشروط مستقل از Z هستند؟ پاسخ مثبت است، زیرا با توجه به سوگیری سکه، دانستن نتیجه X_2 بر توزیع بر روی X_1 تأثیری ندارد، به عنوان مثال،

$$P(X_1 = x_1, X_2 = x_2 | Z = z) = P(X_1 = x_1 | Z = z) P(X_2 = x_2 | Z = z)$$

صرف نظر از آنچه برای X_2 مشاهده می‌کنیم، می‌دانیم که توزیع بر روی X_1 برنولی با بایاس Z داده شده است.

آیا X_1 و X_2 مستقل هستند؟ پاسخ منفی است، زیرا بدون دانستن سوگیری سکه، دانستن نتیجه X_2 چیزی در مورد توزیع برنولی بر روی X_1 به ما می‌گوید. برای مثال، اگر $X_1 = T$ و $X_2 = H$ ، نتیجه دوم نشان می‌دهد که سوگیری ممکن است کاملاً به سمت T منحرف نشود.

$$\begin{aligned} P(X_1 = x_1, X_2 = x_2) &= \sum_{z \in \mathcal{Z}} P(X_1 = x_1, X_2 = x_2 | Z = z) P(Z = z) \\ &= \sum_{z \in \mathcal{Z}} P(X_1 = x_1 | Z = z) P(X_2 = x_2 | Z = z) P(Z = z) \end{aligned}$$

که برابر با $P(X_1 = x_1) P(X_2 = x_2)$ تضمین نشده است، که در آن

$$\begin{aligned} &P(X_1 = x_1) P(X_2 = x_2) \\ &= \left(\sum_{z_1 \in \mathcal{Z}} P(X_1 = x_1 | Z = z_1) P(Z = z_1) \right) \left(\sum_{z_2 \in \mathcal{Z}} P(X_2 = x_2 | Z = z_2) P(Z = z_2) \right) \end{aligned}$$

۱-۴ امیدهای ریاضی و گشتاور

مقدار مورد امید ریاضی، یا میانگین، متغیر تصادفی X ، میانگین X نمونه برداری مکرر در محدوده نمونه گیری است. این لزوماً مقداری نیست که ما امید داریم اغلب آن را ببینیم - که حالت نامیده می‌شود. به طور دقیق تر، با توجه به pmf p یا pdf p برای فضای نتیجه X ، امید X است

$$\mathbb{E}[X] \stackrel{\text{def}}{=} \begin{cases} \sum_{x \in \mathcal{X}} xp(x) & \text{if } X \text{ is discrete} \\ \int_{\mathcal{X}} xp(x)dx & \text{if } X \text{ is continuous} \end{cases}$$

برای تاس انداختن، که در آن هر عدد از 1 تا 6 دارای احتمال یکنواخت است، مقدار مورد انتظار 3.5 است و حالت برای همه اعداد گره خورده است (یعنی چند وجهی است). برای توزیع برنولی، که در آن $X = \{0, 1\}$ مقدار مورد انتظار α است، که حتی نتیجه‌ای نیست که مشاهده شود، اما اگر سکه را بی‌نهایت بار برگردانیم، میانگین 0 و 1 است. حالت در این مورد به α بستگی دارد: اگر $\alpha > 0.5$ باشد، احتمال 1 بیشتر است، آنگاه حالت 1 است. اگر $\alpha < 0.5$ باشد، حالت 0 است. در غیر این صورت، دو وجهی با مدهای 0 و 1 است. برای توزیع گاوسی، مقدار مورد انتظار پارامتر μ است و مد نیز برابر μ است. به طور کلی، ممکن است به مقدار مورد انتظار توابع متغیر تصادفی X علاقه‌مند باشیم. برای مثال، ممکن است بخواهیم $E[X^2]$ یا به طور کلی $E[X^k]$ را برای مقداری $k > 1$ بدانیم. یا ممکن است بخواهیم بدانیم $E[(X - c)^k]$ برای مقداری $k > 1$ و ثابت c . به این گشتاور X می‌گویند. به طور کلی، برای تابع $f: X \rightarrow \mathbb{R}$ ، می‌توانیم $f(X)$ را یک متغیر تصادفی تبدیل شده در نظر بگیریم و امید آن را به صورت تعریف کنیم.

$$\mathbb{E}[f(X)] = \mathbb{E}[f(X)] = \begin{cases} \sum_{x \in \mathcal{X}} f(x)p(x) & \text{if } X \text{ is discrete} \\ \int_{\mathcal{X}} f(x)p(x)dx & \text{if } X \text{ is continuous} \end{cases}$$

اگر $\mathbb{E}[f(X)] = \pm\infty$ می‌گوییم که امید ریاضی وجود ندارد یا به خوبی تعریف نشده است.

یک گشتاور مفید، واریانس است: گشتاور دوم گرایش به مرکز، که در آن گرایش به مرکز نشان دهنده $\mathbb{E}[X] = c$ است. واریانس مقداری را نشان می‌دهد که متغیر تصادفی حول میانگین آن تغییر می‌کند. به عنوان مثال، برای توزیع گاوسی، اگر واریانس σ^2 بزرگ باشد، گاوس بسیار گسترده است، که نشان دهنده چگالی غیر قابل اغماض برای محدوده وسیع‌تری از نقاط x در اطراف μ است. متناوباً، اگر σ^2 تقریباً صفر باشد، گاوسی به شدت در اطراف μ متمرکز می‌شود. همچنین می‌توانیم امیدهای شرطی و امیدها را برای متغیرهای تصادفی چند متغیره در نظر بگیریم. برای دو متغیر تصادفی X و Y و تابع $f: Y \rightarrow \mathbb{R}$ ، امید شرطی است

$$\mathbb{E}[f(Y)|X = x] = \begin{cases} \sum_{y \in \mathcal{Y}} f(y)p(y|x) & \text{if } Y \text{ is discrete} \\ \int_{\mathcal{Y}} f(y)p(y|x)dy & \text{if } Y \text{ is continuous} \end{cases}$$

استفاده از تابع همانی $f(y) = y$ به امید شرطی استاندارد $\mathbb{E}[Y|x]$ منجر می‌شود.

تمرین ۳: قانون کل امید ریاضی را نشان دهید: $\mathbb{E}[Y] = \mathbb{E}[\mathbb{E}[Y|X]]$ ، که در آن امید بیرونی بیش از X و امید درونی بیش از Y است. برای مثال، اگر X و Y هر دو گسسته باشند

$$\begin{aligned} \mathbb{E}[\mathbb{E}[Y|X]] &= \sum_{x \in \mathcal{X}} p(x)\mathbb{E}[Y|X = x] \\ &= \sum_{x \in \mathcal{X}} p(x) \sum_{y \in \mathcal{Y}} yp(y|x) \end{aligned}$$

برای دو متغیر تصادفی X و Y و $f: X \times Y \rightarrow \mathbb{R}$ ، همچنین می‌توانیم امید را بر روی توزیع مشترک تعریف کنیم، با یک متغیر ثابت

$$\mathbb{E}[f(X, y)] = \begin{cases} \sum_{x \in \mathcal{X}} f(x, y) p(x|y) & \text{if } X \text{ is discrete} \\ \int_{\mathcal{X}} f(x, y) p(x|y) dx & \text{if } X \text{ is continuous} \end{cases}$$

یا بیش از هر دو متغیر

$$\mathbb{E}[f(X, Y)] = \begin{cases} \sum_{y \in \mathcal{Y}} p(y) \mathbb{E}[f(X, y)] & \text{if } Y \text{ is discrete} \\ \int_{\mathcal{Y}} p(y) \mathbb{E}[f(X, y)] dy & \text{if } Y \text{ is continuous} \end{cases}$$

به عنوان مثال، اگر X پیوسته و Y گسسته باشد، این نشان می‌دهد

$$\mathbb{E}[f(X, Y)] = \sum_{y \in \mathcal{Y}} p(y) \mathbb{E}[f(X, y)]$$

تمرین ۴: نشان دهید $\int_{\mathcal{X}} (\sum_{y \in \mathcal{Y}} \{f(x, y) p(x, y)\}) dx = \int_{\mathcal{X}} \mathbb{E}[f(x, Y)] p(x) dx$

همانطور که در بالا در مورد واریانس، کوواریانس یک نمونه مهم از این مقادیر مورد انتظار است، با $f(x, y) = (x - \mathbb{E}[X])(y - \mathbb{E}[Y])$ مقدار مورد انتظار در این تابع نشان می‌دهد که چگونه دو متغیر با هم متفاوت هستند. ما از علامت گذاری خاص برای کوواریانس استفاده می‌کنیم، که اغلب استفاده می‌شود

$$\begin{aligned} \text{Cov}[X, Y] &= \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] \\ &= \mathbb{E}[XY] - \mathbb{E}[X] \mathbb{E}[Y] \end{aligned}$$

با $\text{Cov}[X, X] = V[X]$ واریانس متغیر تصادفی X است. همبستگی کوواریانس است که با انحراف استاندارد - ریشه دوم واریانس - هر متغیر تصادفی نرمال شده است.

$$\text{Corr}[X, Y] = \frac{\text{Cov}[X, Y]}{\sqrt{V[X]} \sqrt{V[Y]}}$$

اگر X و Y خود واریانس زیادی داشته باشند، کوواریانس می‌تواند بزرگتر شود. از سوی دیگر، همبستگی بین 1- و 1 تضمین شده است، و به همین ترتیب یک معیار متغیر مقیاس از نحوه تغییر متغیرها با هم است. در بسیاری از شرایط ما نیاز به تجزیه و تحلیل بیش از دو متغیر تصادفی داریم. یک خلاصه دو بعدی ساده از تمام مقادیر کوواریانس زوجی شامل d متغیرهای تصادفی X_1, X_2, \dots, X_d را ماتریس کوواریانس می‌گویند. ماتریس کوواریانس $\Sigma \in \mathbb{R}^{d \times d}$ دارای ورودی (i, j) است که به صورت تعریف شده است.

$$\begin{aligned} \Sigma_{ij} &= \text{Cov}[X_i, X_j] \\ &= \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])] \end{aligned}$$

با ماتریس کامل به صورت نوشته شده است

$$\begin{aligned} \Sigma &= \text{Cov}[\mathbf{X}, \mathbf{X}] \\ &= \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^T] \end{aligned}$$

$$= E[XX^T] - E[X]E[X]^T$$

خط دوم شامل حاصلضرب بیرونی بردار $v = X - E[X] \in \mathbb{R}^d$ است تا $A = vv^T \in \mathbb{R}^{d \times d}$ را تولید کند. این حاصل ضرب بیرونی یک ضرب ماتریس است که ماتریس اول $d \times 1$ و دومی $1 \times d$ است. با استفاده از قوانین ضرب ماتریس، $A_j = v_i v_j$ به دست می‌آید. توجه داشته باشید که عناصر قطری ماتریس کوواریانس $d \times d$ واریانس برای هر متغیر X_i و عناصر خارج از مورب مقادیر کوواریانس بین جفت متغیرها هستند. ماتریس کوواریانس متقارن و مثبت نیمه معین است. به یاد بیاورید که اگر $z^T > \Sigma z \geq 0$ برای همه بردارها $z \neq 0$ یک ماتریس نیمه معین مثبت (که با $\Sigma \geq 0$ نشان داده می‌شود) گفته می‌شود. به طور معادل، مقادیر ویژه ماتریس همگی بزرگتر یا مساوی صفر هستند. اگر ماتریس مثبت قطعی به جای مثبت نیمه قطعی ($\Sigma \geq 0$) باشد، آنگاه مقادیر ویژه کاملاً مثبت هستند و ماتریس کوواریانس رتبه کاملی دارد. یک ماتریس نیمه معین مثبت می‌تواند مقادیر ویژه ای داشته باشد که صفر هستند و بنابراین دارای رتبه کوچکتر از d هستند. برای $d = 1$ ، این خاصیت نیمه معین مثبت ماتریس با خاصیت $\sigma^2 \geq 0$ برای واریانس اسکالر مطابقت دارد.

خواص امید ریاضی

در اینجا به بررسی برخی از خواص مفید انتظارات می‌پردازیم. متغیرهای تصادفی چند متغیره، $X \in \mathbb{R}^d$ و $Y \in \mathbb{R}^m$ را برای $d, m \in N$ ، با متغیرهای تصادفی تک متغیره به عنوان یک مورد خاص در نظر بگیرید. برای ثابت $c \in \mathbb{R}$ ، به این صورت است که:

$$1. E[cX] = cE[X] \in \mathbb{R}^d$$

$$2. E[X + Y] = E[X] + E[Y] \quad \text{وقتی که } d = m$$

$$3. V[c] = 0 \quad \text{واریانس یک ثابت صفر است}$$

$$4. V[X] \stackrel{\text{def}}{=} \text{Cov}[X, X] \geq 0 \quad \text{که در آن برای } d = 1, V[X] \geq 0 \text{ یک اسکالر است.}$$

ما از $V[X]$ به عنوان مخفف $\text{Cov}[X, X]$ استفاده می‌کنیم.

$$5. V[cX] = c^2 V[X] \in \mathbb{R}^{d \times d}$$

$$6. \text{Cov}[X, Y] = E[(X - E[X])(Y - E[Y])^T] = E[XY^T] - E[X]E[Y]^T \in \mathbb{R}^{d \times m}$$

$$7. V[X + Y] = V[X] + V[Y] + 2\text{Cov}[X, Y] \quad \text{وقتی که } d = m$$

علاوه بر این، اگر X و Y متغیرهای تصادفی مستقل باشند، چنین است که:

$$8. E[X_i Y_j] = E[X_i] E[Y_j] \quad \text{برای هر } i, j$$

$$9. V[X + Y] = V[X] + V[Y] \quad \text{وقتی که } d = m$$

$$10. \text{Cov}[X, Y] = 0$$

در نهایت، برای هر متغیر تصادفی d بعدی، $X_1 + X_2 + \dots + X_m$

$$11. \text{Cov}[X_1 + X_2 + \dots + X_m] = \sum_{i=1}^m \sum_{j=1}^m \text{Cov}[X_i, X_j] = \sum_{i=1}^m V[X_i] + 2 \sum_{1 \leq i < j \leq m} \text{Cov}[X_i, X_j]$$

۵-۱ چند متغیره PDF و PMF

ما اکنون افزونه‌هایی را برای تعاریف pmf و pdf برای دسته بندی چند متغیره‌ها در نظر می‌گیریم. توزیع گاوسی چند متغیره تعمیم توزیع گاوسی یا نرمال به حالت d-بعدی با $\Omega = \mathbb{R}^d$ است. به عنوان تعریف شده است

$$p(\omega) \stackrel{\text{def}}{=} \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2} (\omega - \mu)^T \Sigma^{-1} (\omega - \mu)\right)$$

با پارامترهای $\mu \in \mathbb{R}^d$ و ماتریس مثبت-معین Σ که ماتریس کوواریانس است. این تعریف نحوه تغییر متغیره‌ها با هم را در نظر می‌گیرد که توسط ماتریس کوواریانس Σ ارائه شده است. به عنوان مثال، اگر متغیره‌ها مستقل باشند، ماتریس کوواریانس مورب است. علاوه بر این، اگر هر متغیر دارای واریانس واحد باشد، گاوسی کروی است. اگر برخی از ابعاد واریانس بالاتری داشته باشند، گاوسی بیضی شکل است. اگر متغیره‌ها مستقل نباشند، گاوسی برش می‌شود و از میانگین آن به طور متفاوتی منحرف می‌شود که صرفاً با واریانس متغیره‌ها قابل محاسبه است. ما به این توزیع به عنوان $N(\mu, \Sigma)$ اشاره می‌کنیم و گاهی اوقات از $N(\omega|\mu, \Sigma)$ برای نشان دادن pdf برای μ و Σ داده شده استفاده می‌کنیم.

چند اصطلاح جدید در این pdf وجود دارد، بنابراین اجازه دهید این فرمول را بررسی کنیم. اصطلاح $|\Sigma|$ تعیین کننده Σ است. دترمینان برابر حاصل ضرب مقادیر ویژه Σ است. عرض توزیع را منعکس می‌کند. هنگامی که $d = 1$ ، تعیین کننده به سادگی برابر با σ^2 است، و ما طبق معمول با واریانس توزیع مقیاس می‌کنیم. به عنوان مثال دیگر، موردی را در نظر بگیرید که در آن کوواریانس بین همه متغیره‌ها صفر است. تصور کنید یک تابع گاوسی همه یک‌ها را در مورب برای کوواریانس داشته باشد، و دیگری دارای واریانس 5 برای اولین عنصر قطری و 1s در غیر این صورت. تعیین کننده اولی 1 و دومی 5 است. منطقی است که برای گاوسی دوم که واریانس بیشتری در یک بعد دارد، مقیاس بیشتری داشته باشیم. دومین عبارت جدید $(\omega - \mu)^T \Sigma^{-1} (\omega - \mu)$ است. این در واقع با یک محصول نقطه وزنی مطابقت دارد. دوباره یک ماتریس کوواریانس مورب را در نظر بگیرید، با $d = 3$ و 0.5، 2 در قطر. سپس، برای $z = \omega - \mu$

$$z^T \Sigma^{-1} z = \frac{z_1^2}{5} + \frac{z_2^2}{0.5} + \frac{z_3^2}{2}$$

برای متغیر تصادفی گسسته، تعمیم به ابعاد چندگانه برای pmf مستقیم است: pmf چند بعدی به سادگی با جداول احتمال چند بعدی مطابقت دارد. نمونه ای از آن را در جدول ۱،۱ دیدیم. با این حال، مانند حالت تک متغیره، چند pmf با نام وجود دارد، زیرا اغلب از آنها استفاده می‌شود.

یک مثال از یک pmf چند بعدی نامگذاری شده، توزیع طبقه‌ای است که نمونه‌ای از توزیع چندجمله‌ای است. توزیع طبقه‌ای برای مدل سازی یک متغیر تصادفی d-بعدی استفاده می‌شود که در آن هر عنصر می‌تواند 0 یا 1 باشد. این توزیع می‌تواند به طور معادل برای مدل سازی یک متغیر تصادفی اسکالر با d نتایج ممکن استفاده شود. با این حال، وقتی از آن در بخش ۸،۳ استفاده می‌کنیم، مفید خواهد بود که توزیع طبقه ای را به عنوان یک pmf برای یک متغیر تصادفی d بعدی در نظر بگیریم. هر نقطه (k_1, k_2, \dots, k_d) در فضای نتیجه یک بردار باینری است که دقیقاً یک عنصر $k_i = 1$ و بقیه صفر دارد، که نشان می‌دهد که نتیجه i رخ داده است. pmf طبقه بندی شده به این صورت تعریف می‌شود

$$p(k_1, k_2, \dots, k_d) \stackrel{\text{def}}{=} \begin{cases} \alpha_1^{k_1} \alpha_2^{k_2} \alpha_d^{k_d} & \text{if } k_1 + k_2 + \dots + k_d = n \\ 0 & \text{otherwise} \end{cases}$$

که در آن α_i ها ضرایب مثبت هستند به طوری که $\sum_{i=1}^d \alpha_i = 1$ یعنی هر ضریب α_i احتمال نتیجه i را می‌دهد.

مثال ۸: [مشتق ابعاد برای متغیرهای تصادفی گسسته چند بعدی] یکی از راه‌های اجتناب از pdf گسسته کردن متغیرها و سپس تعریف یک pmf است که جدولی از مقادیر احتمال است. اگرچه در برخی موارد معقول است، اما به طور کلی این می‌تواند به طور تصاعدی تعداد پارامترهای توزیع احتمال را افزایش دهد و نمونه ای از مشتق ابعاد است.

برای اینکه بفهمید چرا، مثال زیر را در نظر بگیرید. فرض کنید یک متغیر تصادفی d بعدی داریم که هر ورودی مقادیری بین 0 و 1 دارد (یعنی $X_i = [0, 1]$). شما تصمیم می‌گیرید این را گسسته کنید تا هر ورودی در یکی از سه سطح قرار گیرد و $X_i = \{1, 2, 3\}$ را تغییر دهید. اکنون شما انعطاف زیادی در تعیین این احتمالات در pmf خود دارید، برخلاف گوسی که فرم عملکردی سخت تری دارد. با این حال، متأسفانه، این جدول مقادیر می‌تواند بسیار بزرگ باشد، با ورودی‌های 3^d . برای این کار باید مقادیر احتمالی $3^d - 1$ را مشخص کنید، جایی که یکی از مقادیر به طور خودکار روی یک منهای مجموع همه احتمالات دیگر تنظیم می‌شود تا اطمینان حاصل شود که یک pmf معتبر دارید.

اگر در عوض، یک توزیع گاوسی روی این متغیرها مشخص کرده بودید، تعداد پارامترهای مورد نیاز برای تعریف توزیع فقط $d + d^2$ است که احتمالاً بسیار کمتر است.

فصل ۲

مقدمه‌ای بر بهینه‌سازی

بسیاری از مسائل یادگیری ماشین با یافتن تابع بهینه مطابق با یک هدف، با توابع یادگیری سروکار دارند. به عنوان مثال، ممکن است کسی علاقه‌مند به یافتن تابعی باشد $f: \mathbb{R}^d \rightarrow \mathbb{R}$ که تفاوت‌های مجذور برخی از اهداف را برای همه نمونه‌ها به مقدار مینیمم می‌رساند: $\sum_{i=1}^n (f(x_i) - y_i)^2$. برای یافتن چنین تابعی، باید درک اولیه‌ای از تکنیک‌های بهینه‌سازی داشته باشید. در این فصل، ابزارهای بهینه‌سازی اساسی را برای اهداف هموار عمومی مورد بحث قرار می‌دهیم. بسیاری از الگوریتم‌ها در یادگیری ماشین بر یک رویکرد ساده تکیه دارند: گرادین کاهشی. ابتدا در مورد چگونگی به مقدار مینیمم رساندن اهداف با استفاده از گرادین کاهشی مرتبه اول و دوم بحث می‌کنیم. این نمای کلی تنها بخش کوچکی از بهینه‌سازی را پوشش می‌دهد، اما خوشبختانه، بسیاری از الگوریتم‌های یادگیری ماشین مبتنی بر این رویکردهای بهینه‌سازی ساده هستند. ما بعداً، در فصل ۶، پس‌زمینه‌های بهینه‌سازی بیشتری را ارائه می‌کنیم، زمانی که فرصت استفاده از این دانش اولیه بهینه‌سازی را در دو فصل بعدی داشته باشید.

۲-۱ مسئله بهینه‌سازی اساسی و نقاط ثابت

یک هدف اصلی بهینه‌سازی انتخاب مجموعه‌ای از پارامترهای $w \in \mathbb{R}^d$ برای به مقدار مینیمم رساندن تابع هدف داده شده $c: \mathbb{R}^d \rightarrow \mathbb{R}$ است.

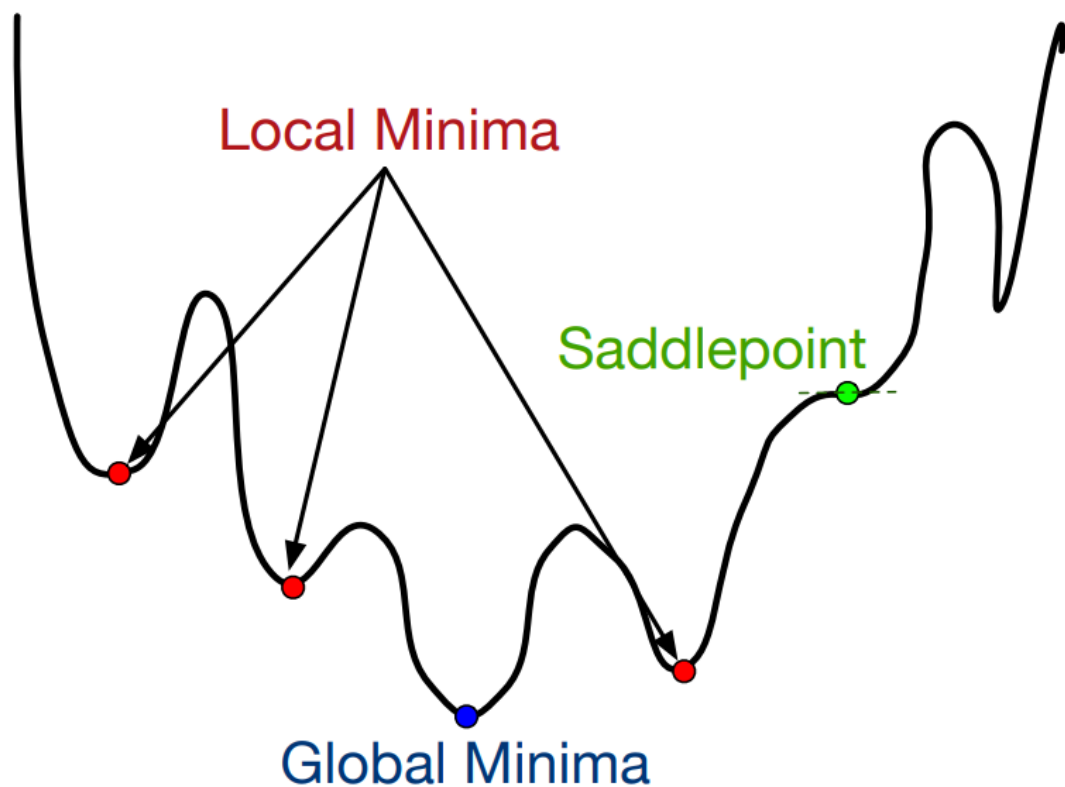
$$\min_{w \in \mathbb{R}^d} c(w)$$

به عنوان مثال، برای به دست آوردن پارامترهای w برای رگرسیون خطی که مجموع مجذور اختلافات را به مقدار مینیمم می‌رساند، از $c(w) = \sum_{i=1}^n (x_i, w_i - y_i)^2$ برای حاصل ضرب نقطه استفاده می‌کنیم.

$$\langle x_i, w_i \rangle = \sum_{j=1}^d x_{ij} w_j$$

ما در اینجا به جای خطا از اصطلاح هدف استفاده می‌کنیم، زیرا خطا مفهوم صریحی دارد که مفهوم آن تابع نادرست است. بعداً خواهیم دید که اهداف شامل هر دو معنای عبارت خطا می‌شوند - که نشان می‌دهد با چه دقتی داده‌ها را بازآفرینی می‌کنند - و همچنین عبارتهایی که اولویت‌های دیگر را در عملکرد ارائه می‌دهند. ترکیب این عبارات با خطا هدف نهایی را ایجاد می‌کند که مایلیم آن را به مقدار مینیمم برسانیم. به عنوان مثال، برای رگرسیون خطی، ما یک هدف منظم را بهینه می‌کنیم، $c(w) =$

پس هدف یافتن w است که آن را به مقدار مینیمم برساند. ساده ترین راه حل می تواند انجام یک جستجوی تصادفی باشد: w تصادفی تولید کنید و $c(w)$ را بررسی کنید. اگر هر w_t جدید تولید شده در تکرار t عملکرد بهتری از بهترین راه حل قبلی w داشته باشد، که در آن $c(w_t) < c(w)$ آنگاه می توانیم w_t را به عنوان راه حل بهینه جدید تنظیم کنیم. ما فرض می کنیم که اهداف ما پیوسته هستند و بنابراین می توانیم از این مزیت برای طراحی استراتژی های جستجوی بهتر استفاده کنیم. به طور خاص، برای توابع هموار، ما قادر خواهیم بود از گرادیان کاهش استفاده کنیم که در قسمت بعدی توضیح می دهیم.



شکل ۱-۲: نقاط ثابت روی یک سطح عملکرد صاف: مینیمم های محلی، مینیمم های مطلق و نقاط زین.

گرادیان کاهش ما را قادر می سازد به نقاط ثابت برسیم: نقاط w که در آن شیب صفر است. ابتدا حالت تک متغیره را در نظر بگیرید. مشتق میزان تغییر سطح تابع در نقطه w را به ما می گوید. هنگامی که مشتق هدف در $w \in \mathbb{R}$ صفر باشد، یعنی $\frac{d}{dw} c(w) = 0$ ، به این معنی است که سطح تابع به صورت محلی صاف است. چنین نقاطی مطابق شکل ۱، ۲ با مینیمم های محلی، ماکزیمم های محلی و نقاط زینی است.

به عنوان مثال، دوباره فرض کنید که ما در حال انجام رگرسیون خطی هستیم، تنها با یک ویژگی و بنابراین فقط یک وزن $w \in \mathbb{R}$. مشتق هدف $c(w) = \sum_{i=1}^n (x_i w - y_i)^2$ به این صورت است.

$$\frac{d}{dw} c(w) = \frac{d}{dw} \sum_{i=1}^n (x_i w - y_i)^2$$

$$= \sum_{i=1}^n \frac{d}{dw} (x_i w - y_i)^2$$

$$= \sum_{i=1}^n 2(x_i w - y_i) x_i$$

جایی که آخرین گام از قانون زنجیره پیروی می‌کند. هدف ما یافتن w به گونه‌ای است که $\frac{d}{dw} c(w) = 0$; هنگامی که چنین نقطه ثابتی را پیدا کردیم، می‌توانیم تعیین کنیم که آیا آن یک مینیمم محلی، ماکسیمم محلی یا نقطه زینی است. از آنجایی که این هدف محدب است، ما در واقع می‌دانیم که تمام نقاط ثابت باید مینیمم جهانی باشند و بنابراین نیازی به انجام این بررسی نداریم. ما در بخش آخر در این مورد بیشتر بحث می‌کنیم، جایی که در مورد برخی از ویژگی‌های اهداف بحث می‌کنیم.

برای حالت چند متغیره، به جای مشتقات، باید گرادیان‌ها را در نظر بگیریم. برای $w \in \mathbb{R}^d$ که $d > 1$ است، باید بپرسیم: بسته به اینکه هر عنصر w چگونه تغییر می‌کند، تابع به صورت محلی چگونه تغییر می‌کند؟ برای تعیین این کمیت، از گرادیان استفاده می‌کنیم که از مشتقات جزئی تشکیل شده است.

$$\nabla c(w) = \left[\frac{\partial c}{\partial w_1}(w) \quad \frac{\partial c}{\partial w_2}(w) \quad \dots \quad \frac{\partial c}{\partial w_d}(w) \right]$$

هر مشتق جزئی $\frac{\partial c}{\partial w_j}(w)$ نحوه تغییر تابع c را نشان می‌دهد، زمانی که فقط w_j تغییر می‌کند و دیگری ثابت نگه داشته می‌شوند. به عنوان مثال، برای $w = (w_1, w_2)$ داریم $c(w) = \frac{1}{2} (x_1 w_1 + x_2 w_2 - y)^2$ ، مشتقات جزئی هستند

$$\frac{\partial c}{\partial w_1}(w) = (x_1 w_1 + x_2 w_2 - y) x_1$$

$$\frac{\partial c}{\partial w_2}(w) = (x_1 w_1 + x_2 w_2 - y) x_2$$

به طور مفید، ما مجبور نیستیم در نظر بگیریم که چگونه کل بردار به طور مشترک در همه متغیرها تغییر می‌کند. بلکه کافی است نقاط ثابت را با یافتن w در جایی که مشتقات جزئی صفر هستند پیدا کنیم.

۲-۲ گرادیان کاهشی

ایده اصلی پشت شیب نزول، تقریب تابع با تقریب سری تیلور است. این تقریب محاسبه جهت نزول را به صورت محلی روی سطح تابع تسهیل می‌کند. ما با در نظر گرفتن تنظیم تک متغیره، با $w \in \mathbb{R}$ شروع می‌کنیم. یک تابع $c(w)$ در همسایگی نقطه w_0 ، می‌تواند با استفاده از سری تیلور به صورت تقریبی باشد.

$$c(w) = \sum_{n=0}^{\infty} \frac{c^{(n)}(w_0)}{n!} (w - w_0)^n$$

که $c^{(n)}(w_0)$ آن n -امین مشتق تابع $c(w)$ است که در نقطه w_0 ارزیابی شده است. این فرض می‌کند که $c(w)$ بی‌نهایت قابل تفکیک است، اما در عمل ما چنین تصاویر تقریبی چند جمله‌ای را برای n محدود می‌گیریم. یک تقریب مرتبه دوم برای این تابع از سه عبارت اول سری به استفاده می‌کند

$$c(w) \approx \hat{c}(w) = c(w_0) + (w - w_0)\dot{c}(w_0) + \frac{1}{2}(w - w_0)^2\ddot{c}(w_0)$$

یک نقطه ثابت از این $\hat{c}(w)$ را می‌توان با پیدا کردن اولین مشتق و صفر کردن آن به راحتی پیدا کرد.

$$\dot{c}(w) \approx \dot{c}(w_0) + (w - w_0)\ddot{c}(w_0) = 0$$

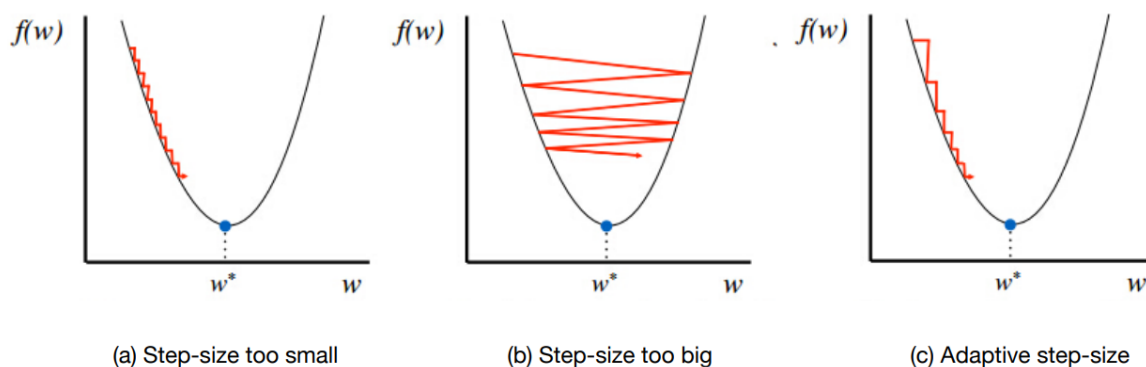
حل این معادله برای w به ما می‌دهد

$$w_1 = w_0 - \frac{c(w_0)}{\ddot{c}(w_0)}$$

به صورت محلی، این w_1 جدید یک پیشرفت در w_0 خواهد بود و یک نقطه ثابت از این تقریب محلی \hat{c} خواهد بود. با این حال، حرکت (به اندازه کافی دور) از w_0 باعث می‌شود این سری محلی تیلور مرتبه دوم نادرست باشد. ما باید تقریب محلی را در این نقطه جدید w_1 بررسی کنیم تا مشخص کنیم که آیا می‌توانیم به صورت محلی بهبود بیشتری داشته باشیم. بنابراین، برای یافتن w بهینه، می‌توانیم به طور مکرر این روش را اعمال کنیم

$$w_{t+1} = w_t - \frac{c(w_t)}{\ddot{c}(w_t)} \quad (2.1)$$

دائماً w_i را بهبود می‌بخشد تا زمانی که به نقطه ای برسیم که مشتق صفر یا تقریباً صفر باشد. این روش را روش نیوتن رافسون یا شیب نزول مرتبه دوم می‌نامند.



نیکل ۲/۲: مسیرهای بهینه سازی مختلف، به دلیل انتخاب های مختلف اندازه.

در گرادینان کاهشی مرتبه اول، تقریب بد است، جایی که ما دیگر از مشتق دوم استفاده نمی‌کنیم. در عوض، هنگام گرفتن تقریب مرتبه اول، می‌دانیم که عبارت‌های $O((w - w_0)^2)$ را نادیده می‌گیریم و بنابراین تقریب محلی تبدیل می‌شود به:

$$c(w) \approx \hat{c}(w) = c(w_0) + (w - w_0)\dot{c}(w_0) + \frac{1}{2\eta}(w - w_0)^2$$

برای مقدار ثابت $\frac{1}{\eta}$ که بزرگی عبارات $O((w - w_0)^2)$ نادیده گرفته شده را منعکس می‌کند. سپس به روزرسانی حاصل برای گام به اندازه η_t است انجام می‌شود

$$w_{t+1} = w_t - \eta_t \dot{c}(w_t) \quad (2.2)$$

از این، می‌توان دریافت که با توجه به دسترسی به مشتق دوم، یک انتخاب معقول برای اندازه مراحل $\eta_t = \frac{1}{\dot{c}(w_t)}$ است. ما می‌توانیم به طور مشابه چنین قوانینی را برای متغیرهای چند متغیره بدست آوریم. به عنوان مثال، شیب نزول برای $c: \mathbb{R}^d \rightarrow \mathbb{R}$ شامل به روزرسانی است

$$w_{t+1} = w_t - \eta_t \nabla c(w_t).$$

به طوری که

$$\nabla c(w_t) = \left(\frac{\partial c}{\partial w_1}(w_1), \frac{\partial c}{\partial w_2}(w_2), \dots, \frac{\partial c}{\partial w_d}(w_d) \right) \in \mathbb{R}^d$$

گرایان تابع c است که در w_t ارزیابی می‌شود. ما در مورد نحوه استخراج این به روز رسانی در تنظیمات چند متغیره در فصل ۶ بحث خواهیم کرد.

۳-۲ انتخاب اندازه گام

بخش مهمی از گرایان کاهشی (مرتبه اول) انتخاب اندازه گام است. اگر اندازه گام خیلی کوچک باشد، برای رسیدن به یک نقطه ثابت نیاز به تکرارهای زیادی است (شکل ۲،۲ (a)). اگر اندازه گام خیلی بزرگ باشد، احتمالاً حول مینیمم نوسان خواهید داشت (شکل ۲،۲ (b)). چیزی که ما واقعاً می‌خواهیم یک اندازه گام تطبیقی است (شکل ۲،۲ (c))، که احتمالاً بزرگ‌تر شروع می‌شود و سپس با نزدیک شدن به یک نقطه ثابت به آرامی در طول زمان کاهش می‌یابد.

روش اصلی برای به دست آوردن اندازه‌های گام تطبیقی استفاده از جستجوی خط است. این ایده از هدف زیر سرچشمه می‌گیرد: ما می‌خواهیم اندازه گام بهینه را مطابق با آن به دست آوریم

$$\min_{\eta \in \mathbb{R}^+} c(w_t - \eta \nabla c(w_t))$$

راه‌حل این بهینه‌سازی مربوط به بهترین اندازه مقیاس اسکالر است که می‌توانیم برای نقطه w_t فعلی با جهت نزول $-\nabla c(w_t)$ انتخاب کنیم. حل این بهینه‌سازی بسیار پرهزینه خواهد بود. با این حال، ما می‌توانیم به سرعت راه‌های تقریبی پیدا کنیم. یک انتخاب طبیعی استفاده از یک جستجوی خط عقبگرد است که بزرگترین اندازه گام معقول η_{\max} را امتحان می‌کند و سپس آن را کاهش می‌دهد تا هدف کاهش یابد. ایده این است که در امتداد خط ممکن $\eta \in (0, \eta_{\max}]$ جستجو کنید، با این تصور که یک گام بزرگ خوب است - تا زمانی که بیش از حد نباشد. اگر بیش از حد بالا برود و اندازه گام خیلی بزرگ شود، و باید کاهش معمولاً طبق قانون $\tau\eta$ برای مقداری $\tau \in [0.5, 0.9]$ است. برای $\tau = 0.5$ ، اندازه گام سریعتر کاهش می‌یابد - در هر مرحله از جستجوی خط عقب نشینی به نصف می‌رسد؛ برای $\tau = 0.9$ ، جستجو آهسته‌تر از η_{\max} عقب می‌نشیند. به محض اینکه اندازه گامی پیدا شد که هدف را کاهش می‌دهد، پذیرفته می‌شود. سپس یک w_t جدید به دست می‌آوریم، دوباره گرایان را محاسبه می‌کنیم و یک بار دیگر جستجوی خط را از η_{\max} شروع می‌کنیم.

می‌توان استراتژی‌های بهتری را برای انتخاب اندازه گام نسبت به این جستجوی ساده تصور کرد. ما در واقع برخی از این موارد را در بخش ۴،۵ مورد بحث قرار خواهیم داد. با این وجود، این جستجوی خط اولیه، شهودی را برای هدف ما در تطبیق اندازه مراحل فراهم می‌کند.

Algorithm 1: Line Search($w_t, c, g = \nabla c(w_t)$)

```

1: Optimization parameters:  $\eta_{max} = 1.0, \tau = 0.7, tolerance \leftarrow 10e^{-4}$ 
2:  $\eta \leftarrow \eta_{max}$ 
3:  $w \leftarrow w_t$ 
4:  $obj \leftarrow c(w)$ 
5: while number of backtracking iterations is less than maximum iterations do
6:    $w \leftarrow w_t - \eta g$ 
7:   // Ensure improvement is at least as much as tolerance
8:   If  $c(w) < obj - tolerance$  then break
9:   // Else, the objective is worse and so we decrease stepsize
10:   $\eta \leftarrow \tau \eta$ 
11: if maximum number of iterations reached then
12:  // Could not improve solution
13:  return  $w_t, \eta = 0$ 
14: return  $w, \eta$ 

```

۴-۲ خواص بهینه سازی

چندین ویژگی بهینه سازی وجود دارد که باید هنگام مطالعه این کتاب در نظر داشت که در اینجا به آنها اشاره می‌کنیم. به ماکسیم رساندن در مقابل به مینیم رساندن^۱ ما تاکنون در مورد هدف به مینیم رساندن یک هدف بحث کرده‌ایم. یک جایگزین معادل، به ماکسیم رساندن منفی این هدف است.

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} c(\mathbf{w}) = \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^d} -c(\mathbf{w})$$

جایی که argmin ، \mathbf{w} را برمی‌گرداند که مینیم مقدار $c(\mathbf{w})$ را تولید می‌کند و argmax ، \mathbf{w} را برمی‌گرداند که ماکسیم مقدار $-c(\mathbf{w})$ را تولید می‌کند. مقادیر واقعی مینیم و ماکسیم یکسان نیستند، زیرا برای یک جواب بهینه داده شده، $c(\mathbf{w}) \neq -c(\mathbf{w})$ است. ما تصمیم می‌گیریم هر یک از بهینه‌سازی‌هایمان را به‌عنوان کمینه‌سازی فرمول‌بندی کنیم و شیب نزول را انجام دهیم. با این حال، فرمول‌بندی بهینه‌سازی‌ها به‌عنوان بیشینه‌سازی و انجام صعود گرادیان به همان اندازه معتبر است.

تابع محدب: به تابع $c: \mathbb{R}^d \rightarrow \mathbb{R}$ محدب گفته می‌شود اگر برای هر $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^d$ و $t \in [0, 1]$

$$c(t\mathbf{w}_1 + (1 - t)\mathbf{w}_2) \leq tc(\mathbf{w}_1) + (1 - t)c(\mathbf{w}_2) \quad (2.4)$$

¹ Maximizing versus minimizing

این تعریف به این معنی است که وقتی بین هر دو نقطه در سطح تابع خطی می‌کشیم، مقادیر تابع بین این دو نقطه همگی زیر این خط قرار می‌گیرند. تحدب یک ویژگی مهم است، زیرا به این معنی است که هر نقطه ثابت یک مینیمم مطلق است. بنابراین، صرف نظر از اینکه نزول شیب خود را از کجا شروع می‌کنیم، با اندازه گام‌های مناسب انتخاب شده و تکرارهای کافی، به یک راه حل بهینه خواهیم رسید.

یک تعریف مربوطه یک تابع مقعر است که دقیقاً برعکس است: همه نقاط بالای خط قرار دارند. برای هر تابع محدب c ، منفی آن تابع $-c$ یک تابع مقعر است.

منحصر به فرد بودن راه حل ما اغلب اهمیت می‌دهیم که بیش از یک راه حل برای مشکل بهینه‌سازی ما وجود داشته باشد. در برخی موارد، ما به قابلیت شناسایی اهمیت می‌دهیم، به این معنی که می‌توانیم راه‌حل واقعی را شناسایی کنیم. اگر بیش از یک راه‌حل وجود داشته باشد، ممکن است تصور شود که مشکل دقیقاً مطرح نشده است. برای برخی از مشکلات، مهم یا حتی ضروری است که قابلیت شناسایی داشته باشیم (به عنوان مثال، تخمین درصد افراد مبتلا به یک بیماری) در حالی که برای برخی دیگر ما صرفاً به یافتن یک عملکرد مناسب (پیش‌بینی‌کننده) اهمیت می‌دهیم که به طور منطقی و دقیق اهداف را پیش‌بینی کند، حتی اگر آن را پیش‌بینی کند. چنین عملکرد منحصر به فردی نیست. ما قابلیت شناسایی را در این سند بیشتر از این در نظر نخواهیم گرفت، اما مهم است که بدانیم آیا هدف شما راه‌حل‌های متعددی دارد یا خیر.

هم ارزی تحت یک جابجایی ثابت جمع یا ضرب در یک ثابت $a \neq 0$ جواب را تغییر نمی‌دهد

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} c(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} a c(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} c(\mathbf{w}) + a$$

با گرفتن گرادیان هر سه هدف و مشاهده صفر بودن گرادیان در شرایط یکسان می‌توانید دلیل آن را ببینید.

$$\nabla a c(\mathbf{w}) = 0 \Leftrightarrow a \nabla c(\mathbf{w}) = 0 \Leftrightarrow \nabla c(\mathbf{w}) = 0$$

9

$$\nabla(c(\mathbf{w}) + a) = 0 \Leftrightarrow \nabla c(\mathbf{w}) = 0$$

فصل ۳

اصول اولیه تخمین پارامتر

در مدل‌سازی احتمالی، معمولاً مجموعه‌ای از مشاهدات به ما ارائه می‌شود و هدف یافتن مدل یا تابعی است که مطابقت خوبی با داده‌ها نشان می‌دهد و الزامات اضافی خاصی را رعایت می‌کند. ما تقریباً این الزامات را به سه گروه دسته‌بندی می‌کنیم: (۱) توانایی تعمیم خوب، (۲) توانایی ترکیب دانش و فرضیات قبلی در مدل‌سازی، و (۳) مقیاس‌پذیری. اول، مدل باید بتواند در آزمون زمان مقاومت کند. یعنی عملکرد آن روی داده‌های دیده نشده قبلی نباید پس از ارائه این داده‌های جدید بدتر شود. گفته می‌شود مدل‌هایی با چنین عملکردی به خوبی تعمیم می‌یابند. دوم، \hat{f} باید بتواند اطلاعات مربوط به فضای مدل \mathcal{F} را که از آن انتخاب شده است، ترکیب کند و فرآیند انتخاب یک مدل باید بتواند "مشاوره" آموزشی را از یک تحلیلگر بپذیرد. در نهایت، زمانی که حجم زیادی از داده در دسترس است، الگوریتم‌های یادگیری باید بتوانند با توجه به منابعی مانند حافظه یا قدرت CPU، راه‌حلی را در زمان معقول ارائه دهند. به طور خلاصه، انتخاب یک مدل در نهایت به مشاهدات در دست، تجربه ما در مدل‌سازی پدیده‌های زندگی واقعی، و توانایی الگوریتم‌ها برای یافتن راه‌حل‌های خوب با توجه به منابع محدود بستگی دارد.

یک راه آسان برای فکر کردن در مورد یافتن "بهترین" مدل از طریق یادگیری پارامترهای یک توزیع است. فرض کنید به ما مجموعه‌ای از مشاهدات داده شده است $\mathcal{D} = \{x_i\}_{i=1}^n$ ، جایی که $x_i \in R$ و می‌دانیم که x_i ها $i.i.d$ هستند. از توزیع گاوسی. در این مورد، مشکل یافتن بهترین مدل را می‌توان به عنوان یافتن بهترین پارامترهای μ^* و σ^* مشاهده کرد: مشکل را می‌توان به عنوان تخمین پارامتر مشاهده کرد. ما به این فرآیند، تخمین می‌گوییم زیرا فرض معمول این است که داده‌ها توسط یک مدل ناشناخته از \mathcal{F} تولید شده است که پارامترهای آن را می‌خواهیم از داده‌ها بازیابی کنیم. ما تخمین پارامتر را با استفاده از تکنیک‌های احتمالی فرمول‌بندی می‌کنیم و متعاقباً راه‌حلی را از طریق بهینه‌سازی پیدا می‌کنیم، گاهی اوقات با محدودیت‌هایی در فضای پارامتر.

۳-۱ نقشه و برآورد ماکسیم احتمال

تصور کنید مجموعه داده‌ای از مشاهدات را مشاهده می‌کنید $\mathcal{D} = \{x_i\}_{i=1}^n$. داده‌ها از مقداری توزیع واقعی p^* گرفته می‌شوند، اما این توزیع برای شما ناشناخته است. در عوض، تنها چیزی که می‌دانید این است که توزیع در مجموعه‌ای از توزیع‌های ممکن، \mathcal{F} است که گاهی اوقات فضای فرضی یا کلاس تابع نامیده می‌شود. به عنوان مثال، \mathcal{F} می‌تواند خانواده همه توزیع‌های گاوسی تک متغیره باشد:

$$\mathcal{F} = \{\mathcal{N}(\mu, \sigma^2) \mid \text{for any } \mu \in \mathbb{R} \text{ and } \sigma \in \mathbb{R}^+\}$$

¹ Independent and identically distributed

توزیع واقعی دارای پارامترهای μ^* و σ^* است. با استفاده از داده‌ها، می‌خواهیم μ و σ را تا حد امکان به تابع هدف نزدیک کنیم. ایده پشت آن ماکسیمم تخمین پسینی^۱ (MAP) یافتن محتمل‌ترین مدل برای داده‌های مشاهده شده است. با توجه به مجموعه داده‌های \mathcal{D} ، راه حل MAP را اینگونه فرمول‌بندی می‌کنیم

$$f_{\text{MAP}} = \underset{f \in \mathcal{F}}{\operatorname{argmax}} p(f|\mathcal{D})$$

که در آن $p(f|\mathcal{D})$ توزیع پسین مدل با توجه به داده‌ها نامیده می‌شود. در فضاهای مدل گسسته، $p(f|\mathcal{D})$ تابع جرم احتمال و تخمین MAP دقیقاً محتمل‌ترین مدل است. همتای آن در فضاهای پیوسته مدلی است که بیشترین مقدار تابع چگالی پسین را دارد. توجه داشته باشید که ما از مدل کلمات که یک تابع است و پارامترهای آن که ضرائب آن تابع هستند تا حدودی به جای هم استفاده می‌کنیم. به عنوان مثال، در بالا، می‌توانیم به طور معادل $\mathcal{F} = \{\mu \in \mathbb{R}, \sigma \in \mathbb{R}^+\}$ را در نظر بگیریم. ما معمولاً به جای اینکه به طور غیرمستقیم در مورد مدل‌ها یا احتمالاتی که آنها پارامتر را انتخاب می‌کنند، به طور مستقیم درباره فضای پارامتر یا فضای تابع استدلال می‌کنیم.

برای محاسبه توزیع پسین ما با اعمال قانون بیز شروع می‌کنیم

$$p(f|\mathcal{D}) = \frac{p(\mathcal{D}|f)p(f)}{p(\mathcal{D})} \quad (3.1)$$

که در آن $p(\mathcal{D}|f)$ تابع درست‌نمایی نامیده می‌شود، $p(f)$ توزیع قبلی مدل و $p(\mathcal{D})$ توزیع حاشیه‌ای داده‌ها است. توجه داشته باشید که ما از \mathcal{D} برای مجموعه داده‌های مشاهده شده استفاده می‌کنیم، اما معمولاً آن را به عنوان تحقق یک متغیر تصادفی چند بعدی \mathcal{D} که بر اساس توزیع $p(\mathcal{D})$ ترسیم می‌شود، در نظر می‌گیریم. با استفاده از فرمول احتمال کل می‌توانیم $p(\mathcal{D})$ را به این صورت بیان کنیم.

$$p(\mathcal{D}) = \begin{cases} \sum_{f \in \mathcal{F}} p(\mathcal{D}|f)p(f) & f: \text{discrete} \\ \int_{\mathcal{F}} p(\mathcal{D}|f)p(f)df & f: \text{continuous} \end{cases}$$

بنابراین، توزیع پسین را می‌توان به طور کامل با استفاده از احتمال و پیشین توصیف کرد. حوزه تحقیق و عملی که شامل روش‌های تعیین این توزیع و مدل‌های بهینه است، آمار استنباطی نامیده می‌شود.

یافتن f_{MAP} را می‌توان تا حد زیادی ساده کرد زیرا $p(\mathcal{D})$ در مخرج بر جواب تاثیر نمی‌گذارد. ما باید معادله (۳،۱) را دوباره بنویسیم

$$\begin{aligned} p(f|\mathcal{D}) &= \frac{p(\mathcal{D}|f)p(f)}{p(\mathcal{D})} \\ &\propto p(\mathcal{D}|f) \cdot p(f) \end{aligned}$$

که در آن \propto نماد تناسب است. بنابراین، با حل مسئله بهینه سازی زیر می‌توانیم راه حل MAP را پیدا کنیم

$$f_{\text{MAP}} = \underset{f \in \mathcal{F}}{\operatorname{argmax}} p(\mathcal{D}|f)p(f)$$

¹ maximum a posteriori

در برخی شرایط ممکن است دلیلی برای ترجیح یک مدل بر مدل دیگر نداشته باشیم و می‌توانیم $p(f)$ را به عنوان یک ثابت نسبت به فضای مدل \mathcal{F} در نظر بگیریم. سپس، MAP به ما کمترین کردن تابع درست‌نمایی کاهش می‌دهد:

$$f_{MLE} = \operatorname{argmax}_{f \in \mathcal{F}} p(\mathcal{D}|f)$$

این راه‌حل را راه‌حل ماکسیمم احتمال¹ (MLE) می‌نامند. به طور رسمی، فرض ثابت بودن $p(f)$ مشکل‌ساز است زیرا یک توزیع یکنواخت را نمی‌توان همیشه تعریف کرد (مثلاً روی \mathbb{R} ، اگرچه راه‌حلهایی برای این موضوع با استفاده از پیشین‌های نامناسب وجود دارد. با این وجود، فکر کردن به MLE به عنوان یک مورد خاص از تخمین MAP مفید است.

مثال ۹: فرض کنید مجموعه داده $D = \{2, 5, 9, 5, 4, 8\}$ یک $i.i.d$ است. نمونه‌ای از توزیع پواسون با پارامتر ثابت اما ناشناخته λ_0 . تخمین ماکسیمم درست‌نمایی λ_0 را پیدا کنید.

تابع جرم احتمال توزیع پواسون به صورت $p(x|\lambda) = \frac{\lambda^x e^{-\lambda}}{x!}$ با پارامتر $\lambda \in \mathbb{R}^+$ بیان می‌شود. ما این پارامتر را به این صورت تخمین می‌زنیم

$$\lambda_{MLE} = \operatorname{argmax}_{\lambda \in (0, \infty)} p(\mathcal{D}|\lambda) \quad (3.2)$$

می‌توانیم تابع احتمال را به صورت زیر بنویسیم

$$p(\mathcal{D}|\lambda) = p(\{x_i\}_{i=1}^n | \lambda) \\ \prod_{i=1}^n p(x_i | \lambda)$$

که در آن احتمال به احتمالات فردی هر x_i تقسیم می‌شود زیرا داده‌ها $i.i.d$ هستند. (متغیرهای تصادفی مستقل هستند). برای یافتن λ که احتمال را به ما کمترین می‌رساند، ابتدا از یک لگاریتم (یک تابع یکنواخت) برای ساده کردن محاسبه استفاده می‌کنیم. سپس اولین مشتق آن را نسبت به λ پیدا کنید. و در نهایت آن را برابر با صفر کنید تا ما کمترین را بدست آورید. به طور خاص، ما احتمال ورود به سیستم $\ln p(\mathcal{D}|\lambda)$ را به این صورت بیان می‌کنیم

$$\begin{aligned} \ln p(\mathcal{D}|\lambda) &= \ln \prod_{i=1}^n p(x_i | \lambda) \\ &= \sum_{i=1}^n \ln p(x_i | \lambda) \\ &= \ln \lambda \sum_{i=1}^n x_i - n\lambda - \sum_{i=1}^n \ln(x_i!) \end{aligned}$$

زیرا

$$\ln p(x_i | \lambda) = \frac{\ln \lambda^{x_i} e^{-\lambda}}{(x_i)!}$$

¹ maximum likelihood

$$= \ln \lambda^{x_i} + \ln e^{-\lambda} - \ln x_i !$$

$$= x_i \ln \lambda - \lambda - \ln x_i !$$

اکنون در این فرم، ساده‌تر به محاسبه مشتق می‌پردازیم

$$\frac{\partial \ln p(\mathcal{D}|\lambda)}{\partial \lambda} = \frac{1}{\lambda} \sum_{i=1}^n x_i - n$$

حل λ به گونه‌ای که $\frac{\partial \ln p(\mathcal{D}|\lambda)}{\partial \lambda} = 0$ یک نقطه ثابت از این مسئله را به ما می‌دهد، و ما $\lambda_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$ را دریافت می‌کنیم. می‌توانیم $n = 6$ و مقادیر \mathcal{D} را جایگزین کنیم تا جواب را به این صورت محاسبه کنیم

$$\lambda_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i = 5.5$$

که به سادگی یک میانگین نمونه است.

مشتق دوم این احتمال \log همیشه منفی است زیرا λ باید مثبت باشد. مشتق دوم $-\lambda^{-2} \sum_{i=1}^n x_i$ است که برای این λ_{MLE} که $0 <$ است. بنابراین، عبارت قبلی در واقع احتمال را به مقدار ماکسیمم می‌رساند. توجه داشته باشید که برای به ماکسیمم رساندن درست این هزینه‌ها، ما همچنین باید از اعمال محدودیت $\lambda \in (0, \infty)$ اطمینان حاصل کنیم. از آنجایی که راه‌حل بالا در مجموعه محدودیت‌ها قرار دارد، می‌دانیم که جواب صحیح معادله (۳،۲) را داریم. با این حال، در موقعیت‌های دیگر، همانطور که بعداً بحث خواهیم کرد، باید به صراحت محدودیت‌هایی را در بهینه‌سازی اعمال کنیم.

مثال ۱۰: فرض کنید $\mathcal{D} = \{2, 5, 9, 5, 4, 8\}$ دوباره یک i.i.d باشد. نمونه‌ای از پواسون (λ_0) ، اما اکنون اطلاعات اضافی نیز به ما داده می‌شود. فرض کنید دانش قبلی در مورد λ_0 را می‌توان با استفاده از توزیع گاما با پارامترهای $k = 3$ و $\theta = 1$ بیان کرد. پیدا کردن λ_0 برای تخمین MAP را بیابید.

ابتدا تابع چگالی احتمال توزیع گاما را برای قبلی خود می‌نویسیم

$$p(\lambda) = \frac{\lambda^{k-1} e^{-\frac{\lambda}{\theta}}}{\theta^k \Gamma(k)}$$

که برای $\lambda > 0$ ، $\Gamma(k)$ تابع گامایی است که تابع فاکتوریل را تعمیم می‌دهد. وقتی k یک عدد صحیح است، $\Gamma(k) = (k-1)!$ داریم. تخمین MAP پارامترها را می‌توان به صورت زیر پیدا کرد

$$\lambda_{MAP} = \underset{\lambda \in (0, \infty)}{\operatorname{argmax}} p(\mathcal{D}|\lambda)p(\lambda)$$

مانند قبل، ما \log را برای ساده کردن محاسبات می‌گیریم تا به دست آوریم

$$\ln p(\mathcal{D}|\lambda)p(\lambda) = \ln p(\mathcal{D}|\lambda) + \ln p(\lambda)$$

$$= \sum_{i=1}^n \ln p(x_i | \lambda) + \ln p(\lambda).$$

ما قبلاً عبارت اول را در مثال قبلی ساده کرده‌ایم. برای ورود به سیستم توزیع قبلی، ما داریم

$$\begin{aligned}\ln p(\lambda) &= \ln(\lambda^{k-1} e^{-\frac{\lambda}{\theta}}) - \ln(\theta^k \Gamma(k)) \\ &= (k-1) \ln \lambda - \frac{\lambda}{\theta} - \ln(\theta^k \Gamma(k)).\end{aligned}$$

جمله آخر با توجه به λ ثابت است. بنابراین وقتی مشتق را می‌گیریم ناپدید می‌شود و می‌توانیم از محاسبه آن اجتناب کنیم. دوباره همه چیز را در جای خود قرار می‌دهیم

$$\ln p(D|\lambda)p(\lambda) = \ln \lambda \sum_{i=1}^n x_i - n\lambda - \sum_{i=1}^n \ln(x_i!) + (k-1) \ln \lambda - \frac{\lambda}{\theta} - \ln(\theta^k \Gamma(k))$$

و گرفتن مشتق آن می‌دهد

$$\frac{\partial \ln p(D|\lambda)p(\lambda)}{\partial \lambda} = \frac{1}{\lambda} \sum_{i=1}^n x_i - n + \frac{k-1}{\lambda} - \frac{1}{\theta}$$

زیرا

$$\frac{\partial \ln p(\lambda)}{\partial \lambda} = \frac{k-1}{\lambda} - \frac{1}{\theta}$$

بار دیگر با صفر کردن مشتق و حل λ به دست می‌آید

$$\lambda_{\text{MAP}} = \frac{k-1 + \sum_{i=1}^n x_i}{n + \frac{1}{\theta}} = 5$$

برای مجموعه داده \mathcal{D}

نگاهی گذرا به λ_{MAP} و λ_{MLE} نشان می‌دهد که با رشد n ، هم اعداد و هم مخرج‌ها در عبارات بالا به طور فزاینده‌ای شبیه‌تر می‌شوند. در واقع، این یک نتیجه شناخته شده است که در حد نمونه‌های نامتناهی، هر دو MAP و MLE به یک مدل، f ، همگرا می‌شوند، تا زمانی که پیشین احتمال (یا چگالی) روی f را نداشته باشد. این نتیجه نشان می‌دهد که برآورد MAP به راه حل MLE برای مجموعه داده‌های بزرگ نزدیک می‌شود. به عبارت دیگر، داده‌های بزرگ از اهمیت دانش قبلی می‌کاهد. این یک نتیجه‌گیری مهم است زیرا دستگاه‌های ریاضی لازم برای استنتاج عملی را ساده می‌کند.

برای بدست آوردن مقداری شهود برای این نتیجه، نشان خواهیم داد که برآوردهای MAP و MLE برای مثال بالا با توزیع پواسون به یک راه حل همگرا می‌شوند. فرض کنید $s_n = \sum_{i=1}^n x_i$ ، که نمونه‌ای از متغیر تصادفی $X_i = \sum_{i=1}^n X_i$ است. اگر $\lim_{n \rightarrow \infty} \frac{s_n}{n^2} = 0$ (یعنی s_n سریعتر از n^2 رشد نمی‌کند)، آنگاه

$$\begin{aligned}|\lambda_{\text{MAP}} - \lambda_{\text{MLE}}| &= \left| \frac{k-1 + s_n}{n + \frac{1}{\theta}} - \frac{s_n}{n} \right| \\ &= \left| \frac{k-1}{n + \frac{1}{\theta}} - \frac{s_n}{n(n + \frac{1}{\theta})} \right|\end{aligned}$$

$$\leq \frac{|k-1|}{n + \frac{1}{\theta}} + \frac{s_n}{n(n + \frac{1}{\theta})} \xrightarrow{n \rightarrow \infty} 0$$

توجه داشته باشید که اگر $\lim_{n \rightarrow \infty} \frac{s_n}{n^2} \neq 0$ ، هر دو برآوردگر به ∞ می‌روند. با این حال، چنین دنباله‌ای از مقادیر احتمال وقوع اساساً صفر است. قضایای سازگاری برای تخمین MLE و MAP بیان می‌کنند که همگرایی با پارامترهای واقعی "قریب به یقین" یا "با احتمال 1" رخ می‌دهد تا نشان دهد که این توالی‌های نامحدود مجموعه‌ای از اندازه‌گیری-صفر را تحت شرایط معقول معینی تشکیل می‌دهند (برای اطلاعات بیشتر ببینید: [۱۹، ۱۳]).

مثال ۱۱: فرض کنید $\mathcal{D} = \{x_i\}_{i=1}^n$ یک $i.i.d$ باشد. نمونه از یک توزیع گاوسی تک متغیره. هدف ما یافتن برآوردهای ماکسیمم احتمال پارامترها است. با تشکیل تابع $\log - likelihood$ شروع می‌کنیم

$$\ln p(\mathcal{D}|\mu, \sigma) = \ln \prod_{i=1}^n p(x_i|\mu, \sigma)$$

$$n \ln \frac{1}{\sqrt{2\pi}} + n \ln \frac{1}{\sigma} - \frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}$$

ما مشتقات جزئی $\log - likelihood$ را با توجه به تمام پارامترها محاسبه می‌کنیم

$$\frac{\partial}{\partial \mu} \ln p(\mathcal{D}|\mu, \sigma) = \frac{\sum_{i=1}^n (x_i - \mu)}{\sigma^2}$$

و

$$\frac{\partial}{\partial \sigma} \ln p(\mathcal{D}|\mu, \sigma) = -\frac{n}{\sigma} + \frac{\sum_{i=1}^n (x_i - \mu)^2}{\sigma^3}$$

از اینجا می‌توانیم استخراج کنیم:

$$\mu_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$$

و

$$\sigma_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_{MLE})^2$$

برآوردهای MAP و MLE را تخمین نقطه‌ای می‌نامند. این تخمین‌ها با تخمین‌های بیزی، که کل توزیع پسین یا فواصل اطمینان را برای پارامترها تخمین می‌زنند، در تضاد هستند. ما در این کتاب اساساً بر تخمین‌های نقطه‌ای تمرکز می‌کنیم و تنها به اختصار رویکردهای بیزی را تقریباً تا پایان بررسی می‌کنیم.

۳-۲ ماکسیم احتمال برای توزیع‌های شرطی

همچنین می‌توانیم مشکلات ماکسیم احتمال را برای توزیع‌های شرطی فرموله کنیم. به یاد بیاورید که یک توزیع شرطی به شکل $p(y|x)$ برای دو متغیر تصادفی X و Y است، که در بالا توزیع حاشیه‌ای $p(x)$ یا $p(y)$ را در نظر گرفتیم. برای توزیع‌های بالا، پرسیدیم: توزیع روی این متغیر چگونه است؟ برای توزیع شرطی، در عوض می‌پرسیم: با توجه به برخی اطلاعات کمکی، اکنون توزیع روی این متغیر چگونه است؟ هنگامی که اطلاعات کمکی تغییر می‌کند، توزیع روی متغیر نیز تغییر می‌کند. برای مثال، ممکن است بخواهیم توزیع را بر فروش یک محصول خاص (Y) با توجه به ماه جاری (X) شرط کنیم. ما انتظار داریم که توزیع بر روی Y بسته به ماه متفاوت باشد.

توزیع‌های شرطی می‌توانند از هر یک از خانواده‌های توزیع مورد بحث در بالا باشند، و ما می‌توانیم به طور مشابه مسائل تخمین پارامتر را فرموله کنیم. با این حال، پارامترها معمولاً به متغیر X مرتبط هستند. ما یک مثال ساده برای نشان دادن این موضوع در زیر ارائه می‌دهیم. بسیاری از فرمول‌های تخمین پارامتری که در ادامه کتاب در نظر می‌گیریم، برای توزیع‌های شرطی هستند، زیرا در یادگیری ماشین معمولاً تعداد زیادی متغیر کمکی (ویژگی‌ها) داریم و سعی می‌کنیم تا اهداف را پیش‌بینی کنیم (یا نوع توزیع را یاد بگیریم). در فصل‌های مربوط به رگرسیون و طبقه‌بندی، نشان خواهیم داد که چند مدل می‌توانند به عنوان ماکسیم احتمال برای توزیع‌های شرطی $p(y|x)$ فرموله شوند.

مثال ۱۲: فرض کنید به شما دو متغیر تصادفی X و Y داده شده است و معتقدید $p(y|x) = N(\mu = x, \sigma^2)$ برای برخی σ ناشناخته. هدف ما تخمین این پارامتر ناشناخته σ است. توجه داشته باشید که توزیع بر روی Y متفاوت است، بسته به اینکه کدام مقدار X مشاهده یا داده شود.

ما دوباره با تشکیل تابع $\log - \text{likelihood}$ شروع می‌کنیم، اکنون برای جفت n نمونه $\mathcal{D} = (x_1, y_1), \dots, (x_n, y_n)$ ما از قانون زنجیره استفاده خواهیم کرد: $p(x_i, y_i) = p(y_i | x_i)p(x_i)$

$$\begin{aligned} \ln p(\mathcal{D}|\sigma) &= \ln \prod_{i=1}^n p(x_i, y_i | \sigma) \\ &= \ln \prod_{i=1}^n p(y_i | x_i, \sigma) p(x_i) \\ &= \sum_{i=1}^n \ln p(y_i | x_i, \sigma) + \ln p(x_i) \\ &= \sum_{i=1}^n \ln \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - x_i)^2}{2\sigma^2}\right) + \ln p(x_i) \\ &= n \ln \frac{1}{\sqrt{2\pi}} + n \ln \frac{1}{\sigma} - \frac{\sum_{i=1}^n (y_i - x_i)^2}{2\sigma^2} + \sum_{i=1}^n \ln p(x_i) \end{aligned}$$

توجه داشته باشید که ما از $\mu = x_i$ برای هر توزیع نرمال $p(y_i | x_i, \sigma)$ استفاده می‌کنیم. اکنون مشتقات جزئی $\log - \text{likelihood}$ را با توجه به پارامتر σ محاسبه می‌کنیم

$$\frac{\partial}{\partial \sigma} \ln p(\mathcal{D}|\sigma) = -\frac{n}{\sigma} + \frac{\sum_{i=1}^n (y_i - x_i)^2}{\sigma^3}$$

توجه کنید که $\frac{\partial}{\partial \sigma} \sum_{i=1}^n \ln p(x_i) = 0$ زیرا σ احتمال $p(x_i)$ پارامتری نمی‌کند. بنابراین برای بدست آوردن σ بهینه، نیازی به دانستن یا تعیین توزیع بر روی متغیر تصادفی X نیست. با قرار دادن مشتق بر روی صفر، برای به دست آوردن یک نقطه ثابت، به دست می‌آوریم.

$$\sigma_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - x_i)^2$$

۳-۳ [پیشرفته] رابطه بین به ماکسیم رساندن احتمال و واگرایی Kullback – Leibler

ما اکنون رابطه بین تخمین ماکسیم درست‌نمایی و واگرایی کول بک لایبلر را بررسی می‌کنیم. واگرایی Kullback – Leibler بین دو توزیع احتمال $p(x)$ و $q(x)$ در $X = R$ به این صورت تعریف شده است.

$$D_{KL}(p||q) = \int_{-\infty}^{\infty} (x) \log \frac{p(x)}{q(x)} dx$$

در تئوری اطلاعات، واگرایی Kullback – Leibler تفسیری طبیعی از بازده فشرده‌سازی سیگنال دارد، زمانی که کد با استفاده از توزیع غیربهینه $q(x)$ به جای توزیع صحیح (اما ناشناخته) $p(x)$ ساخته می‌شود. داده تولید شده است. با این حال، اغلب، واگرایی کولبک-لایبلر به سادگی به عنوان معیاری از واگرایی بین دو توزیع احتمال در نظر گرفته می‌شود. اگرچه این واگرایی یک متریک نیست (متقارن نیست و نابرابری مثلث را برآورده نمی‌کند) اما دارای ویژگی‌های نظری مهمی است که (۱) همیشه غیر منفی است و (۲) برابر با صفر است اگر و تنها اگر $p(x) = q(x)$

اکنون یک واگرایی بین یک توزیع احتمال تخمینی $p(x|\theta)$ و یک توزیع اساسی (درست) $p(x|\theta_0)$ را در نظر بگیرید که بر اساس آن مجموعه داده $\mathcal{D} = \{x_i\}_{i=1}^n$ ایجاد شد. واگرایی کولبک-لایبلر^۱ بین $p(x|\theta_0)$ و $p(x|\theta)$ است.

$$\begin{aligned} D_{KL}(p(x|\theta_0)||p(x|\theta)) &= \int_{-\infty}^{\infty} p(x|\theta_0) \log \frac{p(x|\theta_0)}{p(x|\theta)} dx \\ &= \int_{-\infty}^{\infty} p(x|\theta_0) \log \frac{1}{p(x|\theta)} dx - \int_{-\infty}^{\infty} p(x|\theta_0) \log \frac{1}{p(x|\theta_0)} dx \end{aligned}$$

عبارت دوم در معادله فوق صرفاً آنتروپی (دیفرانسیل) توزیع واقعی است و تحت تأثیر انتخاب ما از مدل θ نیست. از طرف دیگر اصطلاح اول را می‌توان به این صورت بیان کرد

$$\int_{-\infty}^{\infty} p(x|\theta_0) \log \frac{1}{p(x|\theta)} dx = \mathbb{E}[\log p(X|\theta)]$$

بنابراین، به مقدار ماکسیم رساندن $\mathbb{E}[\log p(X|\theta)]$ واگرایی کولبک-لایبلر بین $p(x|\theta_0)$ و $p(x|\theta)$ را به مقدار مینیمم می‌رساند. با استفاده از قانون قوی اعداد بزرگ، این را می‌دانیم که

^۱ Kullback-Leibler (KL)

$$\frac{1}{n} \sum_{i=1}^n \log p(x_i | \theta) \xrightarrow{\text{a.s.}} \mathbb{E}[\log p(X|\theta)]$$

وقتی $n \rightarrow \infty$ ، بنابراین، زمانی که مجموعه داده به اندازه کافی بزرگ باشد، به ماکسیمم رساندن تابع احتمال، واگرایی *Kullback – Leibler* را به مینیمم می‌رساند و به این نتیجه می‌رسد که $p(x|\theta_{MLE}) = p(x|\theta_0)$ ، اگر مفروضات اساسی برآورده شوند. تحت شرایط معقول، می‌توانیم از آن استنباط کنیم که $\theta_{MLE} = \theta_0$. این برای خانواده‌هایی از توزیع‌ها که مجموعه‌ای از پارامترها به طور منحصربه‌فرد توزیع احتمال را تعیین می‌کنند، صادق است. به عنوان مثال، به طور کلی برای مخلوطی از توزیع‌ها صدق نمی‌کند، اما ما بعداً در مورد این وضعیت بحث خواهیم کرد. این نتیجه تنها یکی از بسیاری از ارتباطات بین آمار و نظریه اطلاعات است.

فصل ۴

مقدمه‌ای بر مسائل پیش‌بینی

یادگیری ماشین به بسیاری از مسائل پاسخ می‌دهد، که گاهی اوقات ممکن است باعث ایجاد احساس مشکل شوند. به عنوان یک فهرست غیر جامع، این موارد شامل یادگیری تحت نظارت (با طبقه‌بندی^۱ و رگرسیون^۲) می‌شود. یادگیری نیمه نظارتی؛ یادگیری بدون نظارت؛ تکمیل یادگیری فاقد دسته‌بندی. پیش‌بینی ساختار یافته؛ یادگیری رتبه‌بندی؛ یادگیری رابطه‌ای آماری؛ یادگیری فعال؛ و پیش‌بینی زمانی (با پیش‌بینی سری‌های زمانی و ارزیابی خط‌مشی در یادگیری تقویتی و یادگیری آنلاین). برای برخی از این مجموعه‌ها، مانند یادگیری فعال و یادگیری تقویتی، جمع‌آوری داده‌ها بخش مرکزی الگوریتم است و می‌تواند به طور قابل توجهی کیفیت مدل‌های پیش‌بینی آموخته شده را تعیین کند. اکثر مجموعه‌های دیگر فرض می‌کنند که داده‌ها جمع‌آوری شده‌اند بدون توانایی ما برای تأثیرگذاری بر آن مجموعه - و اکنون ما فقط باید آن داده‌ها را تجزیه و تحلیل کنیم و بهترین پیش‌بینی‌کننده‌ها را یاد بگیریم. در این مجموعه غیرفعال، می‌توانیم فرض کنیم که داده‌ها $i.i.d.$ هستند - که رایج‌ترین حالت است - یا وابستگی‌هایی بین نقاط داده وجود دارد - مانند پیش‌بینی سری‌های زمانی یا یادگیری رابطه‌ای آماری. همچنین زمینه‌هایی وجود دارد که داده‌ها ناقص هستند، مثلاً به این دلیل که کاربر سن خود را پر نکرده است

بنابراین، یک هستی‌شناسی می‌تواند ابعاد زیر را برای دسته‌بندی مسائل یادگیری ماشین در نظر بگیرد:

۱. منفعل در مقابل فعال
۲. $i.i.d.$ در مقابل $non - i.i.d.$
۳. کامل در مقابل ناقص

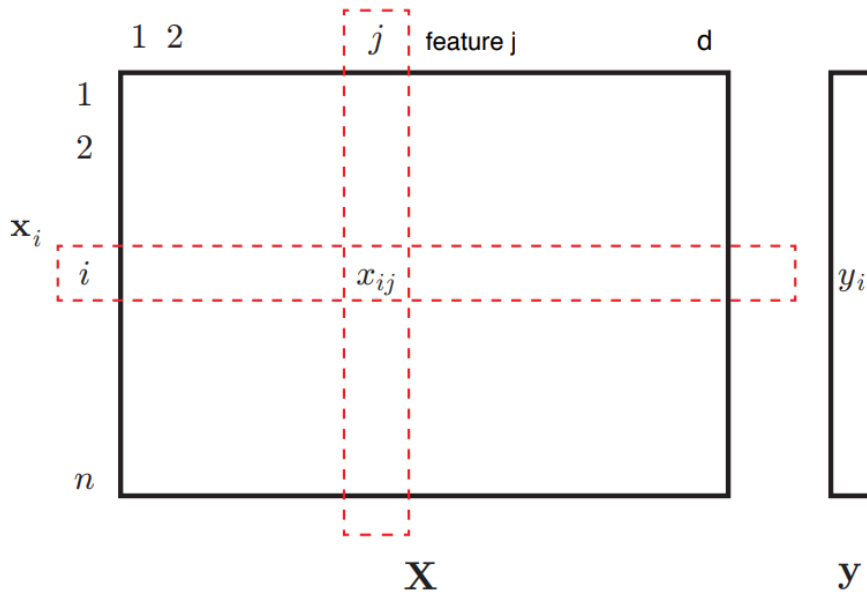
مانند همه هستی‌شناسی‌ها، هر مسئله کاملاً در این دسته‌بندی‌ها قرار نمی‌گیرد. علاوه بر این، این احتمال وجود دارد که اکثر جمع‌آوری داده‌ها کاملاً غیرفعال نباشند (حتی اگر فقط به این دلیل که مدل‌ساز انسانی بر جمع‌آوری داده‌ها تأثیر می‌گذارد)، احتمالاً $i.i.d.$ (حتی اگر قصد داشتیم باشد)، و احتمالاً برخی از اجزای گم شده را دارد. با این وجود، الگوریتم‌ها این مفروضات را به درجات مختلف انجام می‌دهند، حتی اگر داده‌ها آن مفروضات را برآورده نکنند. برای اکثر این یادداشته‌ها، ما روی ساده‌ترین مجموعه‌ها تمرکز می‌کنیم: منفعل، $i.i.d.$ و کامل در این فصل ابتدا طبقه‌بندی و رگرسیون را معرفی می‌کنیم و سپس معیارهای انتخاب توابع برای طبقه‌بندی و رگرسیون را مورد بحث قرار می‌دهیم تا الگوریتم‌های توسعه‌یافته در فصل‌های بعدی را ایجاد کنیم.

¹ Classification

² Regression

۴-۱ مسائل یادگیری تحت نظارت

ما با تعریف یک مجموعه داده شروع می‌کنیم $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$ جایی که $\mathbf{x}_i \in \mathcal{X}$ ورودی یا مشاهدات i -ام و $y_i \in \mathcal{Y}$ هدف ما است. معمولاً فرض می‌کنیم که $\mathcal{X} = \mathbb{R}^d$ ، در این صورت $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{id})$ یک بردار d بعدی است که یک نمونه نامیده می‌شود.



شکل ۴/۱: علامت گذاری برای مجموعه داده. X یک ماتریس n در d است که ردیف‌هایی مربوط به نمونه‌ها و ستون‌های مربوط به ویژگی‌ها است. y یک بردار n در ۱ از اهداف است.

یا یک نمونه هر بعد از \mathcal{X} معمولاً یک ویژگی نامیده می‌شود. ما اغلب مجموعه داده را در یک ماتریس $X \in \mathbb{R}^{n \times d}$ سازماندهی می‌کنیم که در آن هر ردیف با یک نمونه \mathbf{x}_i و هر ستون مربوط به یک ویژگی است (شکل ۴،۱ را ببینید).

تمایز بین \mathbf{x} و y به این دلیل است که فرض می‌کنیم ویژگی‌ها برای هر شی نسبتاً آسان جمع‌آوری می‌شوند (مثلاً با اندازه‌گیری قد یک فرد یا فوت مربع یک خانه)، در حالی که مشاهده یا جمع‌آوری متغیر هدف دشوار یا هزینه بر است. (مثلاً وجود بیماری یا قیمت نهایی فروش خانه قبل از فروش آن). چنین موقعیت‌هایی معمولاً از ساخت یک مدل محاسباتی استفاده می‌کنند که اهداف را از روی مجموعه‌ای از مقادیر ورودی پیش‌بینی می‌کند. این مدل با استفاده از مجموعه‌ای از مشاهدات ورودی که مقادیر هدف قبلاً جمع‌آوری شده‌اند، آموزش داده می‌شود. در استقرار، می‌توانیم از این مدل برای اهداف پیش‌بینی‌هایی از اطلاعات به‌دست‌آمده-مشاهده-در مورد اطلاعاتی که به سختی به دست می‌آیند-استفاده کنیم.

۴-۱-۱ رگرسیون و طبقه‌بندی

تفاوت در الگوریتم‌ها برای مسائل پیش‌بینی، با داده‌های کامل d ، i ، i ، معمولاً از ویژگی‌های ورودی‌ها (مشاهدات) و ویژگی‌های اهداف ناشی می‌شوند. به عنوان مثال، ممکن است لازم باشد مشاهدات متنی - مانند مشاهدات مجموعه‌ای از اسناد - را متفاوت

از یک بردار مشاهداتی با ارزش واقعی ده بعدی از خوانش‌های حسگر که دما و فشار را در یک سیستم فیزیکی منعکس می‌کند، بررسی کنیم. یک استراتژی ساده و نسبتاً رایج برای رسیدگی به این تفاوت‌ها، ترسیم انواع مختلف مشاهدات - زبان، متغیرهای طبقه‌بندی و حتی داده‌های دنباله‌ای - در فضای اقلیدسی است که در آن مشاهدات دوباره به عنوان یک بردار با ارزش واقعی نشان داده می‌شود. بسیاری از الگوریتم‌های پیش‌بینی برای مشاهدات با ارزش واقعی طراحی شده‌اند، بنابراین الگوریتم‌های استاندارد را می‌توان اعمال کرد. این مسئله نمایش داده‌ها به خودی خود یک مشکل اساسی و دشوار است. در فصل ۹ بیشتر درباره آن بحث خواهیم کرد. در حال حاضر، فرض می‌کنیم مشاهدات از قبل به شکل مناسبی هستند، به عنوان یک بردار با ارزش واقعی d بعدی.

ویژگی‌های هدف نیز مهم هستند و منجر به دو تمایز معمولی برای مسائل پیش‌بینی می‌شوند: طبقه‌بندی و رگرسیون. به طور کلی، وقتی y پیوسته است، یک مسئله از نوع رگرسیون داریم و اگر y گسسته باشد، یک مسئله از نوع طبقه‌بندی. در رگرسیون مجموعه هدف ممکن، شامل $y = \mathbb{R}$ یا $y = [0, \infty)$ است. نمونه‌ای از مسئله رگرسیون در جدول ۴,۱ نشان داده شده است.

	size [sqft]	age [yr]	dist [mi]	inc [\$]	dens [ppl/mi ²]	y
\mathbf{x}_1	1250	5	2.85	56,650	12.5	2.35
\mathbf{x}_2	3200	9	8.21	245,800	3.1	3.95
\mathbf{x}_3	825	12	0.34	61,050	112.5	5.10

جدول ۴,۱: مثالی از مسئله رگرسیون: پیش‌بینی قیمت یک خانه در یک منطقه خاص. در اینجا، ویژگی‌ها نشان دهنده اندازه خانه (size) بر حسب فوت مربع، قدمت خانه (age) بر حسب سال، فاصله از مرکز شهر (dist) بر حسب مایل، متوسط درآمد در شعاع یک مایل مربع (inc) در صد هزار دلار، و تراکم جمعیت در همان منطقه (dens). هدف نشان دهنده قیمتی است که یک خانه به آن قیمت فروخته می‌شود، به عنوان مثال.

در طبقه‌بندی، تابعی می‌سازیم که برچسب‌های کلاس گسسته را پیش‌بینی می‌کند. این تابع معمولاً طبقه‌بندی کننده^۱ نامیده می‌شود. کاردینالیت \mathcal{Y} در مسائل طبقه‌بندی معمولاً کوچک است، به عنوان مثال $\mathcal{Y} = \{\text{healthy}, \text{diseased}\}$ نمونه‌ای از یک مجموعه داده برای طبقه‌بندی با تعداد $n = 3$ داده و $d = 5$ ویژگی در جدول ۴,۲ نشان داده شده است.

مسائل طبقه‌بندی را می‌توان بیشتر به مسائل چند کلاسه^۲ و چند برچسبی تقسیم کرد. یک مسئله چند کلاسه شامل ارائه برچسب واحد برای یک ورودی است. به عنوان مثال، برای گروه خونی (ساده) با $Y = \{A, B, AB, O\}$ بیمار فقط می‌تواند با یکی از این برچسب‌ها برچسب گذاری شود. در مسائل چند کلاسه، اگر فقط دو کلاس وجود داشته باشد، طبقه‌بندی باینری نامیده می‌شود، مانند مثال در جدول ۴,۲. در چند برچسب، یک ورودی را می‌توان با بیش از یک برچسب مرتبط کرد. نمونه‌ای از مسئله چند برچسبی، طبقه‌بندی اسناد متنی به دسته‌هایی مانند {ورزش، پزشکی، مسافرت، سیاست} است. در اینجا، یک نمونه واحد ممکن است به بیش از یک مقدار در مجموعه مرتبط باشد. به عنوان مثال، مقاله‌ای در مورد پزشکی ورزشی تابع آموخته شده اکنون می‌تواند چندین خروجی را برگرداند.

¹ classifier

² classifier

به طور معمول، برای سازگاری بیشتر خروجی‌ها بین این دو مجموعه، خروجی برای چند کلاس و چند برچسب یک بردار نشانگر است. برای $|Y| = m$ ، پیش‌بینی گروه‌های خونی ممکن است $[0\ 1\ 0\ 0]$ باشد تا نشان‌دهنده گروه خونی B باشد و پیش‌بینی چهار برچسب مقاله می‌تواند $[1\ 1\ 0\ 0]$ باشد اگر هم مقاله‌ای مربوط به ورزش و هم مربوط به پزشکی باشد.

همانطور که با مشاهده، اهداف ممکن است خود پیچیده باشند، مانند اهداف متنی. یکی از حوزه‌هایی که با اهداف پیچیده‌تر سروکار دارد، پیش‌بینی خروجی ساختاریافته است، جایی که \mathcal{Y} می‌تواند مجموعه‌ای از خروجی‌های ساختاریافته باشد، به عنوان مثال. رشته‌ها، درختان یا نمودارها. کاردینالیت فضای خروجی در مسائل یادگیری ساختار یافته اغلب بسیار زیاد است. برای مثال، هنگام پیش‌بینی عملکرد یک پروتئین، کل درخت هستی‌شناسی باید پیش‌بینی شود، زیرا عملکرد خاصی زیرمجموعه‌ای از عملکردهای دیگر است. همانند مشاهدات، ممکن است بتوانیم بازنمودهای ساده‌تری را برای این اهداف پیدا کنیم تا روش‌های استانداردتری از رگرسیون و طبقه‌بندی را اعمال کنیم. مجدداً، در فصل ۹ بیشتر در این مورد بحث خواهیم کرد. در حال حاضر، اهداف نسبتاً ساده‌ای را فرض می‌کنیم، که بردارهای با ارزش واقعی m بعدی یا تعداد نسبتاً کمی از نتایج گسسته هستند.

	wt [kg]	ht [m]	T [°C]	sbp [mmHg]	dbp [mmHg]	y
\mathbf{x}_1	91	1.85	36.6	121	75	-1
\mathbf{x}_2	75	1.80	37.4	128	85	+1
\mathbf{x}_3	54	1.56	36.6	110	62	-1

جدول ۴/۲: مثالی از یک مسئله طبقه‌بندی دوتایی. پیش‌بینی وضعیت بیماری برای یک بیمار. در اینجا، ویژگی‌ها وزن (wt)، قد (ht)، دما (T)، فشار خون سیستولیک (sbp) و فشار خون دیاستولیک (dbp) را نشان می‌دهند. برچسب‌های کلاس به عنوان وجود یک بیماری خاص را نشان می‌دهد.

۲-۱-۴ تصمیم‌گیری در مورد نحوه فرمول‌بندی کردن مسئله

اگرچه ما مسائل یادگیری تحت نظارت را به دو دسته تقسیم می‌کنیم، اما همیشه مشخص نیست که چگونه یک مسئله باید فرمول‌بندی شود. به عنوان مثال، فضای خروجی $Y = \{0, 1, 2\}$ را در نظر بگیرید. می‌توانیم این مسئله را به عنوان یک مسئله طبقه‌بندی چند کلاسه در نظر بگیریم، یا می‌توانیم $Y = [0, 2]$ را فرض کنیم و یک مدل رگرسیون را یاد بگیریم. سپس می‌توانیم پیش‌بینی‌های برگردانده‌شده توسط مدل رگرسیون را با گرد کردن آن‌ها به نزدیک‌ترین عدد صحیح، آستانه‌ای کنیم.

چگونه تصمیم می‌گیرید که از کدام فرمول برای حل مسئله استفاده کنید؟ اگرچه رویه‌های ریاضی در یادگیری ماشین دقیق هستند، تصمیم‌گیری در مورد چگونگی فرمول‌بندی مسئله دنیای واقعی ظریف است و بنابراین ذاتاً واضح نیست. انتخاب یک روش خاص برای مدل‌سازی به تحلیلگر و دانش او از حوزه و همچنین جنبه‌های فنی یادگیری بستگی دارد. در این مثال، می‌توانید بپرسید: آیا به طور ذاتی ترتیبی برای خروجی‌های $\{0, 1, 2\}$ وجود دارد؟ اگر اینطور نیست، بگوییم که آنها سیب را ترجیح می‌دهند (در اولویت قرار می‌دهند)، پرتقال را ترجیح می‌دهند، یا موز را ترجیح می‌دهند مطابقت دارند، در این صورت ممکن است مدل کردن خروجی به عنوان یک بازه، که اغلب به سفارش دادن اشاره دارد، انتخاب ضعیفی باشد. از سوی دیگر، یادگیری توابع رگرسیون آسان‌تر است و اغلب پیش‌بینی‌های طبقه‌بندی به‌طور شگفت‌آور خوبی تولید می‌کنند. به علاوه، اگر برای این کلاس‌ها، ترتیبی وجود داشته باشد، مثلاً خوب، بهتر، بهترین آنگاه اکثر مدل‌های طبقه‌بندی - که ترتیبی را در خروجی‌ها فرض نمی‌کنند - نمی‌توانند از این ترتیب برای بهبود عملکرد پیش‌بینی استفاده کنند.

فرمول‌بندی کردن مسئله و انتخاب کلاس تابع و هدف گام مهمی در استفاده موثر از یادگیری ماشین است. خوشبختانه، دانش زیادی وجود دارد، به ویژه از نظر تجربی، که می‌تواند این انتخاب را هدایت کند. با یادگیری بیشتر در مورد روش‌ها، همراه با برخی اطلاعات در مورد ساختار دامنه خود، در تشخیص بهتر خواهد شد.

۲-۴ یادگیری بدون نظارت و یادگیری نیمه نظارت

مجموعه داده‌ها همیشه کامل نیستند: در برخی موارد، ما فقط می‌توانیم برچسب‌هایی را برای زیرمجموعه کوچکی از نمونه‌ها دریافت کنیم، یا اصلاً نمی‌توانیم برچسبی دریافت کنیم. برای مثال، هنگام پیش‌بینی اینکه یک گربه در یک تصویر است یا خیر، به یک انسان نیاز داریم تا هر تصویر را بگیرد و آن را با 0 یا 1 برچسب گذاری کند. همه عکس‌های گربه دارای چنین برچسب مرتبطی هستند. استفاده از یادگیری تحت نظارت فقط در این زیرمجموعه برچسب‌گذاری شده در صورتی که احتمال دارد به دلیل داده‌های محدود، پیش‌بینی‌کننده ضعیفی ایجاد کند. یادگیری نیمه نظارت شده با استفاده از تمام داده‌های بدون برچسب، برای تکمیل مجموعه داده کوچک برچسب‌گذاری شده، با یافتن ساختار در ویژگی‌ها سروکار دارد. برای مثال، ویژگی‌ها ممکن است روی یک منی‌فولد با ابعاد پایین‌تر قرار بگیرند. این ساختار را می‌توان از داده‌های بدون برچسب استنباط کرد و به طور بالقوه به طور موثرتری کلاس تابع را برای یادگیری در مجموعه داده برچسب‌گذاری شده محدود کرد. یادگیری بدون نظارت تنها در به دست آوردن این ساختار متمرکز می‌شود، بدون اینکه هدف یادگیری یک تابع برای پیش‌بینی اهداف باشد، زیرا هیچ هدفی ارائه نمی‌شود. یادگیری بدون نظارت در بخش ۹.۲.۲ به عنوان بخشی از یادگیری بازنمایی مورد بحث قرار خواهد گرفت. این دو مجموعه مسائل را می‌توان به عنوان نمونه‌ای از یک مجموعه بزرگتر یادگیری تحت داده‌های از دست رفته در نظر گرفت. به طور کلی، ممکن است نه تنها جمع‌آوری خروجی‌ها، بلکه برخی از ویژگی‌ها نیز دشوار باشد. به عنوان مثال، هنگام جمع‌آوری داده‌های بیمار، این احتمال وجود دارد که برخی از بیماران برخی از اطلاعات را حذف کنند. اگرچه جمع‌آوری اطلاعات «دارای بیماری» دشوار است، اما اطمینان از جمع‌آوری داده‌های ساده‌تر مانند «سن» یا «وزن» نیز می‌تواند دشوار باشد. علاوه بر این، حتی ممکن است پرسیم که چرا بین ویژگی‌ها و اهداف تمایز وجود دارد: همه آنها اطلاعات مرتبط با یک مورد هستند، مانند یک بیمار. با توجه به اینکه یک بیمار واقعاً یک بیماری دارد، ممکن است بخواهید از این ویژگی و سن او برای پیش‌بینی وزن او استفاده کنید - چیزی که آنها خیلی راحت ترجیح دادند فاش نکنند. این روش کلی‌تر برای نزدیک شدن به مسئله می‌تواند زمانی مفید باشد که داده‌ها از دست رفته و منجر به مشکل تکمیل کلی شود. تکنیک‌های متفاوتی اغلب در چنین مجموعه‌هایی استفاده می‌شود، و ما تا بخش ۹.۲.۲ بیشتر از این، به آن نخواهیم پرداخت. در حال حاضر، ما تمرکز خود را بر یادگیری نظارت شده حفظ می‌کنیم، که همچنان می‌توانیم از آن استفاده کنیم، حتی زمانی که برخی ویژگی‌ها از دست رفته‌اند، با استفاده از برخی اکتشاف‌های ساده برای مقابله با این داده‌های از دست رفته.

۳-۴ طبقه‌بندی بهینه و مدل‌های رگرسیون

هدف ما اکنون ایجاد معیارهای عملکردی است که برای ارزیابی پیش‌بینی‌کننده‌های $f: \mathcal{X} \rightarrow \mathcal{Y}$ و متعاقباً تعریف مدل‌های طبقه‌بندی و رگرسیون بهینه استفاده می‌شود. برای انجام این کار، فرض می‌کنیم که به توزیع مشترک واقعی $p(\mathbf{x}, \mathbf{y})$ دسترسی داریم و می‌پرسیم که پیش‌بینی بهینه در این حالت ایده‌آل چه خواهد بود. پیش‌بینی‌کننده بهینه بر اساس هزینه تابع هزینه تعریف می‌شود: $[0, \infty) \rightarrow \mathcal{Y} \times \mathcal{Y}$ ، که در آن هزینه $(\hat{\mathbf{y}}, \mathbf{y})$ هزینه یا جریمه را برای پیش‌بینی $\hat{\mathbf{y}}$ زمانی که هدف واقعی \mathbf{y} است منعکس می‌کند. از آنجایی که \mathbf{X}, \mathbf{Y} تصادفی هستند، هزینه $C = \text{cost}(f(\mathbf{X}), \mathbf{Y})$ نیز یک متغیر تصادفی

است، زیرا تابعی از این متغیرهای تصادفی است. هدف ما به مینیمم رساندن هزینه‌های مورد انتظار است. ابتدا چند نمونه از هزینه‌ها را در نظر می‌گیریم و سپس پیش‌بینی کننده‌های بهینه را استخراج می‌کنیم.

۱-۳-۴ نمونه‌هایی از هزینه‌ها

$$cost(\hat{y}, y) = \begin{cases} 0 & \text{when } y = \hat{y} \\ 1 & \text{when } y \neq \hat{y} \end{cases} \quad (4-1)$$

یک تابع هزینه پیچیده‌تر ممکن است در مجموعه‌هایی ایجاد شود که برخی از پیش‌بینی‌های نادرست مشکل‌سازتر از بقیه هستند. بیایید یک مثال عینی را در حوزه پزشکی در نظر بگیریم. فرض کنید هدف ما این است که تصمیم بگیریم که آیا یک بیمار با مجموعه‌ای از علائم خاص (x) باید برای یک آزمایش آزمایشگاهی اضافی ($y = 1$ اگر بله و $y = -1$ اگر نه) با هزینه C_{lab} فرستاده شود تا تشخیص را بهبود بخشد. با این حال، اگر آزمایش آزمایشگاهی انجام ندهیم و بعداً مشخص شود که بیمار برای درمان مناسب به آزمایش نیاز داشته است، ممکن است جریمه قابل توجهی متحمل شویم، مثلاً $C_{lawsuit}$. اگر $C_{lawsuit} \gg C_{lab}$ ، همانطور که انتظار می‌رود، طبقه‌بندی کننده باید خروجی‌های خود را به طور مناسب تنظیم کند تا تفاوت هزینه را در اشکال مختلف پیش‌بینی نادرست محاسبه کند. در اینجا، هزینه بهتر به صورت جدول نشان داده می‌شود. همیشه نمی‌توان هزینه معناداری را برای تابع تعریف کرد.

		Y	
		-1 (\neg Has Disease)	1 (Has Disease)
\hat{Y}	-1 (\neg Has Disease, No Test)	0	1000
	1 (Has Disease, Do Test)	1	1

جدول ۴/۳: تابع هزینه برای آزمایشگاه پزشکی، هزینه (\hat{y}, y) با $C_{lab} = 1$ و $C_{lawsuit} = 1000$

و بنابراین، یک معیار معقول استفاده از هزینه 0,1 پیش‌فرض در معادله (۴,۱) است. در رگرسیون، هزینه‌های متداول مربع خطا است.

$$cost(\hat{y}, y) = (\hat{y} - y)^2 \quad (4-2)$$

و قدر مطلق خطا

$$cost(\hat{y}, y) = |y - \hat{y}| \quad (4-3)$$

مربع خطا، نسبت به قدر مطلق خطا مقادیر دورتر از y را به شدت جریمه می‌کند. هزینه‌های بسیار دیگری نیز وجود دارد که در بزرگی اهداف نقش دارد، مانند درصد خطا.

۲-۳-۴ استخراج پیش‌بینی کننده‌های بهینه

ابتدا با استخراج طبقه‌بندی کننده بهینه شروع می‌کنیم. می‌توانیم هزینه مورد انتظار را به صورت زیر بیان کنیم، با فرض اینکه ورودی‌ها بردارهای با مقدار حقیقی پیوسته هستند و اهداف از یک مجموعه گسسته Y و $y = f(x)$ برای پیش‌بینی کننده f هستند.

$$\begin{aligned}\mathbb{E}[C] &= \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} \text{cost}(f(x), y) p(x, y) dx \\ &= \int_{\mathcal{X}} p(x) \sum_{y \in \mathcal{Y}} \text{cost}(f(x), y) p(y|x) dx\end{aligned}$$

جایی که ادغام کل فضای ورودی $\mathcal{X} = \mathbb{R}^d$ است. توجه داشته باشید که ما باید یک کلاس را برای هر مشاهده پیش‌بینی کنیم: $f(x)$ تنها می‌تواند یک مقدار \hat{y} را در \mathcal{Y} ارائه دهد. اما مقدار هدف تصادفی است. به همین دلیل طبقه‌بندی کننده بهینه f^* ممکن است نتواند هزینه صفر را به دست آورد. با این حال، به سادگی با نگاه کردن به معادله بالا، می‌توانیم با انتخاب بهترین طبقه‌بندی کننده برای هر x به طور جداگانه، $f^* = \text{argmin } \mathbb{E}[C]$ را بدست آوریم.

$$\begin{aligned}f^*(x) &= \underset{\hat{y} \in \mathcal{Y}}{\text{argmin}} \mathbb{E}[C|X = x] \\ &= \underset{y \in \mathcal{Y}}{\text{argmin}} \sum_{y \in \mathcal{Y}} \text{cost}(\hat{y}, y) p(y|x)\end{aligned}$$

این طبقه‌بندی کننده، طبقه‌بندی کننده ریسک Bayes نامیده می‌شود.

اگر از تابع هزینه $0 - 1$ استفاده کنیم، در معادله (۴،۱)، طبقه‌بندی کننده ریسک Bayes به سادگی تبدیل می‌شود به:

$$\begin{aligned}f^*(x) &= \underset{\hat{y} \in \mathcal{Y}}{\text{argmin}} \sum_{y \in \mathcal{Y}} \text{cost}(\hat{y}, y) p(y|x) \\ &= \underset{\hat{y} \in \mathcal{Y}}{\text{argmax}} \left(1 - \sum_{y \in \mathcal{Y}} (\text{cost}(\hat{y}, y) p(y|x)) \right) \\ &= \underset{\hat{y} \in \mathcal{Y}}{\text{argmax}} \sum_{y \in \mathcal{Y}} (1 - \text{cost}(\hat{y}, y)) p(y|x) \quad \triangleright \text{because } \sum_{y \in \mathcal{Y}} p(y|x) = 1 \\ &= \underset{\hat{y} \in \mathcal{Y}}{\text{argmax}} \sum_{y \in \mathcal{Y}, y \neq \hat{y}} 0 \cdot p(y|x) + \sum_{y \in \mathcal{Y}, y = \hat{y}} 1 \cdot p(y|x) \\ &= \underset{y \in \mathcal{Y}}{\text{argmax}} p(y|x)\end{aligned}$$

بنابراین، اگر $p(y|x)$ شناخته شده باشد یا بتوان به طور دقیق آن را یاد گرفت، ما به طور کامل مجهز به پیش‌بینی هستیم که هزینه کل را به مینیمم برساند. به عبارت دیگر، ما مسئله به مینیمم رساندن هزینه طبقه‌بندی مورد انتظار یا احتمال خطا را به مسئله توابع یادگیری، به ویژه توزیع‌های احتمال یادگیری تبدیل کرده‌ایم.

تجزیه و تحلیل برای رگرسیون مشابه تجزیه و تحلیل طبقه‌بندی است. در اینجا نیز، ما علاقه‌مند به، به مینیمم رساندن هزینه مورد انتظار برای پیش‌بینی هدف واقعی y هستیم که از یک پیش‌بینی کننده $f(x)$ استفاده می‌شود. هزینه مورد انتظار را می‌توان به این صورت بیان کرد:

$$\mathbb{E}[C] = \int_x \int_y \text{cost}(f(x), y) p(x, y) dy dx$$

برای سادگی، مربع خطای معادله (f, y) را در نظر خواهیم گرفت.

$$\text{cost}(f(x), y) = (f(x) - y)^2$$

که منجر به

$$\begin{aligned} \mathbb{E}[C] &= \int_x \int_y (f(x) - y)^2 p(x, y) dy dx \\ &= \int_x p(x) \underbrace{\int_y (f(x) - y)^2 p(y|x) dy}_{g(f(x))} dx \end{aligned}$$

با فرض اینکه $f(x)$ به اندازه کافی منعطف است که به طور جداگانه برای هر واحد حجم dx بهینه شود، می‌بینیم که به مینیمم رساندن $\mathbb{E}[C]$ ما را به مسئله پیدا کردن \hat{y} برای هر x برای کمینه کردن سوق می‌دهد.

$$g(\hat{y}) = \int_y (\hat{y} - y)^2 p(y|x) dy$$

برای یافتن \hat{y} بهینه، می‌توانیم این مسئله کمینه‌سازی را با یافتن یک نقطه ثابت، مینیمم مطلق، حل کنیم. برای این کار، از g نسبت به \hat{y} مشتق می‌گیریم و نقطه‌ای را می‌یابیم که مشتق برابر با صفر است.

$$\frac{\partial g(\hat{y})}{\partial \hat{y}} = 2 \int_y (\hat{y} - y) P(y|x) dy = 0$$

$$\Rightarrow \hat{y} \underbrace{\int_y p(y|x) dy}_{=1} = \int_y p(y|x) dy$$

$$\Rightarrow \hat{y} \underbrace{\int_y p(y|x) dy}_{=1} = \int_y p(y|x) dy$$

$$\Rightarrow \hat{y} \int_y p(y|x) dy = \mathbb{E}[Y|x]$$

بنابراین، پیش‌بینی بهینه اینگونه است

$$f^*(x) = \mathbb{E}[Y|x]$$

بنابراین، مدل رگرسیون بهینه به معنای به مینیمم رساندن مربع خطای بین پیش‌بینی و هدف واقعی، انتظار شرطی $\mathbb{E}[Y|x]$ است.

تمرین ۵: ما می‌توانیم به طور مشابه پیش‌بینی کننده بهینه برای هزینه خطای مطلق را در رابطه (۴,۳) محاسبه کنیم. نشان دهید که پیش‌بینی کننده بهینه برای خطای مطلق، میانه شرطی، میانه $[Y | X = x]$ است.

موارد فوق باعث شده که یادگیری $p(y|x)$ برای طبقه‌بندی معقول باشد تا خطای طبقه‌بندی $0 - 1$ کاهش یابد. یک جایگزین برای یادگیری مستقیم $p(y|x)$ این است که به ترتیب توزیع‌های شرطی کلاس و پیشین، $p(x|y)$ و $p(y)$ را یاد بگیرید. با استفاده از:

$$\begin{aligned} p(y|x) &= \frac{p(x, y)}{p(x)} \\ &= \frac{p(x|y)p(y)}{p(x)} \\ &\propto p(x|y)p(y) \end{aligned}$$

می‌توانیم ببینیم که این دو رویکرد یادگیری در تئوری معادل هستند تا \hat{y} با بالاترین $p(y|x)$ تعیین شود. انتخاب به دانش و/یا ترجیحات قبلی ما بستگی دارد. مدل‌هایی که با تخمین مستقیم $p(y|x)$ به دست می‌آیند، مدل‌های متمایز و مدل‌هایی که با تخمین مستقیم $p(x|y)$ و $p(y)$ به دست می‌آیند، مدل‌های مولد نامیده می‌شوند.

۳-۳-۴ خطای قابل کاهش و کاهش ناپذیر

پس از یافتن مدل رگرسیون بهینه، اکنون می‌توانیم هزینه مورد انتظار را در هر دو مدل بهینه و غیربهینه $f(x)$ بنویسیم. یعنی ما علاقه‌مندیم که $\mathbb{E}[C]$ را بیان کنیم وقتی که

$$\begin{aligned} ۱. \quad f(x) &= \mathbb{E}[Y | x] \\ ۲. \quad f(x) &\neq \mathbb{E}[Y | x] \end{aligned}$$

وقتی $f(x) = \mathbb{E}[Y | x]$ ، هزینه مورد انتظار را می‌توان به سادگی به این صورت بیان کرد

$$\begin{aligned} \mathbb{E}[C] &= \int_{\mathcal{X}} p(x) \int_{\mathcal{Y}} (\mathbb{E}[Y|x] - y)^2 p(y|x) dy dx \\ &= \int_{\mathcal{X}} p(x) V[Y|X=x] dx \end{aligned}$$

به یاد بیاورید که $V[Y | X = x]$ واریانس Y برای x داده شده است. بنابراین هزینه مورد انتظار منعکس کننده هزینه‌های ناشی از نویز یا تغییر در اهداف است. این بهترین سناریو در رگرسیون برای مجذور هزینه خطا است. ما نمی‌توانیم به هزینه کمتر از هزینه مورد انتظار دست یابیم.

موقعیت بعدی زمانی است که $f(x) \neq \mathbb{E}[Y | x]$ در اینجا، ما با تجزیه مربع خطا ادامه خواهیم داد

$$\begin{aligned} (f(x) - y)^2 &= (f(x) - \mathbb{E}[Y | x] + \mathbb{E}[Y | x] - y)^2 \\ &= (f(x) - \mathbb{E}[Y | x])^2 + \underbrace{2(f(x) - \mathbb{E}[Y | x])(\mathbb{E}[Y | x] - y)}_{g(x, y)} + (\mathbb{E}[Y | x] - y)^2 \end{aligned}$$

توجه داشته باشید که مقدار مورد انتظار $g(x, Y)$ برای هر x صفر است زیرا

$$\begin{aligned}
\mathbb{E}[g(\mathbf{x}, Y)] &= \mathbb{E}[(f(\mathbf{x}) - \mathbb{E}[Y|\mathbf{x}])(\mathbb{E}[Y|\mathbf{x}] - Y)|\mathbf{x}] \\
&= (f(\mathbf{x}) - \mathbb{E}[Y|\mathbf{x}])\mathbb{E}[(\mathbb{E}[Y|\mathbf{x}] - Y)|\mathbf{x}] \\
&= (f(\mathbf{x}) - \mathbb{E}[Y|\mathbf{x}]) (\mathbb{E}[Y|\mathbf{x}] - \mathbb{E}[Y|\mathbf{x}]) \\
&= 0
\end{aligned}$$

بنابراین، می‌توانیم نتیجه بگیریم که $\mathbb{E}[g(\mathbf{X}, Y)] = 0$ زمانی که انتظارات را بیش از \mathbf{X} می‌گیریم. اکنون می‌توانیم هزینه مورد انتظار را به صورت این بیان کنیم.

$$\begin{aligned}
\mathbb{E}[C] &= \mathbb{E}[(f(\mathbf{X}) - Y)^2] \\
&= \underbrace{\mathbb{E}[(f(\mathbf{X}) - \mathbb{E}[Y|\mathbf{X}])^2]}_{\text{reducible error}} + \underbrace{\mathbb{E}[(\mathbb{E}[Y|\mathbf{X}] - Y)^2]}_{\text{irreducible error}}
\end{aligned}$$

عبارت اول نشان می‌دهد که مدل آموزش دیده $f(\mathbf{x})$ چقدر از مدل بهینه $\mathbb{E}[Y|\mathbf{x}]$ فاصله دارد. عبارت دوم منعکس کننده تغییرپذیری ذاتی در Y با دادن \mathbf{x} است، همانطور که در معادله (۴،۴) نوشته شده است. این اصطلاحات اغلب خطاهای تقلیل پذیر و غیر قابل تقلیل نیز نامیده می‌شوند. اگر کلاس توابع f را برای پیش‌بینی $\mathbb{E}[Y|\mathbf{x}]$ گسترش دهیم، می‌توانیم اولین خطای مورد انتظار را کاهش دهیم. با این حال، خطای دوم ذاتی یا غیرقابل کاهش است به این معنا که هر چقدر عملکرد را بهبود بخشیم، نمی‌توانیم این عبارت را کاهش دهیم. این به مسئله مشاهده‌پذیری جزئی مربوط می‌شود، جایی که همیشه به دلیل کمبود اطلاعات مقداری تصادفی وجود دارد. این فاصله غیرقابل کاهش به طور بالقوه می‌تواند با ارائه اطلاعات ویژگی‌های بیشتر کاهش یابد (یعنی گسترش اطلاعات در \mathbf{X}). با این حال، برای یک مجموعه داده معین، با ویژگی‌های داده شده، این خطا غیر قابل کاهش است.

به طور خلاصه، ما در اینجا استدلال کردیم که مدل‌های طبقه‌بندی و رگرسیون بهینه به طور انتقادی به دانستن یا یادگیری دقیق توزیع پسین $p(y|\mathbf{x})$ بستگی دارد. این کار را می‌توان به روش‌های مختلفی حل کرد، اما یک رویکرد ساده این است که یک شکل تابعی برای $p(y|\mathbf{x})$ فرض کنیم، مثلاً $p(y|\mathbf{x}, \theta)$ ، که θ مجموعه‌ای از وزن‌ها یا پارامترهایی است که باید از داده‌ها یاد گرفت شوند.

۴-۴ [پیشرفته] مدل‌های بهینه بیز

قبلاً دیدیم که مدل‌های پیش‌بینی بهینه به یادگیری توزیع پسین $p(y|\mathbf{x})$ کاهش می‌یابد که سپس برای به مینیمم رساندن هزینه مورد انتظار (ریسک، هزینه) استفاده می‌شود. با این حال، در عمل، توزیع احتمال $p(y|\mathbf{x})$ باید با استفاده از یک فرم تابعی خاص و مجموعه‌ای از ضرایب قابل تنظیم مدل‌سازی شود. هنگامی که $\mathcal{X} = \mathbb{R}^d$ و $\mathcal{Y} = \{0, 1\}$ ، یکی از این مثالها در رگرسیون لجستیک استفاده می‌شود، که در آن

$$p(1|\mathbf{x}) = \frac{1}{1 + e^{-(w_0 + \sum_{j=1}^d w_j x_j)}}$$

و $p(0|\mathbf{x}) = 1 - p(1|\mathbf{x})$ در اینجا $(w_0, w_1, \dots, w_d) \in \mathbb{R}^d + 1$ مجموعه‌ای از وزن‌ها است که باید از مجموعه داده‌های داده شده \mathcal{D} استنباط شوند و $\mathbf{x} \in \mathbb{R}^d$ یک نقطه داده ورودی است. تعدادی از انواع دیگر روابط عملکردی نیز می‌تواند مورد استفاده قرار گیرد که مجموعه وسیعی از امکانات را برای مدل سازی توزیع‌ها فراهم می‌کند.

برای دقیق تر بودن در مورد این اشکال عملکردی، باید نماد خود را طوری تنظیم کنیم که توزیع بر روی \mathcal{Y} را با \mathcal{X} به عنوان نشان دهیم.

$$p(y|\mathbf{x}) = p(y|\mathbf{x}, f)$$

که در آن f یک تابع خاص از یک تابع (فرضیه) فضای \mathcal{F} است. ما می‌توانیم \mathcal{F} را به عنوان مجموعه‌ای از تمام توابع از یک کلاس مشخص، مثلاً برای همه $(w_0, w_1, \dots, w_d) \in \mathbb{R}^{k+1}$ در نظر بگیریم. در مثال بالا، اما ما همچنین می‌توانیم کلاس عملکردی را فراتر از تغییرات ساده پارامتر گسترش دهیم تا سطوح تصمیم‌گیری غیرخطی را بگنجانیم. ما معمولاً یک تابع را با توجه به داده‌ها انتخاب می‌کنیم - مثلاً ماکسیمم احتمال یا راه حل MAP

در عوض می‌توانیم توزیع $Y|x$ را روی همه توابع قابل قبول f در نظر بگیریم. در یک مسئله یادگیری معمولی، یک مجموعه داده $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ به ما داده می‌شود و از ما خواسته می‌شود که $p(y|x)$ را مدل کنیم. برای این منظور، \mathcal{D} را به عنوان تحقق یک متغیر تصادفی \mathcal{D} در نظر می‌گیریم و فرض می‌کنیم که \mathcal{D} مطابق با توزیع زیربنایی واقعی $p(x, y)$ ترسیم شده است. بنابراین وظیفه ما بیان $p(y|x, \mathcal{D})$ است. با استفاده از قواعد مجموع و محصول، وظیفه اصلی ما برای تخمین $p(y|x)$ به این صورت بازنویسی می‌شود.

$$\begin{aligned} p(y|x, \mathcal{D}) &= \int_{\mathcal{F}} p(y|x, f, \mathcal{D}) p(f|x, \mathcal{D}) df \\ &= \int_{\mathcal{F}} p(y|x, f) p(f|x, \mathcal{D}) df \end{aligned}$$

در اینجا ما از استقلال شرطی بین خروجی Y و مجموعه داده \mathcal{D} استفاده کردیم که یک مدل خاص f بر اساس \mathcal{D} انتخاب شد. بنابراین $p(y|x, f, \mathcal{D}) = p(y|x, f)$ این معادله به ما این احساس را می‌دهد که تصمیم بهینه را می‌توان از طریق مخلوطی از توزیع‌های $p(y|x, f)$ که در آن وزن‌ها به صورت چگالی پسینی $p(f|x, \mathcal{D})$ داده می‌شود، اتخاذ کرد. در فضاهای فرضی متناهی \mathcal{F} که داریم

$$p(y|x, \mathcal{D}) = \sum_{f \in \mathcal{F}} p(y|x, f) p(f|x, \mathcal{D})$$

و $p(f|x, \mathcal{D})$ احتمالات پسینی هستند. همچنین ممکن است فرض کنیم که $p(f|x, \mathcal{D}) = p(f|\mathcal{D})$ در این صورت وزن‌ها را می‌توان بر اساس مجموعه داده‌های \mathcal{D} از قبل محاسبه کرد. این منجر به محاسبات کارآمدتر احتمالات پسین می‌شود.

در طبقه‌بندی، می‌توانیم طبقه‌بندی کننده بهینه خود را به این صورت بازنویسی کنیم

$$f^*(x, \mathcal{D}) = \operatorname{argmax}_{y \in \mathcal{Y}} p(y|x, \mathcal{D})$$

که به آسانی به فرمول زیر منجر می‌شود

$$\begin{aligned} f^*(x, \mathcal{D}) &= \operatorname{argmax}_{y \in \mathcal{Y}} \int_{\mathcal{F}} p(y|x, f, \mathcal{D}) p(f|\mathcal{D}) df \\ &= \operatorname{argmax}_{y \in \mathcal{Y}} \int_{\mathcal{F}} p(y|x, f) p(\mathcal{D}|f) df \end{aligned}$$

می‌توان نشان داد که هیچ طبقه‌بندی کننده‌ای نمی‌تواند از طبقه‌بندی کننده بهینه بیز بهتر عمل کند. جالب توجه است، مدل بهینه بیز همچنین اشاره می‌کند که عملکرد پیش‌بینی بهتری را می‌توان با ترکیب چندین مدل و میانگین‌گیری خروجی‌های آنها به دست آورد. این امر پشتیبانی تئوریک را برای یادگیری گروه و روش‌هایی مانند بسته‌بندی و تقویت فراهم می‌کند.

با توجه به اینکه تابع (فرضیه) فضای \mathcal{F} به طور کلی غیرقابل شمارش است، یک مسئله در طبقه‌بندی بهینه بیز، محاسبه کارآمد $f^*(x, \mathcal{D})$ است. یک روش برای این کار نمونه برداری از توابع از \mathcal{F} با توجه به $p(f)$ و سپس محاسبه $p(f|\mathcal{D})$ یا $p(\mathcal{D}|f)$ است. این را می‌توان تا زمانی محاسبه کرد که $p(y|x, \mathcal{D})$ همگرا شود.

فصل ۵

رگرسیون خطی

با توجه به مجموعه داده $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ هدف یادگیری رابطه بین ویژگی‌ها و هدف است. ما معمولاً با فرضیه‌سازی شکل عملکردی این رابطه شروع می‌کنیم. مثلاً،

$$f(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2$$

که در آن $w = (w_0, w_1, w_2)$ مجموعه‌ای از پارامترهایی است که باید تعیین شوند (آموخته شوند) و $\mathbf{x} = (x_1, x_2)$ از طرف دیگر، ممکن است فرض کنیم که $f(x) = \alpha + \beta x_1 x_2$ ، که در آن $\theta = (\alpha, \beta)$ ، مجموعه دیگری از پارامترها است که باید آموخته شود. در مورد اول، تابع هدف به عنوان ترکیبی خطی از ویژگی‌ها و پارامترها مدل می‌شود، به عنوان مثال.

$$f(x) = \sum_{j=0}^d w_j x_j$$

که \mathbf{x} را به صورت $(x_0 = 1, x_1, x_2, \dots, x_d)$ گسترش دادیم. یافتن بهترین پارامترهای w به عنوان مسئله رگرسیون خطی یاد می‌شود، در حالی که همه انواع دیگر روابط بین ویژگی‌ها و هدف در دسته‌ای از رگرسیون غیر خطی قرار می‌گیرند. در هر یک از شرایط، مسئله رگرسیون را می‌توان به عنوان یک رویکرد مدل‌سازی احتمالی ارائه کرد که به تخمین پارامتر کاهش می‌یابد: به یک مسئله بهینه‌سازی با هدف به ماکسیمم رساندن یا به مینیمم رساندن برخی از معیارهای عملکرد بین مقادیر هدف $\{y_i\}_{i=1}^n$ و پیش‌بینی‌ها $\{f(x_i)\}_{i=1}^n$ ما می‌توانیم یک الگوریتم بهینه‌سازی خاص را به عنوان الگوریتم یادگیری یا آموزش در نظر بگیریم.

۵-۱ فرمول ماکسیمم احتمال

اکنون فرمول آماری رگرسیون خطی را در نظر می‌گیریم. ابتدا مفروضات پشت این فرآیند را بیان می‌کنیم و متعاقباً مسئله را از طریق ماکسیمم کردن تابع احتمال شرطی فرموله می‌کنیم. در بخش بعدی نحوه حل بهینه‌سازی و تجزیه و تحلیل راه‌حل و ویژگی‌های آماری اساسی آن را نشان خواهیم داد.

فرض کنید که مجموعه داده مشاهده شده \mathcal{D} محصول یک فرآیند تولید داده است که در آن n نقطه داده به طور مستقل و بر اساس توزیع یکسان $p(\mathbf{x})$ ترسیم شده است. همچنین فرض کنید که متغیر هدف Y یک رابطه خطی زیربنایی با ویژگی‌های $X = (X_1, X_2, \dots, X_d)$ دارد که با مقداری خطای ε اصلاح شده است که از توزیع گاوسی میانگین صفر پیروی می‌کند،

یعنی $\varepsilon: \mathcal{N}(0, \sigma^2)$. یعنی برای یک ورودی \mathbf{x} هدف y تحقق یک متغیر تصادفی Y است که به این صورت تعریف شده است.

$$Y = \sum_{j=0}^d \omega_j X_j + \varepsilon$$

که در آن $\boldsymbol{\omega} = (\omega_0, \omega_1, \dots, \omega_d)$ مجموعه‌ای از ضرایب ناشناخته است که ما به دنبال بازیابی آن از طریق تخمین هستیم. به طور کلی، فرض نرمال بودن عبارت خطا معقول است (قضیه حد مرکزی را به یاد بیاورید)، اگرچه ممکن است استقلال بین \mathbf{x} و ε در عمل برقرار نباشد. با استفاده از چند ویژگی ساده انتظارات، می‌توانیم ببینیم که Y نیز از توزیع گاوسی پیروی می‌کند، یعنی چگالی شرطی آن $p(y|\mathbf{x}, \boldsymbol{\omega}) = \mathcal{N}(\boldsymbol{\omega}^T \mathbf{x}, \sigma^2)$ است.

در رگرسیون خطی، ما به دنبال تقریب هدف به صورت $f(x) = \mathbf{w}^T \mathbf{x}$ هستیم، جایی که وزن‌های \mathbf{w} تعیین می‌شوند. ابتدا تابع درستنمایی شرطی را برای یک جفت واحد (\mathbf{x}, y) به این صورت می‌نویسیم

$$p(y|\mathbf{x}, \boldsymbol{\omega}) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(y - \sum_{j=0}^d \omega_j x_j)^2}{2\sigma^2}\right)$$

که در آن از نماد $\exp(a) = e^a$ استفاده می‌کنیم تا خوانش توان را آسان‌تر کنیم. توجه کنید که تنها تغییر تابع چگالی شرطی Y این است که به جای $\boldsymbol{\omega}$ از ضرایب \mathbf{w} استفاده می‌شود. با ترکیب کل مجموعه داده $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ اکنون می‌توانیم تابع درستنمایی شرطی را به صورت $p(\mathbf{y}|\mathbf{X}, \mathbf{w})$ بنویسیم و وزن‌ها را به صورت

$$\mathbf{w}_{\text{MLE}} = \arg \max_{\mathbf{w}} \{p(\mathbf{y}|\mathbf{X}, \mathbf{w})\}.$$

از آنجایی که n نمونه مستقل هستند و به طور یکسان توزیع شده‌اند ($i.i.d.$)، ما داریم

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \mathbf{w}) &= \prod_{i=1}^n p(y_i|\mathbf{x}_i, \mathbf{w}) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(y_i - \sum_{j=0}^d \omega_j x_{ij})^2}{2\sigma^2}\right) \end{aligned}$$

به دلایل راحتی ریاضی، به لگاریتم (تابع یکنواخت) تابع درستنمایی نگاه می‌کنیم و لگاریتم احتمال را به این صورت بیان می‌کنیم.

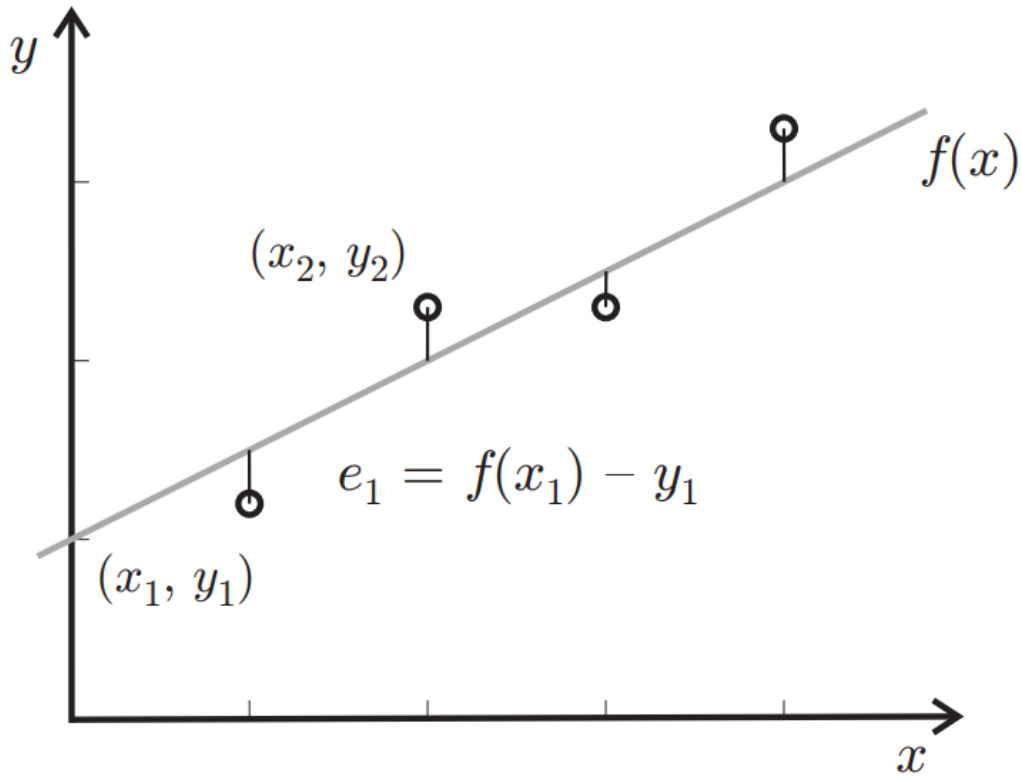
$$\ln(p(\mathbf{y}|\mathbf{X}, \mathbf{w})) = -\sum_{i=1}^n \log(\sqrt{2\pi} \sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^n \left(y_i - \sum_{j=0}^d \omega_j x_{ij}\right)^2$$

با توجه به اینکه جمله اول در سمت راست مستقل از \mathbf{w} است، ماکسیمم کردن تابع درستنمایی دقیقاً با مینیمم کردن مجموع مجذور خطاها مطابقت دارد.

$$Err(\mathbf{w}) = \sum_{i=1}^n (f(x_i) - y_i)^2 \quad \triangleright \quad f(x_i) = \sum_{j=0}^d \omega_j x_{ij} = \hat{y}_i$$

از نظر هندسی، این خطا مربع فاصله اقلیدسی بین بردار پیش‌بینی‌های $\hat{\mathbf{y}} = (f(x_1), f(x_2), \dots, f(x_n))$ و بردار مقادیر هدف مشاهده شده $\mathbf{y} = (y_1, y_2, \dots, y_n)$ یک مثال ساده که مسئله رگرسیون خطی را نشان می‌دهد در شکل ۵,۱ نشان داده شده است.

برای اینکه واضح‌تر ببینید که چرا راه‌حل ماکسیمم درست‌نمایی با کمینه کردن $Err(w)$ مطابقت دارد، توجه کنید که به ماکسیمم رساندن احتمال معادل به ماکسیمم رساندن $\log - \text{relihood}$ است.



شکل ۵,۱: یک راه‌حل رگرسیون خطی در مجموعه داده $\mathcal{D} = \{(1, 1.2), (2, 2.3), (3, 2.3), (4, 3.3)\}$ وظیفه فرآیند بهینه‌سازی یافتن بهترین تابع خطی $f(x) = \omega_0 + \omega_1 x$ است به طوری که مجموع مربعات خطاهای $e_1^2 + e_2^2 + e_3^2 + e_4^2$ به مینیمم برسد.

(چون \log یکنواخت است) که معادل به مینیمم رساندن احتمال $\log - \text{likelihood}$ منفی است. بنابراین، ماکسیمم احتمال \mathbf{w}_{MLE} می‌شود

$$\begin{aligned} \mathbf{w}_{MLE} &= \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} -\ln(p(\mathbf{y}|\mathbf{X}, \mathbf{w})) \\ &= \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{i=1}^n \log(\sqrt{2\pi\sigma^2}) + \frac{1}{2\sigma^2} \sum_{i=1}^n \left(y_i - \sum_{j=0}^d w_j x_{ij} \right)^2 \\ &= \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{i=1}^n \left(y_i - \sum_{j=0}^d w_j x_{ij} \right)^2 \end{aligned}$$

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} Err(\mathbf{w})$$

در بخش‌های بعدی به چگونگی حل این بهینه‌سازی و خواص راه‌حل می‌پردازیم.

توجه داشته باشید که ما می‌توانستیم به سادگی با برخی از توابع خطا (تعریف شده توسط متخصص) شروع کنیم، همانطور که در ابتدا برای OLS و با استفاده از $Err(w)$ انجام شد. با این حال، چارچوب آماری بینش‌هایی را در مورد مفروضات پشت رگرسیون OLS ارائه می‌دهد. به طور خاص، مفروضات شامل این است که داده D به صورت $i.i.d.$ ترسیم شده است. یک رابطه خطی اساسی بین ویژگی‌ها و هدف وجود دارد. که نویز (اصطلاح خطا) گاوسی صفر و مستقل از ویژگی‌ها و عدم وجود نویز در مجموعه ویژگی‌ها است.

۵-۲ رگرسیون مینیمم مربعات معمولی^۱ (OLS).

برای به مینیمم رساندن مجموع مجذور خطاها، ابتدا $Err(\mathbf{w})$ را دوباره می‌نویسیم

$$\begin{aligned} Err(\mathbf{w}) &= \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2 \\ &= \sum_{i=1}^n \left(\sum_{j=0}^d w_j x_{ij} - y_i \right)^2 \end{aligned}$$

جایی که، دوباره، هر نقطه داده \mathbf{x}_i را با $x_{i0} = 1$ گسترش دادیم تا عبارت را ساده کنیم.

اکنون گرادیان $\nabla Err(w)$ را محاسبه می‌کنیم. یافتن وزن‌هایی که $\nabla Err(w) = 0$ باعث ایجاد یک نقطه ثابت می‌شود. برای اطمینان از اینکه این نقطه ثابت یک مینیمم مطلق است، به اطلاعات بیشتری نیاز داریم. می‌توانیم به مشتق دوم نگاه کنیم. این مستلزم درک هسین^۲ است، بنابراین ما آن را بعداً در یادداشتهای مثال ۱۶ لحاظ می‌کنیم. اما خوشبختانه، در اینجا حتی ساده‌تر است، زیرا می‌دانیم که این هدف در \mathbf{w} محدب است. بنابراین، هر نقطه ثابت یک مینیمم مطلق خواهد بود.

اکنون مشتقات جزئی را برابر 0 قرار می‌دهیم و معادلات هر وزن w_j را حل می‌کنیم

$$\begin{aligned} \frac{\partial Err}{\partial w_0} &= 2 \sum_{i=1}^n \left(\sum_{j=0}^d w_j x_{ij} - y_i \right) x_{i0} = 0 \\ \frac{\partial Err}{\partial w_1} &= 2 \sum_{i=1}^n \left(\sum_{j=0}^d w_j x_{ij} - y_i \right) x_{i1} = 0 \\ &\vdots \end{aligned}$$

¹ Ordinary Least-Squares

² Hessian

$$\frac{\partial \text{Err}}{\partial w_d} = 2 \sum_{i=1}^n \left(\sum_{j=0}^d w_j x_{ij} - y_i \right) x_{id} = 0$$

این منجر به سیستمی از معادلات خطی $d + 1$ با مجهولات $d + 1$ می‌شود که می‌توانند به طور معمول حل شوند (به عنوان مثال با استفاده از حذف گاوسی)

در حالی که این فرمول‌بندی مفید است، به ما اجازه نمی‌دهد که یک راه‌حل به شکل بسته برای \mathbf{w} بدست آوریم یا در مورد وجود یا تعداد این ترکیبات بحث کنیم. برای پرداختن به اولین نکته، مقداری محاسبات ماتریسی را اعمال خواهیم کرد، در حالی که بقیه نکات بعداً مورد بحث قرار خواهند گرفت. ابتدا مجموع مربعات خطاها را با استفاده از نماد ماتریس به این صورت می‌نویسیم

$$\text{Err}(\mathbf{w}) = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

که در آن $\|\mathbf{v}\|_2 = \sqrt{\mathbf{v}^T \mathbf{v}} = \sqrt{v_1^2 + v_2^2 + \dots + v_3^2}$ طول بردار \mathbf{v} است. به آن نرم ℓ_2 نیز می‌گویند. اکنون می‌توانیم مسئله رگرسیون خطی مینیمم مربعات معمولی (OLS) را به این صورت فرمول‌بندی کنیم

$$\mathbf{w}_{MLE} = \arg \min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

با یافتن $\nabla \text{Err}(\mathbf{w})$ کار را ادامه می‌دهیم. تابع گرادیان $\nabla \text{Err}(\mathbf{w})$ مشتقی از یک اسکالر با توجه به یک بردار است. با این حال، مراحل میانی محاسبه گرادیان به مشتقات بردارها با در نظر گرفتن بردارها نیاز دارد (برخی از قوانین چنین مشتقاتی در جدول B.1 نشان داده شده است). به کارگیری قوانین جدول B.1 نتیجه می‌دهد

$$\nabla \text{Err}(\mathbf{w}) = 2\mathbf{X}^T \mathbf{X}\mathbf{w} - 2\mathbf{X}^T \mathbf{y}$$

و بنابراین، از $\nabla \text{Err}(\mathbf{w}) = 0$ متوجه می‌شویم که

$$\mathbf{w}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

اکنون می‌توانیم مقادیر هدف پیش‌بینی شده را به این صورت بیان کنیم

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}_{MLE} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

ماتریس $\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ ماتریس طرح‌ریزی نامیده می‌شود. به بخش C.1 مراجعه کنید تا بفهمید چگونه \mathbf{y} را به فضای ستون \mathbf{X} می‌فرستد

مثال ۱۳: دوباره مجموعه داده $\mathcal{D} = \{(1, 1.2), (2, 2.3), (3, 2.3), (4, 3.3)\}$ را از شکل ۵.۱ در نظر بگیرید. ما می‌خواهیم ضرائب بهینه مینیمم مربعات متناسب با $f(x) = w_0 + w_1 x$ را پیدا کنیم و سپس مجموع مربعات خطاهای \mathcal{D} را پس از برازش محاسبه کنیم.

$$\mathbf{x} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix}, \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}, \mathbf{y} = \begin{bmatrix} 1.2 \\ 2.3 \\ 2.3 \\ 3.3 \end{bmatrix}$$

که در آن یک ستون از یک‌ها به \mathbf{x} اضافه شد تا یک وقفه غیر صفر مجاز باشد ($y = w_0$ وقتی $x = 0$). جایگزینی \mathbf{x} و \mathbf{y} به معادله (۵،۱) به $w = (0.7, 0.63)$ می‌رسد و مجموع خطاهای مربع برابر $\text{Err}(\mathbf{w}) = 0.223$ است.

همانطور که در مثال بالا مشاهده شد، اضافه کردن ستونی از یک‌ها به ماتریس داده \mathbf{X} به منظور اطمینان از اینکه خط برازش شده، یا به طور کلی یک ابر صفحه، لازم نیست از مبدأ سیستم مختصات عبور کند، یک روش استاندارد است. اما این تأثیر را می‌توان از راه‌های دیگری نیز به دست آورد. اولین جزء بردار گرادیان را در نظر بگیرید

$$\frac{\partial \text{Err}}{\partial w_0} = 2 \sum_{i=1}^n \left(\sum_{j=0}^d w_j x_{ij} - y_i \right) x_{i0} = 0$$

جایی که طبق تعریف، چون $x_{i0} = 1$ ، آن را اینگونه بدست می‌آوریم

$$0 = \sum_{i=1}^n \left(\sum_{j=0}^d w_j x_{ij} - y_i \right) = \sum_{i=1}^n \left(w_0 + \sum_{j=1}^d w_j x_{ij} - y_i \right)$$

که می‌دهد

$$\sum_{i=1}^n w_0 = \sum_{i=1}^n y_i - \sum_{j=1}^d w_j \sum_{i=1}^n x_{ij}$$

وقتی همه ویژگی‌ها (ستون‌های \mathbf{X}) نرمال می‌شوند تا میانگین صفر داشته باشند، یعنی وقتی $\sum_{i=1}^n x_{ij} = 0$ برای هر ستون j ، نتیجه می‌شود که

$$w_0 = \frac{1}{n} \sum_{i=1}^n y_i$$

اکنون می‌بینیم که اگر متغیر هدف به میانگین صفر نیز نرمال شود، نتیجه آن این است که $w_0 = 0$ و ستون یک‌ها مورد نیاز نیست.

۵-۲-۱ تابع خطای وزنی

در برخی از کاربردها، به مینیمم رساندن تابع خطای وزنی مفید است

$$\text{Err}(\mathbf{w}) = \sum_{i=1}^n c_i \left(\sum_{j=0}^d w_j x_{ij} - y_i \right)^2$$

جایی که $c_i > 0$ هزینه‌ای برای نقطه داده i است. با بیان آن به صورت ماتریسی، هدف این است که $(\mathbf{X}\mathbf{w} - \mathbf{y})^T \mathbf{C} (\mathbf{X}\mathbf{w} - \mathbf{y})$ را به مینیمم برسانیم، جایی که $\mathbf{C} = \text{diag}(c_1, c_2, \dots, c_n)$. با استفاده از روشی مشابه با بالا، می‌توان نشان داد که راه‌حل مینیمم مربعات وزنی \mathbf{w}_C را می‌توان به این صورت بیان کرد.

$$\mathbf{w}_C = (\mathbf{X}^T \mathbf{C} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{C} \mathbf{y}$$

علاوه بر این، می‌توان این را استخراج کرد که

$$\mathbf{w}_C = \mathbf{w}_{MLE} + (\mathbf{X}^T \mathbf{C} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{1} - \mathbf{C}) (\mathbf{X} \mathbf{w}_{MLE} - \mathbf{y})$$

جایی که \mathbf{w}_{MLE} توسط معادله ارائه شده است. (۵,۱). می‌توانیم ببینیم که راه‌حل‌ها زمانی که $\mathbf{C} = \mathbf{I}$ ، و همچنین زمانی که $\mathbf{X} \mathbf{w}_{MLE} = \mathbf{y}$ یکسان هستند، یکسان هستند.

۵-۲-۲ پیش‌بینی چندین خروجی به طور همزمان

گسترش چندین خروجی ساده است، جایی که اکنون هدف یک بردار m بعدی، $\mathbf{y} \in \mathbb{R}^m$ است، نه یک اسکالر، که ماتریس هدف $\mathbf{Y} \in \mathbb{R}^{n \times m}$ را می‌دهد. به همین ترتیب، وزن‌های $\mathbf{W} \in \mathbb{R}^{d \times m}$ برای دادن $\mathbf{W}^T \mathbf{x} \in \mathbb{R}^m$ با خطا

$$\begin{aligned} \text{Err}(\mathbf{W}) &= \|\mathbf{X}\mathbf{W} - \mathbf{Y}\|_F^2 = \sum_{i=1}^n \|\mathbf{X}_{i,:} \mathbf{W} - \mathbf{Y}_{i,:}\|_2^2 &> \text{Frobenius norm} \\ &= \text{trace} \left((\mathbf{X}\mathbf{W} - \mathbf{Y})^T (\mathbf{X}\mathbf{W} - \mathbf{Y}) \right) \end{aligned}$$

و راه‌حل آن

$$\mathbf{w}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

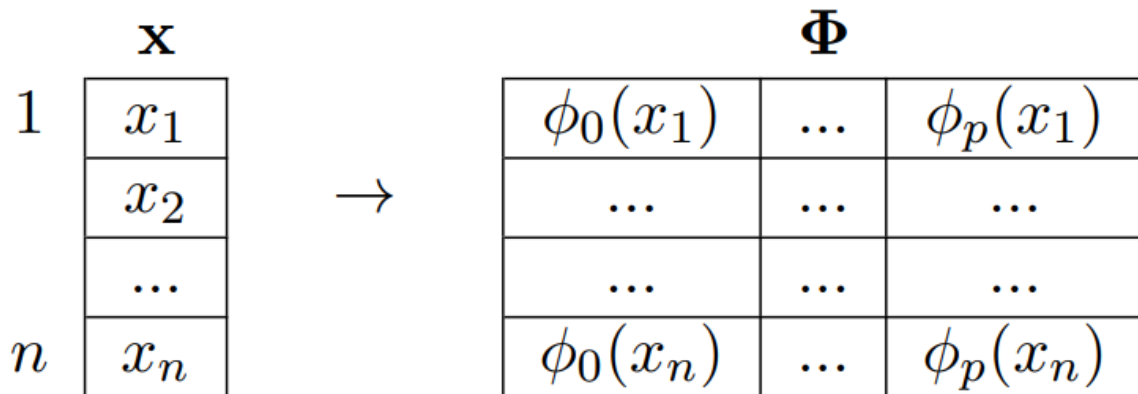
تمرین ۶: این راه‌حل را با گرفتن مشتقات جزئی یا ترجیحاً با استفاده از قوانین گرادیان برای متغیرهای ماتریس استخراج کنید. یک منبع خوب برای گرادیان‌های ماتریس، کتابچه راهنما ماتریس^۱ [۱۶] است.

برای به دست آوردن بینش بیشتر در مورد راه‌حل رگرسیون خطی معمولی، چشم انداز جبری را در پیوست C.1 ببینید. بینش بیشتری در مورد منحصر به فرد بودن راه‌حل و فضای راه‌حل‌های ممکن ارائه می‌دهد و بهینه‌سازی رگرسیون خطی را به سیستم‌های حل متصل می‌کند.

۵-۳ رگرسیون خطی برای مسائل غیر خطی

در ابتدا، ممکن است به نظر برسد که کاربرد رگرسیون خطی برای مسائل زندگی واقعی بسیار محدود است. در کل، مشخص نیست که آیا آن را واقع‌بینانه (بیشتر اوقات) فرض کنیم

¹ matrix cookbook



شکل ۵/۲: تبدیل یک ماتریس داده x با اندازه $n \times 1$ به یک ماتریس Φ با اندازه $n \times (p + 1)$ با استفاده از مجموعه ای از توابع پایه ϕ_j و $j = 0, 1, \dots, p$

که متغیر هدف ترکیبی خطی از ویژگی‌ها است. خوشبختانه، کاربرد رگرسیون خطی گسترده‌تر است زیرا می‌توانیم از آن برای به دست آوردن توابع غیر خطی استفاده کنیم. ایده اصلی این است که قبل از مرحله برازش، یک تبدیل غیر خطی به ماتریس داده **X** اعمال شود، که سپس یک برازش غیر خطی را امکان پذیر می‌کند. به دست آوردن چنین نمایش ویژگی مفیدی یک مشکل اصلی در یادگیری ماشین است. ما در این مورد در فصل ۹ به تفصیل بحث خواهیم کرد. در اینجا، ابتدا یک نمایش ساده‌تر را بررسی خواهیم کرد که یادگیری غیرخطی را امکان پذیر می‌کند: برازش منحنی چند جمله‌ای^۱

۱-۳-۵ برازش منحنی چند جمله‌ای

ما با داده‌های یک بعدی شروع می‌کنیم. در رگرسیون *OLS*، ما به دنبال تناسب در شکل زیر هستیم

$$f(x) = w_0 + w_1 x$$

که در آن x نقطه داده و $\mathbf{w} = (w_0, w_1)$ بردار وزن است. برای دستیابی به تناسب چند جمله‌ای درجه p ، عبارت قبلی را به آن تغییر می‌دهیم

$$f(x) = \sum_{j=0}^p w_j x^j$$

که در آن p درجه چند جمله‌ای است. ما این عبارت را با استفاده از مجموعه‌ای از توابع پایه بازنویسی می‌کنیم

$$f(x) = \sum_{j=0}^p w_j \phi^j(x) = \mathbf{w}^T \boldsymbol{\phi}$$

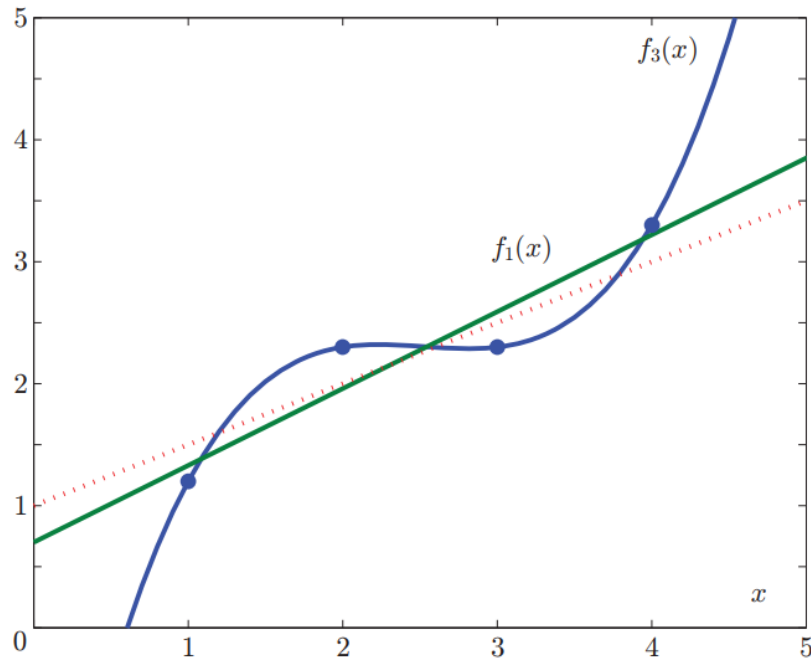
که در آن $\phi_j(x) = x^j$ و $\boldsymbol{\phi} = (\phi_0(x), \phi_1(x), \dots, \phi_p(x))$. اعمال این تبدیل برای هر نقطه داده در x منجر به ماتریس داده‌ای جدید Φ می‌شود، همانطور که در شکل ۵/۲ نشان داده شده است.

^۱ Polynomial curve fitting

در ادامه بحث از بخش ۵,۲، مجموعه بهینه وزن‌ها به صورت محاسبه می‌شود

$$\mathbf{w}_{MLE} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

مثال ۱۴: در شکل ۵,۱ نمونه‌ای از یک مجموعه داده با چهار نقطه داده را ارائه کردیم. چیزی که ما ذکر نکردیم این بود که با توجه به مجموعه‌ای $\{x_1, x_2, x_3, x_4\}$ اهداف تولید شدند.



شکل ۵/۳: نمونه‌ای از تناسب خطی در مقابل چند جمله‌ای در مجموعه داده نشان داده شده در شکل ۵/۱. تناسب خطی، $f_1(x)$ ، به عنوان یک خط سبز ثابت نشان داده شده است، در حالی که تناسب چند جمله‌ای مکعبی، $f_3(x)$ ، به عنوان یک خط آبی ثابت نشان داده شده است. خط قرمز نقطه‌چین مفهوم خطی هدف را نشان می‌دهد.

با استفاده از تابع $1 + \frac{x}{2}$ و سپس اضافه کردن یک خطای اندازه‌گیری $e = (-0.3, 0.3, -0.2, 0.3)$ مشخص شد که ضرائب بهینه $\mathbf{w}_{MLE} = (0.7, 0.63)$ به ضرائب واقعی $\omega = (1, 0.5)$ نزدیک است، حتی اگر عبارات خطا نسبتاً معنی‌دار بودند. اکنون سعی خواهیم کرد ضرائب یک تناسب چند جمله‌ای را با درجات $p = 2$ و $p = 3$ تخمین بزنیم. همچنین مجموع خطاهای مجذور \mathcal{D} را پس از برازش و همچنین روی یک مجموعه گسسته بزرگ از مقادیر $x \in \{0, 0.1, 0.2, \dots, 10\}$ محاسبه خواهیم کرد. که در آن مقادیر هدف با استفاده از تابع واقعی $1 + \frac{x}{2}$ تولید می‌شوند.

استفاده از یک برازش چند جمله‌ای با درجات $p = 2$ و $p = 3$ به ترتیب به $\mathbf{w}_3 = (0.575, 0.755, -0.025)$ و $\mathbf{w}_2 = (-3.1, 6.6, -2.65, 0.35)$ منجر می‌شود. مجموع مجذور خطاها در \mathcal{D} برابر است با $\text{Err}(\mathbf{w}_2) = 0.221$ و $\text{Err}(\mathbf{w}_3) \approx 0$. بنابراین، بهترین تناسب با چند جمله‌ای مکعبی به دست می‌آید. با این حال، مجموع خطاهای مجذور در مجموعه داده‌های بیرونی توانایی تعمیم ضعیف مدل مکعبی را نشان می‌دهد زیرا ما $\text{Err}(\mathbf{w}) = 26.9$ ، $\text{Err}(\mathbf{w}_2) = 3.9$ ، و $\text{Err}(\mathbf{w}_3) = 22018.5$ به دست می‌آوریم. این اثر بیش از حد برازش^۱ نامیده می‌شود. به طور کلی، برازش بیش از حد با تفاوت قابل توجهی در تناسب بین مجموعه داده‌ای که مدل بر روی آن آموزش داده شده است و مجموعه داده‌های خارجی که انتظار می‌رود مدل بر روی آن اعمال شود نشان داده می‌شود (شکل ۵,۳). در این مورد، بیش از حد برازش رخ داد زیرا پیچیدگی مدل به طور قابل توجهی افزایش یافت، در حالی که اندازه مجموعه داده‌ها کوچک باقی ماند.

^۱ overfitting

یکی از نشانه‌های برازش بیش از حد، افزایش بزرگی ضرائب است. به عنوان مثال، در حالی که مقادیر مطلق همه ضرایب در \mathbf{w} و \mathbf{w}_2 کمتر از یک بود، مقادیر ضرایب در \mathbf{w}_3 با علائم متناوب به طور قابل توجهی بزرگ‌تر شدند (که نشان دهنده جبران بیش از حد است). در بخش ۵,۴,۲ به عنوان رویکردی برای جلوگیری از این اثر، قانون گذاری را مورد بحث قرار خواهیم داد.

برازش منحنی چند جمله‌ای تنها یکی از راه‌های برازش غیرخطی است، زیرا انتخاب توابع پایه نباید به توان‌های X محدود شود. از جمله توابع پایه غیر خطی که معمولاً مورد استفاده قرار می‌گیرند تابع سیگموئید^۱ هستند.

$$\varphi_j(x) = \frac{1}{1 + e^{-\frac{x - \mu_j}{s_j}}}$$

یا تابع نمایی به سبک گاوسی

$$\varphi_j(x) = e^{-\frac{(x - \mu_j)^2}{2\sigma_j^2}}$$

که در آن μ_j ، s_j و σ_j ثابت‌هایی هستند که باید تعیین شوند. با این حال، این رویکرد فقط برای یک ورودی یک بعدی X کار می‌کند. برای ابعاد بالاتر، این رویکرد را می‌توان با استفاده از توابع پایه شعاعی تعمیم داد. برای جزئیات بیشتر به بخش ۹,۱ مراجعه کنید.

۴-۵ ثبات و مبادله^۲ بایاس واریانس^۳

راه حل OLS می‌تواند ناپایدار باشد. در این بخش، ما نشان می‌دهیم که چرا، و بحث می‌کنیم که چگونه می‌توان از منظم‌سازی برای کاهش این مشکل استفاده کرد. سپس یک مفهوم اساسی در یادگیری ماشین را مورد بحث قرار خواهیم داد: مبادله بایاس واریانس.

۱-۴-۵ حساسیت راه حل OLS

راه حل OLS ناپایدار است اگر $\mathbf{X}^T \mathbf{X}$ معکوس پذیر نباشد. این می‌تواند به دو دلیل اصلی رخ دهد: ویژگی‌های وابسته به خطی و مجموعه داده‌های کوچک. مجموعه داده‌ها اغلب شامل تعداد زیادی ویژگی است که گاهی یکسان، مشابه یا تقریباً خطی وابسته هستند. اگر مجموعه داده کوچک باشد، ممکن است برخی از ویژگی‌ها در نمونه‌ها یکسان باشند، و دوباره منجر به رتبه پایین \mathbf{X} می‌شود. وقتی $\mathbf{X}^T \mathbf{X}$ معکوس پذیر نیست یا شرایط نامناسبی دارد، راه حل OLS به اغتشاش‌های کوچک در \mathbf{y} و \mathbf{X} بسیار حساس است.

برای اینکه بفهمیم چرا، به تجزیه مقدار یکتای \mathbf{X} نگاه خواهیم کرد. مانند ساختارهای جبر خطی قبلی، به ما این امکان را می‌دهد که به راحتی ویژگی‌های \mathbf{X} را بررسی کنیم. بیایید حالت رانج را در نظر بگیریم، که در آن $n > d$: تعداد نمونه‌ها بیشتر است. نسبت به بعد ورودی تجزیه مقدار یکتا $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ برای ماتریس‌های متعامد $\mathbf{U} \in \mathbb{R}^{n \times n}$ ، $\mathbf{V} \in \mathbb{R}^{d \times d}$

^۱ sigmoid function

^۲ trade-of

^۳ bias-variance

و ماتریس قطری غیر منفی (مستطیل شکل) $\Sigma \in \mathbb{R}^{n \times d}$ ورودی‌های قطری در Σ مقادیر یکتای هستند که معمولاً آنها را به ترتیب نزولی $\sigma_1, \sigma_2, \dots, \sigma_d$ مرتب می‌کنیم. که میدهد:

$$\Sigma \stackrel{\text{def}}{=} \begin{bmatrix} \sigma_1 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & 0 & \dots & 0 \\ & & \vdots & & \\ 0 & 0 & \dots & 0 & \sigma_d \\ 0 & 0 & \dots & 0 & 0 \\ & & \vdots & (n-d) \text{ rows of zeros} & \\ 0 & 0 & \dots & 0 & 0 \end{bmatrix} = \begin{bmatrix} \Sigma_d \\ 0 \end{bmatrix} \text{ where } \Sigma_d \stackrel{\text{def}}{=} \begin{bmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ & & \vdots & \\ 0 & 0 & \dots & \sigma_d \end{bmatrix}$$

هر ماتریس $\mathbf{X} \in \mathbb{R}^{n \times d}$ را می‌توان به تجزیه مقدار یکتای آن تجزیه کرد، زیرا هر تبدیل خطی را می‌توان به یک چرخش (ضرب در \mathbf{V}^T) تجزیه کرد، به دنبال آن یک مقیاس (ضرب در Σ) و به دنبال آن دوباره چرخش (ضرب در \mathbf{U}).

این تجزیه، تجزیه و تحلیل خواص یک ماتریس را ساده می‌کند. برای مثال، تعداد مقادیر غیر صفر یکتا، رتبه \mathbf{X} را تشکیل می‌دهد. برای اینکه ببینید چرا، $\sigma_d = 0$ و $\sigma_{d-1} > 0$ را فرض کنید، به این معنی که \mathbf{X} دارای رتبه $d - 1$ است. هر بردار $\mathbf{w} \in \mathbb{R}^d$ را در نظر بگیرید، و $\mathbf{X}\mathbf{w}$ را در نظر بگیرید. ما می‌توانیم این محصول را به صورت $\mathbf{U}\Sigma\mathbf{V}^T\mathbf{w} = \mathbf{U}\Sigma\tilde{\mathbf{w}}$ برای $\tilde{\mathbf{w}} = \mathbf{V}^T\mathbf{w}$ بنویسیم. محصول $\Sigma\tilde{\mathbf{w}}$ را تنظیم می‌کند

آخرین بعد $\tilde{\mathbf{w}}$ به صفر، به طور موثر آن بعد را حذف می‌کند و بنابراین $\tilde{\mathbf{w}}$ را به فضایی با ابعاد پایین‌تر ($d - 1$) نشان می‌دهد. سپس آن بردار پیش‌بینی شده را با استفاده از \mathbf{U} می‌چرخاند، اما نمی‌تواند آن پیش‌بینی را در فضایی با ابعاد پایین‌تر خنثی کند. بنابراین، $\mathbf{X}\mathbf{w}$ فقط می‌تواند $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$ را حاصل کند که در یک صفحه $d - 1$ بعدی قرار دارد که در \mathbb{R}^{d-1} چرخیده است. بنابراین، این تجزیه می‌تواند به ما در درک فضای پیش‌بینی‌های ممکن برای رگرسیون خطی $\mathbf{X}\mathbf{w}$ کمک کند.

اکنون می‌توانیم راه‌حل مینیمم مربعات را بر حسب تجزیه مقدار یکتای \mathbf{X} مورد بحث قرار دهیم. توجه کنید که

$$\mathbf{X}^T\mathbf{X} = \mathbf{V}\Sigma^T\mathbf{U}^T\mathbf{U}\Sigma\mathbf{V}^T = \mathbf{V}\Sigma_d^2\mathbf{V}^T$$

زیرا \mathbf{U} متعامد است و بنابراین $\mathbf{U}^T\mathbf{U} = \mathbf{I}$ ماتریس همانی است (\mathbf{I} یک ماتریس قطری با یک‌هایی در قطر است). معکوس $\mathbf{X}^T\mathbf{X}$ وجود دارد اگر \mathbf{X} رتبه کامل باشد، یعنی Σ_d هیچ صفری در قطر نداشته باشد، زیرا $(\mathbf{X}^T\mathbf{X})^{-1} = \mathbf{V}\Sigma_d^{-2}\mathbf{V}^T$. راه‌حل حاصل برای \mathbf{w} اینگونه به نظر می‌رسد

$$\mathbf{w} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} = \mathbf{V}\Sigma^{-1}\mathbf{U}^T\mathbf{y} = \sum_{j=1}^d \frac{\mathbf{u}_j^T \mathbf{y}}{\sigma_j} \mathbf{v}_j \quad (5.2)$$

جایی که $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_n] \in \mathbb{R}^{n \times n}$ ماتریس متعامد متشکل از بردارهای یکتای سمت چپ است، $\mathbf{U}^d = [\mathbf{u}_1, \dots, \mathbf{u}_d] \in \mathbb{R}^{n \times d}$ اولین d بردار یکتای سمت چپ است و $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_d] \in \mathbb{R}^{d \times d}$ ماتریس متعامد متشکل از بردارهای یکتای راست است.

راه‌حل در معادله (5.3) روشن می‌کند که چرا راه‌حل رگرسیون خطی می‌تواند به آشفتگی‌ها حساس باشد. برای مقادیر کوچک یکتا، σ_j^{-1} بزرگ است و هر گونه تغییر در \mathbf{y} را تقویت می‌کند. به عنوان مثال، برای مولفه نویز کمی متفاوت $\epsilon_i \in$ برای نمونه i ، بردار راه‌حل \mathbf{w} می‌تواند بسیار متفاوت باشد. یک استراتژی رایج برای مقابله با این بی‌ثباتی، حذف یا کوتاه کردن مقادیر کوچک تکی است. این یک شکل منظم‌سازی است که در بخش بعدی به آن می‌پردازیم.

نکته: در حالت کلی، جایی که \mathbf{X} رتبه کامل نیست، هنوز هم می‌توانیم یک راه‌حل مینیمم مربعی برای $\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$ بدست آوریم. اکنون، راه‌حل‌های بی‌نهایت زیادی وجود دارد. انتخاب رایج انتخاب راه‌حل مینیمم واریانس است که مربوط به حذف مولفه‌ها (بردارهای یکتا) برای مقادیر صفر یکتا است:

$$\mathbf{w} = \sum_{j=1}^{\text{rank of } \mathbf{X}} \frac{\mathbf{u}_j^T \mathbf{y}}{\sigma_j} \mathbf{v}_j \quad (5.3)$$

مثال ۱۵: [تقریباً وابسته به خطی] بیایید به یک مثال ساده نگاه کنیم که چرا $\mathbf{X} \in \mathbb{R}^{n \times d}$ ممکن است مقادیر تکی کوچکی داشته باشد. اول، $d = 2$ و $x_2 = x_1$ را فرض کنید، یعنی ویژگی دوم یک کپی از اولی و به سادگی یک زائد است. سپس $\mathbf{X} = \mathbf{U}_2 \mathbf{\Sigma}_2 \mathbf{V}^T$ یک SVD باریک است، جایی که \mathbf{U}_2 فقط دو ستون اول SVD کامل را دارد. ما می‌توانیم این SVD باریک را بنویسیم زیرا $\mathbf{X} = \mathbf{U}_2 \mathbf{\Sigma}_2 \mathbf{V}^T \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ ، که در آن مقادیر یکتای صفر، ستون‌های باقی‌مانده \mathbf{U} را صفر می‌کند.

SVD فقط اولین ستون $\mathbf{x}_1 \in \mathbb{R}^{n \times 1}$ ساده است: $\mathbf{x}_1 = \mathbf{u}_1 \sigma_1 \mathbf{v}_1$ ، که در آن $\mathbf{u}_1 = \mathbf{x}_1 / \|\mathbf{x}_1\|$ ، $\sigma_1 = \|\mathbf{x}_1\|$ و $\mathbf{v}_1 =$ 1. بنابراین، SVD هست $\mathbf{X} = [\mathbf{x}_1 \mathbf{x}_2]$ برای هر n است. بردار واحد بعدی \mathbf{u}_2 که متعامد به \mathbf{u}_1 است و بردارهای یکتا راست $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^2$

$$\mathbf{X} = [\mathbf{u}_1 \mathbf{u}_2] \mathbf{\Sigma} [\mathbf{v}_1 \mathbf{v}_2]^T = [\mathbf{u}_1 \mathbf{u}_2] \begin{bmatrix} 2\sigma_1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0.5 & 0.5 \\ -0.5 & 0.5 \end{bmatrix} = u_1 \sigma_1 [1.0 \quad 1.0]$$

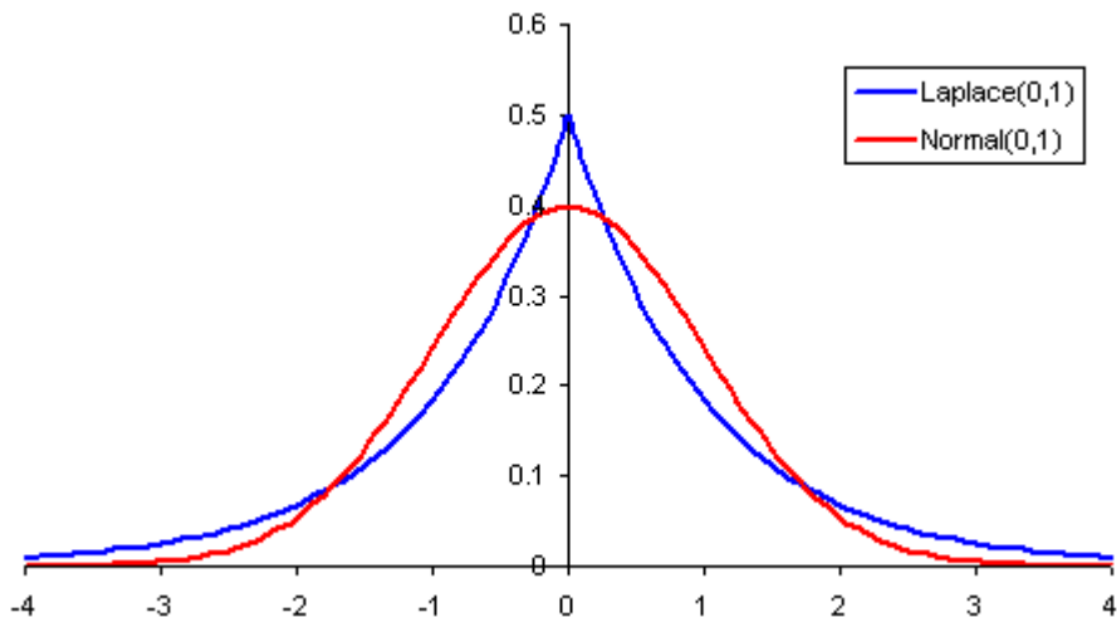


Figure ۵/۴: مقایسه بین پیشین‌های گاوسی و لاپلاس. هر دو ترجیح می‌دهند مقادیر نزدیک به صفر باشند، اما لاپلاس قبلی به شدت ترجیح می‌دهد مقادیر برابر با صفر باشد.

که در آن ما \mathbf{v}_1 را به دو بعدی گسترش دادیم (از آنجایی که $d = 2$)، و \mathbf{v}_2 را متعامد به آن بردار تعریف کردیم، و مجبور بودیم σ_1 را تغییر مقیاس دهیم تا بردارهای واحد یکتا را حفظ کنیم. بنابراین چون \mathbf{x}_2 به \mathbf{x}_1 وابسته است، وقتی آن را به عنوان یک ستون و مقدار یکتای $\sigma_2 = 0$ را اضافه می‌کنیم، رتبه افزایش نمی‌یابد.

اگر به جای $x_2 = x_1 + \epsilon$ برای یک بردار نویز کوچک $\epsilon \in \mathbb{R}^n$ باشد، در عوض می‌بینیم که σ_2 دیگر صفر نخواهد بود، بلکه بسیار نزدیک به صفر خواهد بود، زیرا u_1 و اولین مقدار یکتای σ_1 تا حد زیادی برای بازسازی x_2 قادر خواهند بود.

۲-۴-۵ منظم سازی^۱

تا اینجا در مورد رگرسیون خطی از نظر ماکسیمم احتمال بحث کردیم. اما، مانند قبل، ما همچنین می‌توانیم یک هدف MAP را پیشنهاد کنیم. به جای تعیین هیچ پیش از w ، می‌توانیم یکی پیش از آن را انتخاب کنیم تا به تنظیم بیش از حد برازش داده‌های مشاهده‌شده کمک کند. ما دو پیشین متداول (تنظیم کننده) را مورد بحث قرار خواهیم داد: پیشین گاوسی (نرم ℓ_2) و پیشین لاپلاس (نرم ℓ_1)، که در شکل ۵,۴ نشان داده شده است.

با گرفتن گزارش گاوسی میانگین صفر قبل، $\mathcal{N}(0, \lambda^{-1} I)$ ، به دست می‌آوریم

$$-\ln p(w) = \frac{1}{2} \ln(2\pi |\lambda^{-1} I|) + \frac{w^T w}{2\lambda^{-1}} = \frac{1}{2} \ln(2\pi) - d \ln(\lambda) + \frac{\lambda}{2} w^T w$$

زیرا $|I| = \lambda^{-d}$ ، که در آن $|A|$ تعیین کننده ماتریس A است. مانند قبل، می‌توانیم اولین ثابت را حذف کنیم که بر انتخاب w تأثیری ندارد.

اکنون می‌توانیم احتمال ورود منفی و گزارش منفی قبلی را با هم ترکیب کنیم. سپس با نادیده گرفتن ثابت‌ها، می‌توانیم \log -likelihood منفی و \log منفی را با قبل جمع کنیم تا به دست بیاوریم.

$$\begin{aligned} \operatorname{argmin}_w -\ln(p(y|X, w)) - \ln p(w) &= \operatorname{argmin}_w \frac{1}{2\sigma^2} \sum_{i=1}^n \left(y_i - \sum_{j=0}^d \omega_j x_{ij} \right)^2 + \frac{\lambda}{2} w^T w \\ &= \operatorname{argmin}_w \sum_{i=1}^n \left(y_i - \sum_{j=0}^d \omega_j x_{ij} \right)^2 + \frac{\lambda \sigma^2}{2} w^T w \end{aligned}$$

بنابراین، اگر فرض کنیم که وزن‌ها دارای میانگین گاوسی صفر $\mathcal{N}(0, \lambda^{-1} \sigma^2 I)$ قبلی هستند آنگاه مسئله رگرسیون پشته زیر را دریافت می‌کنیم:

$$c(w) = (Xw - y)^T (Xw - y) + \lambda w^T w \quad \triangleright \quad ||w||_2^2 = w^T w$$

که در آن λ یک پارامتر انتخاب شده توسط کاربر است که به آن پارامتر منظم‌سازی می‌گویند. ایده این است که ضرائب وزنی بیش از حد بزرگ را جریمه کنیم. هرچه λ بزرگتر باشد، وزنه‌های بزرگ بیشتری جریمه می‌شوند. به همین ترتیب، λ بزرگتر مربوط به یک کوواریانس کوچکتر در قبلی است، که وزن‌ها را نزدیک به صفر می‌کند. بنابراین، برآورد MAP باید بین این قبل از وزن‌ها و برازش داده‌های مشاهده‌شده تعادل برقرار کند.

اگر این معادله را به روش قبلی حل کنیم، به دست می‌آوریم

¹ Regularization

$$\mathbf{w}_{\text{MAP}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

این اثر خوب جابجایی مقادیر مجذور یکتای در Σ_d^2 در λ را دارد، تا زمانی که λ خود به اندازه کافی بزرگ باشد، مسائل پایداری را با تقسیم بر مقادیر کوچک حذف می‌کند.

اگر توزیع لاپلاس را انتخاب کنیم، هدف جریمه شده ℓ_1 دریافت می‌کنیم

$$c(\mathbf{w}) = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \|\mathbf{w}\|_1$$

که اغلب به آن Lasso می‌گویند. این هدف را می‌توان به طور مشابه با هدف منظم شده ℓ_2 به دست آورد، اما در عوض از توزیع لاپلاس با پارامتر λ برای قبلی استفاده کرد. همانند تنظیم کننده ℓ_2 برای رگرسیون برآمدگی، این تنظیم کننده مقادیر بزرگ در \mathbf{w} را جریمه می‌کند. با این حال، راه‌حل‌های پراکنده‌تری را نیز تولید می‌کند، جایی که ورودی‌های \mathbf{w} صفر هستند. این ترجیح را می‌توان در شکل ۵,۴ مشاهده کرد، جایی که توزیع لاپلاس بیشتر حول صفر متمرکز شده است. با این حال، در عمل، این ترجیح حتی قوی‌تر از آن چیزی است که توزیع نشان می‌دهد، به دلیل اینکه چگونه از دست دادن مینیمم مربعات کروی و تنظیم کننده «۱» با هم تعامل دارند.

اجباری کردن ورودی‌های \mathbf{w} به صفر تأثیر انتخاب ویژگی دارد، زیرا صفر کردن ورودی‌ها در \mathbf{w} معادل حذف ویژگی مربوطه است. هر بار که یک پیش‌بینی انجام می‌شود، محصول نقطه‌ای را در نظر بگیرید،

$$\mathbf{x}^T \mathbf{w} = \sum_{j=0}^d x_j w_j = \sum_{j: w_j \neq 0} x_j w_j$$

این معادل به سادگی حذف ورودی در \mathbf{X} و \mathbf{w} است که در آن $w_j = 0$ است.

برای Lasso، ما دیگر راه‌حلی به شکل بسته نداریم. ما یک راه‌حل شکل بسته نداریم، زیرا نمی‌توانیم \mathbf{w} را در فرم بسته‌ای که نقطه ثابتی را ارائه می‌دهد حل کنیم. در عوض، از گرادینان نزول برای محاسبه راه‌حل \mathbf{w} استفاده می‌کنیم. با این حال، تنظیم کننده ℓ_1 در 0 غیر قابل تمایز است. درک چگونگی بهینه‌سازی این هدف به پیشینه بهینه‌سازی کمی بیشتری نیاز دارد، بنابراین ما این الگوریتم را در فصل بعدی، در الگوریتم ۴ ارائه می‌کنیم.

۳-۴-۵ انتظار و واریانس برای راه‌حل منظم

یک سوال طبیعی که پیش می‌آید این است که چگونه می‌توان این پارامتر منظم‌سازی را انتخاب کرد و تأثیر آن بر بردار حل نهایی. انتخاب این پارامتر منظم‌سازی منجر به مبادله بایاس واریانس می‌شود. برای درک این مبادله، باید بدانیم که منظور از بایاس دار بودن راه‌حل چیست و چگونه می‌توان واریانس راه‌حل را در مجموعه داده‌های ممکن مشخص کرد.

اجازه دهید با درک بایاس و واریانس با یک راه‌حل غیرقانونی شروع کنیم، با این فرض که مفروضات توزیعی پشت رگرسیون خطی درست هستند. این بدان معنی است که یک پارامتر واقعی ω وجود دارد به طوری که برای هر یک از نقاط داده $Y_i = \sum_{j=0}^d \omega_j X_{ij} + \varepsilon_i$ ، جایی که ε_j ، i, d ، i است. متغیرهای تصادفی بر اساس $\mathcal{N}(0, \sigma^2)$ ترسیم شده‌اند. ما می‌توانیم بردار حل (تخمین‌گر) \mathbf{w}_{MLE} را به عنوان یک متغیر تصادفی مشخص کنیم، جایی که تصادفی بودن در میان مجموعه‌های داده ممکن است، که می‌توانستند مشاهده شوند. از این نظر، ما مجموعه داده \mathcal{D} را یک متغیر تصادفی در نظر می‌گیریم و راه‌حل $\mathbf{w}_{\text{ML}}(\mathcal{D})$ از آن مجموعه داده را به عنوان تابعی از این متغیر تصادفی در نظر می‌گیریم.

اجازه دهید اکنون به مقدار مورد انتظار (با توجه به مجموعه داده‌های آموزشی \mathcal{D}) برای بردار وزن w_{ML} با $\varepsilon = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)$ نگاه کنیم:

$$\begin{aligned}\mathbb{E}[w_{ML}(\mathcal{D})] &= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \omega + \varepsilon)] \\ &= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X})] + \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \varepsilon] \\ &= \mathbb{E}[\omega] + \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T] \mathbb{E}[\varepsilon] \\ &= \omega\end{aligned}$$

که در آن تساوی سوم از این واقعیت ناشی می‌شود که عبارتهای نویز ε مستقل از ویژگی‌ها و آخرین تساوی هستند، زیرا ω یک بردار ثابت (غیر تصادفی) است و $\mathbb{E}[\varepsilon] = 0$. برآوردگر که مقدار مورد انتظار آن مقدار واقعی است. پارامتر را برآوردگر بی طرف می‌نامند. ماتریس کوواریانس برای مجموعه بهینه پارامترها را می‌توان به این صورت بیان کرد

$$\begin{aligned}\text{Cov}[w_{ML}(\mathcal{D})] &= \mathbb{E}[(w_{ML}(\mathcal{D}) - \omega)(w_{ML}(\mathcal{D}) - \omega)^T] \\ &= \mathbb{E}[w_{ML}(\mathcal{D}) w_{ML}(\mathcal{D})^T] - \omega \omega^T\end{aligned}$$

ما $\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ را میگیریم و داریم $w_{ML}(\mathcal{D}) = \omega + \mathbf{X}^\dagger \varepsilon$ در نتیجه

$$\begin{aligned}\text{Cov}[w_{ML}(\mathcal{D})] &= \mathbb{E}[(\omega + \mathbf{X}^\dagger \varepsilon)(\omega + \mathbf{X}^\dagger \varepsilon)^T] - \omega \omega^T \\ &= \omega \omega^T + \mathbb{E}[\mathbf{X}^\dagger \varepsilon \varepsilon^T \mathbf{X}^{\dagger T}] - \omega \omega^T\end{aligned}$$

زیرا $\mathbb{E}[\omega + \mathbf{X}^\dagger \varepsilon] = \mathbb{E}[\mathbf{X}^\dagger] \mathbb{E}[\varepsilon] \omega^T = 0$. اکنون چون عبارات نویز مستقل از ورودی‌ها هستند، یعنی $\mathbb{E}[\varepsilon \varepsilon^T | \mathbf{X}] = \sigma^2 \mathbf{I}$ می‌توانیم از قانون احتمال کل (که قانون برج نیز نامیده می‌شود) استفاده کنیم.

$$\begin{aligned}\mathbb{E}[\mathbf{X}^\dagger \varepsilon \varepsilon^T \mathbf{X}^{\dagger T}] &= \mathbb{E}[\mathbb{E}[\mathbf{X}^\dagger \varepsilon \varepsilon^T \mathbf{X}^{\dagger T} | \mathbf{X}]] \\ &= \mathbb{E}[\mathbf{X}^\dagger \mathbb{E}[\varepsilon \varepsilon^T | \mathbf{X}] \mathbf{X}^{\dagger T}] \\ &= \sigma^2 \mathbb{E}[\mathbf{X}^\dagger \mathbf{X}^{\dagger T}]\end{aligned}$$

بنابراین، ما داریم

$$\text{Cov}[w_{ML}(\mathcal{D})] = \sigma^2 \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1}]$$

می‌توان نشان داد که برآوردگر $w_{ML}(\mathcal{D}) = \mathbf{X}^\dagger \mathbf{y}$ دارای کمترین واریانس در بین همه برآوردگرهای بی طرف است (قضیه گاوس-مارکوف).

با این حال، متأسفانه، همانطور که در بالا توضیح داده شد، ماتریس $\mathbf{X}^T \mathbf{X} = \mathbf{V} \mathbf{\Sigma} \mathbf{V}^T$ می‌تواند شرایط ضعیفی داشته باشد، با مقادیری صفر یا نزدیک به صفر. در نتیجه، این ماتریس کوواریانس می‌تواند شرط بدی شود، با مقادیر کوواریانس با بزرگی بالا. این نشان می‌دهد که در بین مجموعه‌های داده، راه‌حل $w_{ML}(\mathcal{D})$ می‌تواند بسیار متفاوت باشد. این نوع رفتار حکایت از تناسب بیش از حد دارد و مطلوب نیست. اگر راه‌حل ما می‌تواند در چندین زیر مجموعه تصادفی مختلف داده بسیار متفاوت باشد، نمی‌توانیم به هیچ یک از این راه‌حل‌ها اطمینان داشته باشیم.

از سوی دیگر، راه حل منظم شده، بسیار کمتر احتمال دارد که کوواریانس بالایی داشته باشد، اما دیگر بی طرف نخواهد بود. اجازه دهید $\mathbf{w}_{MAP}(\mathcal{D})$ تخمین MAP برای مسئله منظم ℓ_2 با مقدار $\lambda > 0$ باشد. با استفاده از تحلیلی مشابه با بالا، مقدار مورد انتظار $\mathbf{w}_{MAP}(\mathcal{D})$ برابر است با

$$\begin{aligned}\mathbb{E}[\mathbf{w}_{MAP}(\mathcal{D})] &= \mathbb{E}\left[(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T(\mathbf{X}\omega + \varepsilon)\right] \\ &= \mathbb{E}\left[(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}(\mathbf{X}^T\mathbf{X})\omega\right] \\ &\neq \omega\end{aligned}$$

همانطور که $\lambda \rightarrow 0$ ، راه حل MAP به بی طرف بودن نزدیک و نزدیکتر می شود. کوواریانس برابر است با

$$\text{Cov}[\mathbf{w}_{MAP}(\mathcal{D})] = \sigma^2 \mathbb{E}\left[(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}(\mathbf{X}^T\mathbf{X})(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\right]$$

این کوواریانس بسیار کمتر مستعد ابتلا به شرایط نامناسب $\mathbf{X}^T\mathbf{X}$ است، زیرا همانطور که در بالا بحث شد، تغییر λ شرایط را بهبود می بخشد. در نتیجه، ما انتظار داریم که \mathbf{w}_{MAP} واریانس کمتری در بین مجموعه داده های مختلف داشته باشد که می توان مشاهده کرد. این به طور متناظر نشان می دهد که ما کمتر به مجموعه داده ای اضافه می کنیم. توجه داشته باشید که با $\lambda \rightarrow \infty$ ، واریانس به صفر کاهش می یابد، اما بایاس به مقدار ماکسیمم آن افزایش می یابد (یعنی نرم وزن های واقعی). همانطور که در شکل ۵،۵ نشان داده شده است، یک انتخاب بهینه از λ وجود دارد که این مبادله بایاس واریانس را به مینیمم می رساند - اگر بتوانیم آن را پیدا کنیم.

دلیل اینکه ما به بایاس و واریانس اهمیت می دهیم این است که میانگین مجذور خطای مورد انتظار نسبت به وزن های واقعی را می توان به بایاس و واریانس تجزیه کرد. تا ببینیم چرا

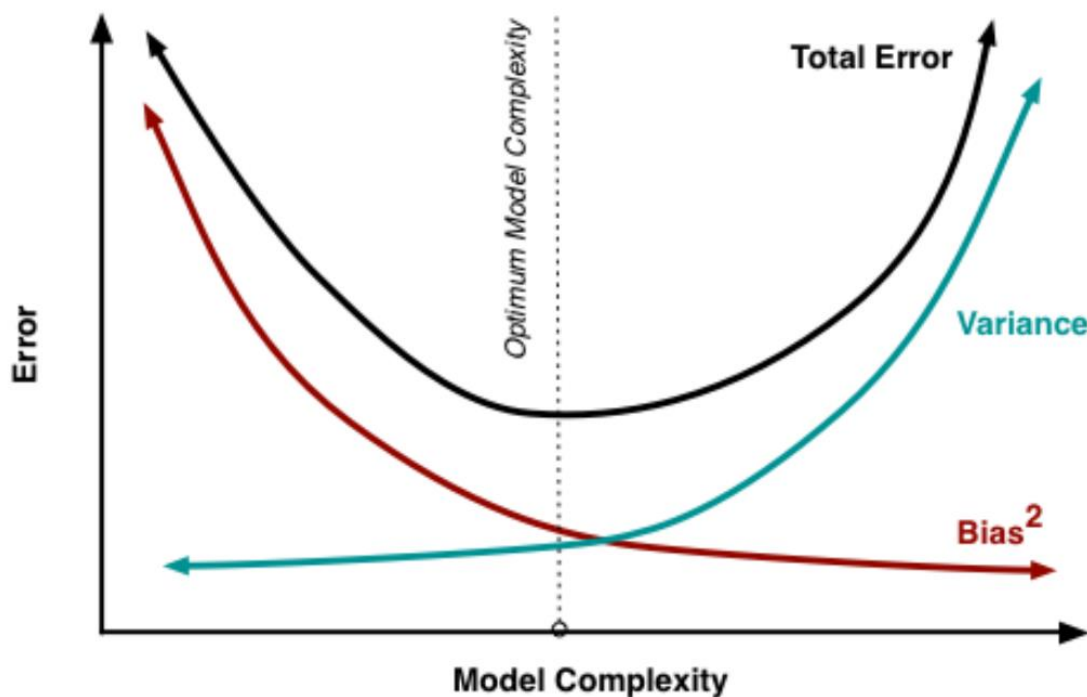
$$\begin{aligned}\mathbb{E}\left[\|\mathbf{w}(\mathcal{D}) - \omega\|_2^2\right] &= \mathbb{E}\left[\sum_{j=1}^d \{w_j(\mathcal{D}) - w_j\}^2\right] \\ &= \sum_{j=1}^d \mathbb{E}\left[(w_j(\mathcal{D}) - w_j)^2\right]\end{aligned}$$

که در آن ما می توانیم این اصطلاح درونی را بیشتر ساده کنیم

$$\begin{aligned}\mathbb{E}\left[(w_j(\mathcal{D}) - w_j)^2\right] &= \mathbb{E}\left[w_j(\mathcal{D}) - \mathbb{E}[w_j(\mathcal{D})] + \mathbb{E}[w_j(\mathcal{D})] - w_j\right]^2 \\ &= \mathbb{E}\left[(w_j(\mathcal{D}) - \mathbb{E}[w_j(\mathcal{D})])^2\right] + \mathbb{E}\left[(\mathbb{E}[w_j(\mathcal{D})] - w_j)^2\right]\end{aligned}$$

جایی که مرحله دوم از این واقعیت ناشی می شود که

$$\begin{aligned}-2\mathbb{E}\left[(w_j(\mathcal{D}) - \mathbb{E}[w_j(\mathcal{D})])(\mathbb{E}[w_j(\mathcal{D})] - w_j)\right] &= (\mathbb{E}[w_j(\mathcal{D})] - w_j)\mathbb{E}\left[w_j(\mathcal{D}) - \mathbb{E}[w_j(\mathcal{D})]\right] \\ &= 0\end{aligned}$$



شکل ۵/۵: مبادله بایاس واریانس. منبع: <http://scott.fortmann-roe.com/docs/BiasVariance.html>

جمله اول بالا در $\mathbb{E}(w_j(\mathcal{D}) - \omega_j)^2$ واریانس وزن j و جمله دوم بایاس وزن j است که در آن $\mathbb{E}(\mathbb{E}[w_j(\mathcal{D})] - \omega_j)^2 = 0$ زیرا هیچ چیز در این عبارت تصادفی نیست بنابراین انتظار بیرونی حذف می‌شود. که یعنی

$$\mathbb{E}[(f_{\mathcal{D}}(\mathbf{x}) - f(\mathbf{x}))^2] = (\mathbb{E}[f_{\mathcal{D}}(\mathbf{x})] - f(\mathbf{x}))^2 + V[f_{\mathcal{D}}(\mathbf{x})]$$

نشان می‌دهد که خطای میانگین مربعات مورد انتظار برای بردار وزن واقعی ω به بایاس مربع تجزیه می‌شود - که در آن بایاس $\mathbb{E}[w_j(\mathcal{D})] - \omega_j$ است - و واریانس. مبادله واریانس بایاس این واقعیت را منعکس می‌کند که تا زمانی که واریانس بیشتر از بایاس مجذور کاهش یابد، می‌توانیم به طور بالقوه خطای میانگین مربع را با اعمال برخی سوگیری کاهش دهیم.

نکته: ما به طور مستقیم مبادله سوگیری-واریانس را بهینه نمی‌کنیم. ما در واقع نمی‌توانیم بایاس را اندازه‌گیری کنیم، بنابراین مستقیماً این اصطلاحات را به مینیمم نمی‌رسانیم. در عوض، این تجزیه نحوه انتخاب مدل‌ها را راهنمایی می‌کند.

تمرین ۷: فرمول کوواریانس را برای $w_{MAP}(\mathcal{D})$ استخراج کنید

۵-۵ مبادله بایاس واریانس^۱

در بالا فرض کردیم که مدل واقعی خطی است، و بنابراین تنها سوگیری معرفی شده از قانون‌گذاری بود. این فرض بر این بود که فضای فرضی توابع خطی شامل تابع واقعی است، و این که سوگیری معرفی شده تنها به دلیل قاعده‌سازی است. در واقع، هنگام استفاده از رگرسیون خطی با منظم‌سازی، هم از انتخاب یک کلاس تابع ساده‌تر و هم از نظم‌دهی، بایاس را معرفی می‌کنیم. اگر تابع درست خطی نباشد، نمی‌توان وزن‌های آموخته‌شده را برای یک تابع خطی مستقیماً با تابع واقعی مقایسه کرد.

¹ Bias-Variance

اگر از یک مبنای قدرتمند برای تبدیل داده‌ها استفاده شود، آنگاه می‌توانیم توابع غیرخطی را یاد بگیریم حتی اگر راه حل از رگرسیون خطی استفاده کند. در این حالت، امکان پذیر است که این کلاس تابع به اندازه کافی قدرتمند باشد و تابع واقعی را شامل شود، و این سوگیری بیشتر به دلیل منظم شدن است. اما، به طور کلی، تضمین اینکه ما یک کلاس تابعی را که شامل تابع $true$ باشد، سخت خواهد بود، و مقایسه مستقیم پارامترهای ما با پارامترهای واقعی (که ممکن است حتی از یک بعد هم نباشند) دشوار خواهد بود.

با در نظر گرفتن خطای تقلیل پذیر، می‌توانیم به طور کلی‌تر در مورد سوگیری و واریانس صحبت کنیم. در واقع، مبادله بایاس واریانس تماماً در مورد کاهش خطای قابل کاهش است. (به یاد داشته باشید، ما نمی‌توانیم خطای تقلیل‌ناپذیر را کاهش دهیم - نام همه چیز را می‌گوید - با بهبود نحوه تخمین تابع.) ما می‌توانیم یک تجزیه بایاس واریانس کلی‌تر تعریف کنیم که خروجی‌های تابع را به جای بردارهای پارامتر مقایسه می‌کند. معادله خطای تقلیل پذیر را به $\mathbb{E}[(f_D(X) - f(X))^2]$ به یاد بیاورید، که در آن $f(X)$ تابع بهینه است، یعنی $f(x) = \mathbb{E}[Y|x]$ برای مربع هزینه. ما قبلاً در مورد این خطای تقلیل‌پذیر برای یک تابع ثابت بحث کرده‌ایم، با انتظار فقط بیش از X . اما اکنون به علاوه این واقعیت را در نظر می‌گیریم که f_D تصادفی است، و می‌توانیم در مورد انتظار و واریانس آن برای x معین استدلال کنیم.

بباید فقط با در نظر گرفتن خطای میانگین مربع مورد انتظار، برای یک ورودی داده شده x شروع کنیم. با استفاده از مراحل مشابه تجزیه بالا، دریافت می‌کنیم

$$\begin{aligned} & \mathbb{E}[(f_D(x) - f(x))^2] \\ &= (\mathbb{E}[f_D(x)] - f(x))^2 + V[f_D(x)] \end{aligned}$$

توجه کنید که در خط دوم، انتظار کنونی در مجذور فاصله است. این عبارت با بایاس مربعی مطابقت دارد. بایاس در اینجا خروجی تابع تخمینی $f_D(x)$ را در تمام مجموعه‌های داده D منعکس می‌کند. اصطلاح واریانس نشان می‌دهد که پیش‌بینی x چقدر می‌تواند متفاوت باشد، اگر در مجموعه داده‌های $i.i.d.$ مختلف یاد بگیریم. این تجزیه خطای میانگین مربع به بایاس و واریانس مجذور آشکار نیست، اما مراحل مشابه بالا را دنبال می‌کند. به عنوان تمرین باقی می‌ماند

تعمیم بالا نشان می‌دهد که یکی از راه‌هایی که ما بایاس و واریانس را متعادل می‌کنیم، در واقع انتخاب کلاس تابع است. اگر یک کلاس تابع ساده را انتخاب کنیم، کلاس احتمالاً به اندازه کافی بزرگ نیست - به اندازه کافی قدرتمند نیست - تا تابع واقعی را نشان دهد. این مقداری بایاس را معرفی می‌کند، اما احتمالاً واریانس کمتری نیز دارد، زیرا آن کلاس تابع ساده‌تر احتمال کمتری دارد که به هر مجموعه داده اضافه شود. اگر این کلاس خیلی ساده باشد، می‌توانیم بگوییم که تابع ما زیر پارامتر است و زیر برازش است. از طرف دیگر، اگر یک کلاس تابع قدرتمندتر را انتخاب کنیم که حاوی تابع واقعی باشد، ممکن است هیچ گونه سوگیری نداشته باشیم، اما به دلیل توانایی یافتن تابعی در کلاس بزرگ شما که بیش از حد با یک مجموعه داده معین مطابقت دارد، واریانس بالایی داشته باشیم. در این تنظیمات، ممکن است بگوییم که تابع بیش از حد پارامتر شده است، و اگرچه ما توانایی یادگیری یک تابع بسیار دقیق را داریم، یافتن آن تابع در این کلاس بزرگ‌تر مشکل خواهد بود. در عوض، احتمالاً مدلی انتخاب می‌شود که با داده‌های داده‌شده بیش از حد برازش می‌کند و به داده‌های جدید تعمیم نمی‌یابد (یعنی در داده‌های جدید ضعیف عمل می‌کند).

یافتن تعادل بین بایاس و واریانس، و بین عدم تناسب و برازش بیش از حد، یک مشکل اصلی در یادگیری ماشین است. ما در فصل ۱۰ راه‌هایی را برای بررسی تئوری و تجربی این مبادله مورد بحث قرار می‌دهیم.

فصل ۶

اصول بهینه‌سازی پیشرفته‌تر

با توجه به پیشینه بهینه‌سازی در فصل ۲، و مشاهده چگونگی مفید بودن آن در فصل‌های بعدی، اکنون می‌توانیم به رویکردهای بهینه‌سازی پیشرفته‌تر روی آوریم. اکنون به طور عمیق‌تر بحث خواهیم کرد که چگونه به روزرسانی مرتبه دوم گرادیان کاهشی را برای موارد چند متغیره بدست آوریم. سپس برخی از پیشرفت‌های محاسباتی در این روش‌ها را مورد بحث قرار می‌دهیم، به‌ویژه از طریق استفاده از تکنیک‌های انتخاب اندازه گام بهبودیافته، با استفاده از گرادیان کاهشی تصادفی و برخی تغییرات کوچک برای مقابله با نقاط غیر قابل تمایز. در نهایت، ما همچنین برخی از اصول اولیه را در مورد بهینه‌سازی محدود ارائه خواهیم کرد. هنگام حرکت به حالت چند متغیره، عادت کردن به حساب چند متغیره مفید خواهد بود. ما برخی از قوانین اساسی را در بخش B.1 ارائه می‌کنیم. مرجع کامل‌تری برای این قوانین را می‌توان در کتابچه راهنمای ماتریسی (بسیار مفید) [۱۶] یافت.

۱-۶ گرادیان کاهشی در توابع چند متغیره

می‌توانیم بحث در مورد به‌روزرسانی گرادیان کاهشی در بخش ۲،۲ را از حالت تک متغیره به حالت چند متغیره با استفاده از تقریب سری تیلور چند متغیره تعمیم دهیم. تقریب تیلور مرتبه دوم برای یک تابع حقیقی از چندین متغیر می‌تواند به این صورت نوشته شود.

$$c(\mathbf{w}) \approx \hat{c}(\mathbf{w}) = c(\mathbf{w}_0) + \nabla c(\mathbf{w}_0)^T (\mathbf{w} - \mathbf{w}_0) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_0)^T H_{c(\mathbf{w}_0)} (\mathbf{w} - \mathbf{w}_0)$$

که

$$\nabla c(\mathbf{w}_0) = \left(\frac{\partial c}{\partial w_1}(\mathbf{w}_0), \frac{\partial c}{\partial w_2}(\mathbf{w}_0), \dots, \frac{\partial c}{\partial w_d}(\mathbf{w}_0) \right) \in \mathbb{R}^d$$

گرادیان تابع c است که در \mathbf{w}_0 و ارزیابی می‌شود

$$\mathbf{H}_{c(w_0)} = \begin{bmatrix} \frac{\partial^2 c}{\partial w_1^2}(w_0) & \frac{\partial^2 c}{\partial w_1 \partial w_2}(w_0) & \dots & \frac{\partial^2 c}{\partial w_1 \partial w_d}(w_0) \\ \vdots & \frac{\partial^2 c}{\partial w_2^2}(w_0) & & \\ & \vdots & \ddots & \\ \frac{\partial^2 c}{\partial w_d \partial w_1}(w_0) & \dots & & \frac{\partial^2 c}{\partial w_d^2}(w_0) \end{bmatrix} \in \mathbb{R}^{d \times d}$$

ماتریس هسین^۱ تابع c است که در w_0 ارزیابی شده است. در بخش بعدی مقداری شهود برای هسین ارائه می‌کنیم، اما در اینجا می‌توان آن را به طور شهودی مشابه مشتق دوم در نظر گرفت. مانند مشتق دوم، اطلاعاتی در مورد انحنای تابع ارائه می‌دهد، و بنابراین اطلاعات مفیدی در مورد میزان گام برداشتن در جهت گرادیان برای هر w_i ارائه می‌دهد.

به عنوان یادآوری در مورد ضرب ماتریس-بردار، \mathbf{H} حاصل ضرب یک ماتریس $d \times d$ و $d \times 1$ بردار \mathbf{w} یک بردار $d \times 1$ است. سپس، گرفتن $\mathbf{w}^T \mathbf{H} \mathbf{w}$ حاصل ضرب نقطه‌ای بین \mathbf{w} بردار $1 \times d$ و $d \times 1$ بردار $\mathbf{H} \mathbf{w}$ است که منجر به یک اسکالر می‌شود. برای ضرب ماتریس بردار داریم

$$\mathbf{H} \mathbf{w} = \begin{bmatrix} H_{1,:} \\ H_{2,:} \\ \vdots \\ H_{d,:} \end{bmatrix} \mathbf{w} = \begin{bmatrix} H_{1:w} \\ H_{2:w} \\ \vdots \\ H_{d:w} \end{bmatrix} = \begin{bmatrix} \langle H_{1,:}, \mathbf{w} \rangle \\ \langle H_{2,:}, \mathbf{w} \rangle \\ \vdots \\ \langle H_{d,:}, \mathbf{w} \rangle \end{bmatrix}$$

هنگام انجام ضرب ماتریس-بردار، فقط می‌توانید تصور کنید که بردار \mathbf{w} به طرفین بچرخد و هر ردیف \mathbf{H} را ضرب کند. برای ضرب ماتریس-ماتریس، \mathbf{AB} ، باید اطمینان حاصل کنید که بعد دوم ماتریس \mathbf{A} برابر با بعد اول ماتریس \mathbf{B} است. -ضرب ماتریس‌ها-بردار برای هر ستون \mathbf{B} تجزیه می‌شود.

مانند قبل، برای دریافت به روزرسانی افزایشی، می‌توانیم گرادیان این تقریب را گرفته و نقطه ثابت (محلی) را بدست آوریم. با استفاده از قوانین اساسی خلاصه شده در زیر در بخش B.1، گرادیان $\hat{c}(\mathbf{w})$ برابر است با

$$\nabla \hat{c}(\mathbf{w}) = \nabla c(\mathbf{w}_0) + \mathbf{H}_{c(w_0)}(\mathbf{w} - \mathbf{w}_0)$$

باز هم می‌خواهیم \mathbf{w}_1 را به گونه‌ای پیدا کنیم که این گرادیان صفر باشد. اگر هنوز با معکوس یک ماتریس آشنا نیستید، در بخش‌های بعدی این یادداشت‌ها (به ویژه برای رگرسیون خطی در فصل ۵) بیشتر مورد بحث قرار خواهد گرفت. در حال حاضر، برای حل $\mathbf{H}_{c(w_0)}(\mathbf{w} - \mathbf{w}_0) = -\nabla c(\mathbf{w}_0)$ ، می‌توان معکوس $\mathbf{H}_{c(w_0)}^{-1}$ را محاسبه کرد و هر دو طرف معادله را در این معکوس ضرب کرد. این کار دوباره مشابه معکوس یک اسکالر است: $h^{-1} h = 1$. به روزرسانی چند متغیره مربوطه، که فراتر از معادله (۲،۱) برای حالت اسکالر گسترش یافته است،

$$\mathbf{w}_{i+1} = \mathbf{w}_i - \mathbf{H}_{c(w_i)}^{-1} \nabla c(\mathbf{w}_i) \quad (6.1)$$

در معادله ۶،۱، هم گرادیان و هم هسین در نقطه \mathbf{w}_i ارزیابی می‌شوند

اندازه هسین باعث می‌شود که انتخاب بین گرادیان کاهشی مرتبه اول و مرتبه دوم در حالت چند متغیره کمتر آشکار شود. بر خلاف مجموعه اسکالر، محاسبه خود هسین سنگین است (در اندازه \mathbf{w} درجه دوم) و محاسبه معکوس هسین حتی سنگین‌تر است. به همین دلیل، به‌روزرسانی‌های مرتبه اول سبک‌تر اغلب ترجیح داده می‌شوند. برای مثال، اگر محاسبه هسین مانند هدف

¹ Hessian

رگرسیون خطی هزینه $O(d^2n)$ داشته باشد، پیچیدگی محاسباتی گرادیان کاهشی مرتبه دوم $O(d^3 + d^2n)$ در هر تکرار است، با فرض زمان $O(d^3)$ برای یافتن معکوس‌های ماتریس. از سوی دیگر، دوباره برای رگرسیون خطی، پیچیدگی محاسباتی برای گرادیان کاهشی مرتبه اول فقط $O(dn)$ در هر تکرار است.

به‌روزرسانی مرتبه اول برای حالت چند متغیره تقریبی حتی بزرگ‌تر است، زیرا کل هسین با یک عدد اسکالر $\frac{1}{\eta}$ تقریب می‌یابد (که تقریب هسین را به یک ماتریس مورب با $\frac{1}{\eta}$ در قطر تبدیل می‌کند). سپس گرادیان تقریب مرتبه اول تبدیل می‌شود به

$$\nabla \hat{c}(w) = \nabla c(w_0) + \frac{1}{\eta}(w - w_0)$$

و در نتیجه به روزرسانی مرتبه اول

$$w_{i+1} = w_i - \eta_i \nabla c(w_i)$$

انتخاب این اندازه گام یک ملاحظه مهم است. ما قبلاً یک استراتژی اساسی را برای انتخاب اندازه گام مورد بحث قرار داده‌ایم. در بخش ۵،۶، چند مورد دیگر را مورد بحث قرار می‌دهیم.

۲-۶ خواص هسین

مانند مشتق دوم، هسین انحنای تابع را در نقطه w_0 منعکس می‌کند. هر ورودی نشان می‌دهد که چگونه مشتق جزئی برای w_j با تغییر w_i تغییر می‌کند.

برای شهود بیشتر، مشتق جهت را در نظر بگیرید. مشتق جهتی نشان می‌دهد که چگونه یک تابع (چند متغیری) با گام برداشتن مقدار کوچک t در یک جهت ثابت u تغییر می‌کند.

$$\lim_{t \rightarrow 0} \frac{c(w + tu) - c(w)}{t}$$

هنگامی که خود را محدود به تغییر تابع در این یک جهت می‌کنیم، تصور آن آسان‌تر است و به ما اجازه می‌دهد از قوانین مشتق دوم آشنا برای تنظیم تک متغیره استفاده کنیم.

$$w(t) = w + tu$$

$$g(t) = c(w(t))$$

ما می‌توانیم از قانون زنجیره‌ای در $g(t)$ برای محاسبه مشتق با توجه به t استفاده کنیم.

$$\dot{g}(t) = \nabla c(w(t))^T \frac{\partial(w(t))}{\partial t}$$

$$= \nabla c(w(t))^T u$$

$$\dot{g}(0) = \nabla c(w(t))^T u$$

$$= \nabla c(w)^T u = 0$$

جایی که آخرین تساوی رخ می‌دهد چون w یک نقطه ثابت است و بنابراین $\nabla c(w) = 0$ مشتق دوم است

$$\begin{aligned}\dot{\mathbf{g}}(t) &= \frac{\partial(\mathbf{w}(t))^T}{\partial t} \mathbf{H}_{c(\mathbf{w}(t))} \frac{\partial(\mathbf{w}(t))}{\partial t} \\ &= \mathbf{u}^T \mathbf{H}_{c(\mathbf{w}(t))} \mathbf{u} \\ \dot{\mathbf{g}}(0) &= \mathbf{u}^T \mathbf{H}_{c(\mathbf{w})} \mathbf{u}\end{aligned}$$

برای اینکه این نقطه ثابت \mathbf{w} (مطابق با $t = 0$) مینیمم محلی باشد، $\dot{\mathbf{g}}(0)$ باید آزمون مشتق دوم را برآورده کند: $\dot{\mathbf{g}}(0) > 0$. این آزمون تنها در صورتی انجام می‌شود که $\mathbf{H}_{c(\mathbf{w})}$ باشد. قطعی مثبت، با تعریف ماتریس قطعی مثبت. به یاد بیاورید که یک ماتریس قطعی مثبت \mathbf{H} همانی است که با توجه به $\mathbf{u} \neq 0$ ، $\mathbf{u}^T \mathbf{H} \mathbf{u} > 0$ یا معادل آن، همه مقادیر ویژه بزرگتر از صفر باشد. از آنجایی که \mathbf{u} یک جهت دلخواه دور از \mathbf{w} بود، هسین باید مثبت - معین باشد تا اطمینان حاصل شود که $\dot{\mathbf{g}}(0) > 0$ برای همه $\mathbf{u} \neq 0$.

بنابراین، مقادیر ویژه هسین، انحنای تابع را به صورت محلی منعکس می‌کند. اگر $\mathbf{H}_{c(\mathbf{w})}$ مقدار ویژه λ بسیار کوچکی داشته باشد، آنگاه بردار ویژه مربوطه \mathbf{u} که $\mathbf{H}_{c(\mathbf{w})} \mathbf{u} = \lambda \mathbf{u}$ را برآورده می‌کند - جهتی دورتر از \mathbf{w} جایی است که تابع تقریباً مسطح است. این به این دلیل است که $\dot{\mathbf{g}}(0) = \mathbf{u}^T \mathbf{H}_{c(\mathbf{w})} \mathbf{u} = \lambda \|\mathbf{u}\|_2^2 = \lambda$ بسیار کوچک است.

مثال ۱۶: اکنون می‌توانیم هسین $\mathbf{H}_{c(\mathbf{w})}$ را برای حل رگرسیون خطی در نظر بگیریم. این هسین ما را قادر می‌سازد تا بررسی کنیم که آیا واقعاً یک مینیمم محلی پیدا کرده‌ایم یا نه، اگر در عوض یک نقطه ثابت پیدا کرده‌ایم که ماکسیمم محلی یا یک نقطه زینتی است. هسین است

$$\mathbf{H}_{c(\mathbf{w})} = 2\mathbf{X}^T \mathbf{X}$$

این هسین ماتریس نیمه معین مثبت است. برای اینکه ببینید چرا، برای هر بردار $\mathbf{w} \neq 0$ در نظر بگیرید،

$$\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} = (\mathbf{X} \mathbf{w})^T \mathbf{X} \mathbf{w} = \|\mathbf{X} \mathbf{w}\|_2^2 \geq 0$$

که در آن فقط می‌تواند برابری اتفاق بیفتد - برای برخی از \mathbf{w} - اگر ستون‌های \mathbf{X} به صورت خطی وابسته باشند. از آنجایی که هسین برای هر \mathbf{w} نیمه معین است، این امر تحدب $c(\mathbf{w})$ را تأیید می‌کند. علاوه بر این، اگر ستون‌های \mathbf{X} به صورت خطی مستقل باشند، هسین مثبت قطعی است، که نشان می‌دهد مینیمم مطلق منحصر به فرد است.

۳-۶ مدیریت مجموعه داده‌های بزرگ

یکی از رویکردهای رایج برای مدیریت مجموعه داده‌های بزرگ استفاده از تقریب تصادفی است که در آن نمونه‌ها به صورت تدریجی پردازش می‌شوند. برای اینکه ببینیم چگونه این کار انجام می‌شود، اجازه دهید گرادیان تابع هدف، $\nabla c(\mathbf{w})$ را دوباره بررسی کنیم. ما یک محلول شکل بسته برای $\nabla c(\mathbf{w}) = 0$ به دست آوردیم. با این حال، برای بسیاری از توابع هدف دیگر، حل $\nabla c(\mathbf{w}) = 0$ به صورت بسته امکان پذیر نیست. در عوض، از مقدار اولیه \mathbf{w}_0 (معمولاً تصادفی) شروع می‌کنیم و سپس در جهت منفی گرادیان قدم می‌گذاریم تا به مینیمم محلی برسیم. این رویکرد شیب نزولی نامیده می‌شود و در الگوریتم ۲ خلاصه می‌شود. توجه کنید که در اینجا گرادیان با تعداد نمونه‌های n نرمال می‌شود، زیرا $\mathbf{X}^T (\mathbf{X} \mathbf{w} - \mathbf{y})$ با تعداد نمونه‌ها رشد می‌کند و انتخاب اندازه مراحل را دشوارتر می‌کند.

Algorithm 2: Batch Gradient Descent($c, \mathbf{X}, \mathbf{y}$)

```

1: // A non – optimized, basic implementation of batch gradient descent
2:  $\mathbf{w} \leftarrow$  random vector in  $\mathbb{R}^d$ 
3:  $err \leftarrow \infty$ 
4:  $tolerance \leftarrow 10e^{-4}$ 
5:  $max\ iterations \leftarrow 10e^5$ 
6: while  $|c(\mathbf{w}) - err| > tolerance$  and  $havenot\ reached\ max\ iterations$  do
7:  $err \leftarrow c(\mathbf{w})$   $\triangleright$  for linear regression,  $c(\mathbf{w}) = \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$ 
8:  $\mathbf{g} \leftarrow \nabla c(\mathbf{w})$   $\triangleright$  for linear regression,  $\nabla c(\mathbf{w}) = \frac{1}{n} \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y})$ 
9: // The step – size  $\eta$  could be chosen by line – search, as in Algorithm 1
10:  $\eta \leftarrow linesearch(\mathbf{w}, c, \mathbf{g})$ 
11:  $\mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{g}$ 
12: return  $\mathbf{w}$ 

```

با این حال، برای تعداد زیادی نمونه n ، محاسبه گرادیان در تمام نمونه‌ها می‌تواند سنگین یا غیرممکن باشد. یک جایگزین این است که گرادیان را با دقت کمتری با نمونه‌های کمتر تقریبی کنیم. در تقریب تصادفی، ما معمولاً شیب را با یک نمونه تقریب می‌زنیم، مانند الگوریتم ۳. اگرچه این رویکرد ممکن است بیش از حد تقریبی به نظر برسد، یک تاریخچه نظری و تجربی طولانی وجود دارد که اثربخشی آن را نشان می‌دهد (برای مثال [۵، ۶] را ببینید). با افزایش روزافزون اندازه مجموعه داده‌ها برای بسیاری از سناریوها، کلیت تقریب تصادفی آن را می‌توان به روشی مدرن برای برخورد با داده‌های بزرگ تبدیل کرد. برای سناریوهای تخصصی، البته رویکردهای دیگری نیز وجود دارد. برای مثال، [۱۷] را ببینید.

الگوریتم آموزشی برای گرادیان کاهشی تصادفی اکنون می‌تواند مورد بازنگری قرار گیرد تا به طور تصادفی یک نقطه داده را در هر زمان از \mathcal{D} ترسیم کند و سپس وزن‌های فعلی را با استفاده از معادله قبلی به روز کند. به طور معمول، در عمل، این مستلزم تکرار یک یا چند بار در مجموعه داده به ترتیب (با فرض تصادفی بودن، با نمونه‌های $i, i, d.$) است. هر تکرار بر روی مجموعه داده یک دوره نامیده می‌شود. شرایط همگرایی معمولاً شامل شرایط اندازه گام‌ها می‌شود که نیاز به کاهش آنها در طول زمان دارد. مانند نزول شیب دسته ای، این به روزرسانی‌های شیب نزولی تصادفی همگرا می‌شوند، هرچند با نوسان بیشتر حول بردار وزن حقیقی، با کاهش اندازه گام به تدریج این نوسانات را هموار می‌کند.

Algorithm 3: Stochastic Gradient Descent($c, \mathbf{X}, \mathbf{y}$)

```

1:  $\mathbf{w} \leftarrow$  random vector in  $\mathbb{R}^d$ 
2: for  $i = 1, \dots, \text{number of epochs}$  do
3:   Shuffle data points from  $1, \dots, n$ 
4:   for  $j = 1, \dots, n$  do
5:      $\mathbf{g} \leftarrow \nabla c_j(\mathbf{w})$   $\triangleright$  for linear regression,  $\nabla c_j(\mathbf{w}) = (x_j^T \mathbf{w} - y_j) x_j$ 
6:     // For convergence, the step – size  $\eta_t$  needs to decrease with time, such as
7:     //  $\eta_t = \eta_0 t^{-1/2}$  or  $\eta_t = \eta_0 i^{-1}$  for an initial  $\eta_0$  (e.g.,  $\eta_0 = 1.0$ ).
8:     // In practice, it is common to pick a fixed, small stepsize
9:      $\eta_t \leftarrow i^{-1}$ 
10:     $\mathbf{w} \leftarrow \mathbf{w} - \eta_t \mathbf{g}$ 
11: return  $\mathbf{w}$ 

```

۶-۴ بهینه‌سازی غیر هموار اما همچنان مستمر

ما در سراسر این یادداشت‌ها فرض می‌کنیم که اهداف ما مستمر هستند. با این حال، این به معنای هموار بودن آنها نیست: در برخی موارد، این اهداف پیوسته ممکن است دارای نقاط غیر قابل تمایز باشند. به عنوان مثال، تنظیم کننده ℓ_1 در 0 غیر قابل تمایز است، و $\|Xw - y\|_2^2 + \lambda \|w\|_1$ را غیر قابل تمایز می‌کند. یک استراتژی استفاده از زیر گرادیان کاهشی است. این به معنای انتخاب یک انتخاب معقول برای گرادیان در نقطه غیر قابل تمایز است. برای مثال، در اینجا، می‌توانیم مشتق جزئی ℓ_1 برای w_j را صفر در صفر، -1 برای $w_j > 0$ و 1 برای $w_j < 0$ در نظر بگیریم. متأسفانه، این کاهشی است زیرا تمایل به پرش در اطراف صفر وجود دارد. برخلاف ℓ_2 ، گرادیان به تدریج نزدیک به صفر کاهش نمی‌یابد و به آرامی w_j را کاهش می‌دهد، بلکه بین دو مقدار بزرگ -1 و 1 می‌پرد. با چنین گرادیان بزرگی، کاهش تدریجی w_j به صفر دشوار است، حتی اگر این راه‌حل بهینه باشد.

یک جایگزین برای چنین اهداف غیر هموار، استفاده از روش‌های پروگزیمال است. ایده ساده است: از شیب نزول برای مولفه هموار بهینه‌سازی (اصطلاح خطای $\|Xw - y\|_2^2$) استفاده کنید، و سپس برای مقادیر w که نزدیک به صفر هستند، آنها را روی صفر قرار دهید. این ایده آستانه‌سازی، اگرچه ساده است، اما از لحاظ نظری یک رویکرد صحیح برای بهینه‌سازی با «۱» غیر هموار است. این عملگر آستانه‌ای، عملگر پروگزیمال^۱ نامیده می‌شود و می‌توان آن را به عنوان یک عملگر پروجکشن^۲ دید. هر بار که w با گرادیان به روز می‌شود، آن را از یک راه‌حل پراکنده دور می‌کند. سپس عملگر پروگزیمال w را به فضای راه‌حل‌های پراکنده باز می‌گرداند. عملگر پروگزیمال برای ℓ_1 از نظر عنصر به w اعمال می‌شود، و بنابراین در هر w_i به عنوان، با اندازه گام η و پارامتر تنظیم λ ، تعریف می‌شود.

$$\text{prox}_{\eta\lambda\ell_1}(w_i) = \begin{cases} w_i - \eta\lambda & \text{if } w_i > \eta\lambda \\ 0 & \text{if } |w_i| \leq \eta\lambda \\ w_i + \eta\lambda & \text{if } w_i < -\eta\lambda \end{cases}$$

عملگر پروگزیمال در کل بردار w از نظر عنصر اینگونه تعریف می‌شود: $\text{prox}_{\eta\lambda\ell_1}(w) = [\text{prox}_{\eta\lambda\ell_1}(w_1), \dots, \text{prox}_{\eta\lambda\ell_1}(w_d)]$. به خوبی، این تئوری بیان می‌کند که اندازه گام نباید بزرگ‌تر از معکوس ثابت لپشیتز^۳ برای بخش هموار هدف باشد، جایی که به طور شهودی ثابت‌های لپشیتز نشان‌دهنده سرعت تغییر تابع است. در الگوریتم ۴، ما یک الگوریتم گرادیان کاهشی برای به روزرسانی افزایشی با منظم کننده ℓ_1 ارائه می‌دهیم که به عنوان الگوریتمی به نام $ISTA$ [۴] معرفی شده است. به طور کلی‌تر، روش‌های پروگزیمال برای اهداف غیر هموار دیگر استفاده می‌شود، اگرچه در این یادداشت‌ها ما فقط Lasso را در نظر می‌گیریم.

Algorithm 4: Batch gradient descent for ℓ_1 regularized linear regression (X, y, λ)

- 1: $w \leftarrow 0 \in \mathbb{R}^d$
- 2: $err \leftarrow \infty$
- 3: $tolerance \leftarrow 10e^{-4}$
- 4: // Precomputing these matrices, to avoid recomputing them in the loop
- 5: $XX \leftarrow \frac{1}{n} X^T X$

¹ proximal

² projection

³ Lipschitz

```

6:  $\mathbf{X}y \leftarrow \frac{1}{n} \mathbf{X}^T \mathbf{y}$ 
7: // This stepsize is specific to the least – squares loss for linear regression
8:  $\eta \leftarrow 1 / (2 \|\mathbf{X}\mathbf{X}\|_F)$ 
9: while  $|c(\mathbf{w}) - err| > tolerance$  and have not reached max iterations do
10:    $err \leftarrow c(\mathbf{w})$ 
11:   // Proximal operator projects back into the space of sparse solutions given by  $\ell_1$ 
12:    $\mathbf{w} \leftarrow prox_{\eta\lambda\ell_1}(\mathbf{w} - \eta\mathbf{X}\mathbf{X}\mathbf{w} + \eta\mathbf{X}\mathbf{y})$ 
13: return  $\mathbf{w}$ 

```

۵-۶ روش‌های بیشتر برای انتخاب اندازه گام‌ها

از آنجایی که انتخاب اندازه گام بخش مهمی از یک الگوریتم فرود موثر است، راه‌های زیادی برای انجام این کار وجود دارد. علاوه بر جستجوی خط، یکی از رایج‌ترین روش‌ها استفاده از روش‌های شبه مرتبه دوم (یا شبه نیوتن) است. همانطور که دیدیم، معکوس هسین راه خوبی برای انتخاب اندازه گام ارائه می‌دهد، اما معمولاً برای محاسبه بسیار سنگین است چه برسد به معکوس کردن. روش‌های شبه مرتبه دوم تقریباً هسین را با کمترین میزان ذخیره‌سازی و محاسبات ممکن انجام می‌دهند. یکی از ساده‌ترین این تقریب‌ها تقریب فقط قطر هسین و معکوس کردن آن است که فقط برای محاسبه $O(d)$ و فضا هزینه دارد. چنین تقریبی معمولاً حتی برای مورب هسین معکوس بسیار ضعیف است و بنابراین معمولاً استفاده نمی‌شود. در عوض، محبوب‌ترین روش‌ها عبارتند از [14] LBFGS ، [21] Adadelta و [11] Adam .

فصل ۷

مدل‌های خطی تعمیم یافته^۱

در بخش‌های قبلی، دیدیم که چارچوب آماری بینش‌های ارزشمندی را در مورد رگرسیون خطی ارائه می‌کند، به‌ویژه با توجه به بیان صریح بیشتر مفروضات در سیستم (تصویر کامل را تنها زمانی خواهیم دید که از فرمول بیزی استفاده شود). این مفروضات برای برآورد دقیق پارامترهای مدل ضروری بودند، که سپس می‌توان از آن برای پیش‌بینی نقاط داده‌ای که قبلاً دیده نشده بود استفاده کرد.

در این بخش، مدل‌های خطی تعمیم‌یافته (GLM) را معرفی می‌کنیم که رگرسیون مینیمم مربعات معمولی را فراتر از توزیع‌های احتمال گاوسی و وابستگی‌های خطی بین ویژگی‌ها و هدف گسترش می‌دهند. این تعمیم همچنین شما را با طیف وسیع‌تری از توابع از دست دادن، به نام واگرایی برگمن^۲ آشنا می‌کند.

ابتدا نکات اصلی رگرسیون مینیمم مربعات معمولی را مرور خواهیم کرد. در آنجا، ما فرض کردیم که مجموعه‌ای از $i.i.d.$ نقاط داده با اهداف خود $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ بر اساس توزیع $p(x, y)$ ترسیم شدند. ما همچنین فرض کردیم که یک رابطه اساسی بین ویژگی‌ها و هدف خطی است، یعنی

$$Y = \sum_{j=0}^d \omega_j X_j + \varepsilon$$

که در آن ω مجموعه‌ای از وزن‌های مجهول بود و ε یک متغیر تصادفی با میانگین صفر با واریانس σ^2 بود. به منظور ساده‌سازی تعمیم، این مدل را کمی دوباره فرموله می‌کنیم. به‌ویژه، جدا کردن رابطه خطی اساسی بین ویژگی‌ها و هدف از این واقعیت که Y به طور معمول توزیع شده است مفید خواهد بود. یعنی خواهیم نوشت

$$1. \mathbb{E}[y|x] = \omega^T x$$

$$2. p(y|x) = \mathcal{N}(\mu, \sigma^2)$$

با $\mu = \omega^T x$ دو عبارت را به هم متصل می‌کند. این روش فرمول‌بندی رگرسیون خطی به ما امکان می‌دهد (1) چارچوب را به روابط غیرخطی بین ویژگی‌ها و هدف تعمیم دهیم و همچنین (2) از توزیع‌های خطا به غیر از گاوسی استفاده کنیم.

¹ Generalized Linear Models

² Bregman divergences

۷-۱ انتقال نمایی و توزیع پواسون

ابتدا با مثالی از GLM شروع می‌کنیم، قبل از اینکه به کلاس عمومی و تعریف کلی برویم. فرض کنید که نقاط داده مطابق با شهرهای جهان است که با برخی از ویژگی‌های عددی توصیف شده‌اند - و متغیر هدف تعداد روزهای آفتابی مشاهده شده در یک سال خاص است. متغیر هدف y با توجه به ویژگی‌های x ممکن است شبیه توزیع پواسون باشد. بنابراین طبیعی‌تر است که $p(y|x) = \text{Poisson}(\lambda)$ را مدل کنیم، که در آن $\lambda > 0$ پارامتر (میانگین) توزیع پواسون است: $E[y|x] = \lambda$. با این حال، چون $\lambda \in \mathbb{R}^+$ ، مدل کردن λ با $\omega^T x \in \mathbb{R}$ مناسب نخواهد بود. بلکه می‌خواهیم پیش‌بینی خطی خود را با تابع f انتقال دهیم تا محدوده ترکیب خطی ویژگی‌ها را به دامنه پارامترهای توزیع احتمال تنظیم کنیم.

ما می‌توانیم این کار را با معرفی یک انتقال نمایی برای این توزیع پواسون و به طور کلی‌تر، بعداً هر تابع انتقال معکوس f انجام دهیم. اگر بجای آن بتوانیم ω را به گونه‌ای تخمین بزنیم که $\lambda = e^{\omega^T x}$ ، آنگاه می‌توانیم تضمین کنیم که تخمین‌هایمان در محدوده صحیح هستند. از طرف دیگر، می‌توان در نظر گرفت که ما در حال یادگیری یک وزن دهی خطی از ویژگی‌ها برای یادگیری یک پارامتر تبدیل شده، $\log(\lambda) = \omega^T x$ هستیم. این اصلاح ساده به همین دلیل است که این مدل‌ها را مدل‌های خطی کلی‌شده می‌نامند، زیرا مؤلفه اصلی وزن‌دهی خطی است. ما انواع توزیع‌ها و انتقال‌هایی را که می‌توان در بخش‌های زیر در نظر گرفت، فرمول‌بندی می‌کنیم، اما ابتدا این مثال را با رگرسیون پواسون برای ارائه یک مثال ملموس به پایان می‌رسانیم. برای ایجاد مدل GLM برای رگرسیون پواسون، (۱) یک انتقال نمایی بین انتظار هدف و ترکیب خطی ویژگی‌ها، و (۲) توزیع پواسون برای متغیر هدف را فرض می‌کنیم.

$$1. E[y|x] = \exp(\omega^T x) \text{ or } \log(E[y|x]) = \omega^T x$$

$$2. p(y|x) = \text{Poisson}(\lambda)$$

با استفاده از این واقعیت که $E[y|x] = \lambda$ ، دو فرمول را با استفاده از $\lambda = e^{\omega^T x}$ به هم وصل می‌کنیم. توزیع احتمال اینگونه حاصل میشود

$$p(y|x) = \frac{e^{\omega^T x y} \cdot e^{-e^{\omega^T x}}}{y!} \quad \text{for any } y \in \mathbb{N}$$

برای یافتن پارامترهای مدل رگرسیون می‌توانیم از تخمین ماکسیمم درست‌نمایی استفاده کنیم. تابع $\log - \text{likelihood}$ که به شکل

$$ll(\mathbf{w}) = \sum_{i=1}^n ll_i(\mathbf{w})$$

$$ll_i(\mathbf{w}) = \mathbf{w}^T \mathbf{x}_i y_i - e^{\mathbf{w}^T \mathbf{x}_i} - \ln y_i!$$

هدف ما به مقدار مینیمم رساندن احتمال ورود به سیستم منفی است: $\min_{\mathbf{w}} -ll(\mathbf{w})$. به راحتی می‌توان فهمید که $\nabla ll(\mathbf{w}) = 0$ راه‌حل یک شکل بسته ندارد. بنابراین، برخلاف رگرسیون خطی، باید از گرادینت کاهشی استفاده کنیم. می‌توانیم از شیب نزول مرتبه اول یا دوم و گرادینت کاهشی دسته‌ای یا تصادفی استفاده کنیم. مرحله کلیدی در هر یک از اینها این است که ابتدا گرادینت را برای یک نمونه محاسبه کنید. ما با استخراج مشتق جزئی لگاریتم احتمال منفی برای یک نمونه شروع می‌کنیم

$$\begin{aligned}
-\frac{\partial l_i(\mathbf{w})}{\partial w_j} &= e^{\mathbf{w}^T \mathbf{x}_i} x_{ij} - x_{ij} y_i \\
&= x_{ij} (e^{\mathbf{w}^T \mathbf{x}_i} - y_i)
\end{aligned}$$

گرادیان برای یک نمونه اینگونه است

$$-\nabla l_i(\mathbf{w}) = \mathbf{x}_i \cdot (p_i - y_i)$$

که در آن $p_i = e^{\mathbf{w}^T \mathbf{x}_i}$ پیش‌بینی است. توجه داشته باشید که $p_i - y_i$ مربوط به یک خطای پیش‌بینی برای نمونه i است. گرادیان دسته‌ای است

$$\begin{aligned}
-\nabla l(\mathbf{w}) &= -\sum_{i=1}^n \nabla l_i(\mathbf{w}) \\
&= \sum_{i=1}^n \mathbf{x}_i (p_i - y_i) \\
&= \mathbf{X}^T (\mathbf{p} - \mathbf{y})
\end{aligned}$$

که در آن \mathbf{p} یک بردار با عناصر $p_i = e^{\mathbf{w}^T \mathbf{x}_i}$ است، که در آن $\mathbf{p} - \mathbf{y}$ یک بردار خطا است.

معمولاً، اکنون فقط به صورت تصادفی یا شیب نزولی دسته‌ای انجام می‌شود. برای گرادیان کاهش تصادفی، هر مرحله شامل استفاده از گرادیان برای یک نمونه (یعنی $l_i(\mathbf{w})$) و برای شیب نزولی دسته‌ای، هر مرحله شامل استفاده از گرادیان برای همه نمونه‌ها (یعنی $l(\mathbf{w})$) است. علاوه بر این، می‌توانیم ماتریس هسین را در نظر بگیریم، هم برای ارزیابی ویژگی‌های نقاط ثابت و هم برای گرادیان کاهش مرتبه دوم، اگرچه، اگر d بزرگ باشد، احتمالاً خیلی سنگین است. دومین مشتق جزئی تابع احتمال ورود به سیستم منفی برای یک نمونه است

$$\begin{aligned}
-\frac{\partial^2 l_i(\mathbf{w})}{\partial w_j \partial w_k} &= x_{ij} e^{\mathbf{w}^T \mathbf{x}_i} x_{ik} \\
&= x_{ij} p_i x_{ik}
\end{aligned}$$

با

$$-\frac{\partial^2 l_i(\mathbf{w})}{\partial w_j \partial w_k} = \sum_{i=1}^n \frac{\partial^2 l_i(\mathbf{w})}{\partial w_j \partial w_k}$$

برای \mathbf{P} یک ماتریس مورب $n \times n$ با p_i در قسمت مورب، ماتریس هسین اینگونه است

$$\mathbf{H}_{-l}(\mathbf{w}) = \mathbf{X}^T \mathbf{P} \mathbf{X}$$

اگر \mathbf{X} رتبه پایین نباشد، این ماتریس مثبت است، که به این معنی است که تنها یک نقطه ثابت وجود دارد و آن مینیمم مطلق است. در واقع، می‌دانیم که هدف رگرسیون پواسون محدب است، حتی اگر \mathbf{X} رتبه کامل نباشد، و بنابراین همه نقاط ثابت مینیمم‌های مطلق هستند.

تمرین ۸: آیا اگر \mathbf{X} رتبه کامل نباشد یک مینیمم مطلق وجود دارد؟

تمرین ۹: به روز رسانی مرتبه دوم برای رگرسیون پواسون چیست؟

۷-۲ توزیع‌های خانوادگی نمایی

در بخش قبل، از یک مثال خاص برای نشان دادن چگونگی تعمیم فراتر از توزیع‌های گاوسی استفاده کردیم. این رویکرد به طور کلی به هر خانواده توزیع نمایی گسترش می‌یابد. برای سادگی، در اینجا ما روی خانواده نمایی طبیعی تمرکز می‌کنیم که برای اکثر مدل‌های خطی تعمیم یافته کافی است. خانواده نمایی طبیعی کلاسی از توزیع‌های احتمال با شکل زیر است

$$p(x|\theta) = \exp(\theta x - a(\theta) + b(x))$$

در جایی که $\theta \in \mathbb{R}$ پارامتر توزیع است، $a: \mathbb{R} \rightarrow \mathbb{R}$ یک تابع نرمال ساز \log است و $b: \mathbb{R} \rightarrow \mathbb{R}$ تابعی از x است که معمولاً در بهینه سازی ما نادیده گرفته می‌شود زیرا تابع θ نیست. بسیاری از توزیع‌ها (خانواده‌های) که اغلب با آن‌ها مواجه می‌شوند، اعضای خانواده نمایی هستند. به عنوان مثال، توزیع‌های نمایی، گاوسی، گاما، پواسون یا دوجمله‌ای. بنابراین، مطالعه کلی خانواده نمایی برای درک بهتر مشترکات و تفاوت‌های بین توابع تک تک اعضا مفید است.

مثال ۱۷: توزیع پواسون را می‌توان به صورت بیان کرد

$$p(x|\lambda) = \exp(x \log \lambda - \lambda - \log x!)$$

که $\mathcal{X} = \mathbb{N}_0$ و $\lambda \in \mathbb{R}^+$

اکنون اجازه دهید بینش بیشتری در مورد ویژگی‌های پارامترهای خانواده نمایی و اینکه چرا این کلاس برای تخمین راحت است، به دست آوریم. تابع $a(\theta)$ معمولاً تابع \log - partitioning یا به سادگی یک \log - normalizer نامیده می‌شود. به این دلیل نامیده می‌شود

$$a(\theta) = \log \int_{\mathcal{X}} \exp(\theta x + b(x)) dx$$

و بنابراین نقش اطمینان از داشتن چگالی معتبر را ایفا می‌کند: $\int_{\mathcal{X}} p(x) dx = 1$ نکته مهم این است که برای بسیاری از GLM ‌های رایج، مشتق a با تابع انتقال مطابقت دارد. برای مثال، برای رگرسیون پواسون، تابع انتقال $f(\theta) = \exp(\theta)$ ، و مشتق e^θ برابر a است. بنابراین، همانطور که در پایین بحث می‌کنیم، \log - normalizer برای یک خانواده نمایی نشان می‌دهد که چه انتقال f باید استفاده شود.

ویژگی‌های این \log - normalizer برای تخمین مدل‌های خطی تعمیم یافته نیز کلیدی هستند. می‌توان آن را استخراج کرد که

$$\frac{\partial a(\theta)}{\partial \theta} = \mathbb{E}[X]$$

$$\frac{\partial^2 a(\theta)}{\partial \theta^2} = \text{V}[X]$$

۷-۳ فرمول بندی کردن مدل های خطی تعمیم یافته

اکنون مدل های خطی تعمیم یافته را فرمول بندی می کنیم. دو جزء کلیدی GLM ها را می توان به صورت بیان کرد

$$1. E[y|x] = f(\omega^T x) \text{ or } g(E[y|x]) = \omega^T x \text{ where } g = f^{-1}$$

$$2. p(y|x) \quad \text{توزیع خانواده نمایی است}$$

تابع f را تابع انتقال و g را تابع پیوند می نامند. برای رگرسیون پواسون، f تابع نمایی است، و همانطور که برای رگرسیون لجستیک^۱ خواهیم دید، f تابع سیگموئید^۲ است. تابع انتقال محدوده $\omega^T x$ را به دامنه Y تنظیم می کند. به دلیل این رابطه، توابع پیوند معمولاً مستقل از توزیع Y انتخاب نمی شوند. تعمیم به خانواده نمایی از توزیع گاوسی که در رگرسیون مینیمم مربعات معمولی استفاده می شود، به ما امکان مدل سازی طیف وسیع تری از توابع هدف را می دهد. GLM ها شامل سه مدل پرکاربرد می باشند: رگرسیون خطی، رگرسیون پواسون و رگرسیون لجستیک که در فصل بعدی در مورد آنها صحبت خواهیم کرد.

برای ارتباط واضح تر اینها با توزیع های خانواده نمایی، باید توزیع های شرطی را در نظر بگیریم. هر $p(y|x)$ یک توزیع خانواده نمایی با پارامتر $\theta = x^T w$ است. هنگام یادگیری w - با به ماکسیمم رساندن احتمال - ما در حال یادگیری پارامتر θ_i برای هر نمونه هستیم (x_i, y_i) . $\log - \text{likelihood}$ منفی کلی است

$$\begin{aligned} -ll(w) &= -\log \prod_{i=1}^n e^{\theta_i y_i - a(\theta) + b(y_i)} \\ &= - \sum_i (\theta_i y_i - a(\theta) + b(y_i)) \\ &= - \sum_i ll_i(w) \end{aligned}$$

با گرادیان های

$$\begin{aligned} - \frac{\partial ll_i(w)}{\partial w_j} &= \frac{\partial a(\theta_i)}{\partial w_j} - \frac{\partial \theta_i}{\partial w_j y_i} \\ &= \frac{\partial a(\theta_i)}{\partial \theta_i} \frac{\partial \theta_i}{\partial w_j} - \frac{\partial \theta_i}{\partial w_j} y_i \\ &= \left(\frac{\partial a(\theta_i)}{\partial \theta_i} - y_i \right) \frac{\partial \theta_i}{\partial w_j} \end{aligned}$$

همانطور که برای رگرسیون پواسون واضح بود، هیچ تضمینی برای راه حل به شکل بسته برای w وجود ندارد. بنابراین، فرمول های GLM معمولاً از تکنیک های تکرار شونده مانند گرادیان کاهشی استفاده می کنند. از این رو، یک مکانیسم واحد را می توان برای طیف گسترده ای از توابع پیوند و توزیع احتمال، با استفاده از این گرادیان های بالا استفاده کرد.

¹ logistic

² sigmoid

این به روزرسانی را می‌توان با استفاده از رایج‌ترین تنظیمات برای GLM ها ملموس‌تر کرد. نکته مهم این است که این تنظیم فقط به دانش تابع انتقال f نیاز دارد، بدون اینکه صریحاً نیاز به دانستن $\log - \text{normalized } a$ داشته باشد. این ساده سازی از ارتباط بین انتقال f و $\log - \text{normalizer } a$ که در بالا اشاره شد ناشی می‌شود. ما بحث کردیم که تابع انتقال f برای منعکس کردن محدوده متغیر خروجی y انتخاب شده است. با این حال، انتخاب باید ویژگی‌های دیگری نیز داشته باشد. به طور خاص، ما می‌خواهیم اطمینان حاصل کنیم که g یک احتمال Log منفی محدب و هموار ارائه می‌کند تا بهینه‌سازی را ساده کنیم. به طور مفید، پارامتر a از توزیع خانواده نمایی دقیقاً چنین انتخابی را در اختیار ما قرار می‌دهد: $f = \nabla a$. از آنجایی که $\frac{\partial \theta_i}{\partial w_j} = x_{ij}$ برای $\theta_i = x_i^T w$ ، دریافت می‌کنیم که

$$\begin{aligned} -\frac{\partial ll_i(w)}{\partial w_j} &= \left(\frac{\partial a(\theta_i)}{\partial \theta_i} - y_i \right) \frac{\partial \theta_i}{\partial w_j} \\ &= (f(\theta_i) - y_i) x_{ij} \\ &= (f(x_i^T w) - y_i) x_{ij} \end{aligned}$$

بنابراین، با توجه به انتقال مناسب f برای توزیع خانواده نمایی مورد نظر، به‌روزرسانی گرادیان کاهشی تصادفی به سادگی انجام می‌شود.

$$w_{t+1} = w_t - \eta_t (f(x_i^T w_t) - y_i) x_i$$

و به روزرسانی گرادیان کاهشی دسته‌ای است

$$\begin{aligned} w_{t+1} &= w_t - \eta_t \sum_{i=1}^n (f(x_i^T w_t) - y_i) x_i \\ &= w_t - \eta_t X^T (p - y) \end{aligned}$$

که در آن $p_i = f(x_i^T w_t)$ برای بررسی هسین، دومین مشتق جزئی تابع احتمال ورود به سیستم منفی برای یک نمونه است.

$$-\frac{\partial^2 ll_i(w)}{\partial w_j \partial w_k} = x_{ij} \frac{\partial f(\theta_i)}{\partial \theta_i} x_{ik}$$

برای D یک ماتریس مورب $n \times n$ با $\frac{\partial f(\theta_i)}{\partial \theta_i}$ در قطر، ماتریس هسی بنابراین

$$H_{-ll(w)} = -X^T D X. \quad (7.3)$$

همانطور که در رگرسیون پواسون، این ماتریس تضمین شده است که مثبت نیمه معین است، و اگر X رتبه پایین نباشد، مثبت قطعی است.

نکته: تنظیم رایج $f = \nabla a$ برای GLM ها با اهداف کاربردی به نام واگرایی برگمن¹ ارتباط دارد. این واگرایی‌ها به صورت $D_a(\hat{y}||y)$ نوشته می‌شوند، که نشان دهنده تفاوت بین \hat{y} و y است، جایی که واگرایی با a پارامتر می‌شود. به مینیمم رساندن این واگرایی برگمن با به مینیمم رساندن احتمال ورود به سیستم منفی خانواده نمایی مربوطه مطابقت دارد:

¹ Bregman divergences

$$\operatorname{argmin}_{\theta} D_a(x||g(\theta)) = \operatorname{argmin}_{\theta} -\ln p(x|\theta)$$

برای جزئیات بیشتر در مورد این رابطه به [۲۰، بخش ۲،۲] و [۲] مراجعه کنید.

توجه داشته باشید که پیوند انتخاب شده لزوماً نباید با مشتق a مطابقت داشته باشد. در عوض، این مکانیسمی را برای اطمینان از یک تابع از دست دادن خوب فراهم می‌کند، زیرا واگرایی‌های برگمن ویژگی‌های خوبی دارند، از جمله محدب بودن در آرگومان اول. با این حال، این بدان معنا نیست که هر پیوند دیگری لزوماً منجر به عملکرد زیان نامطلوب می‌شود.

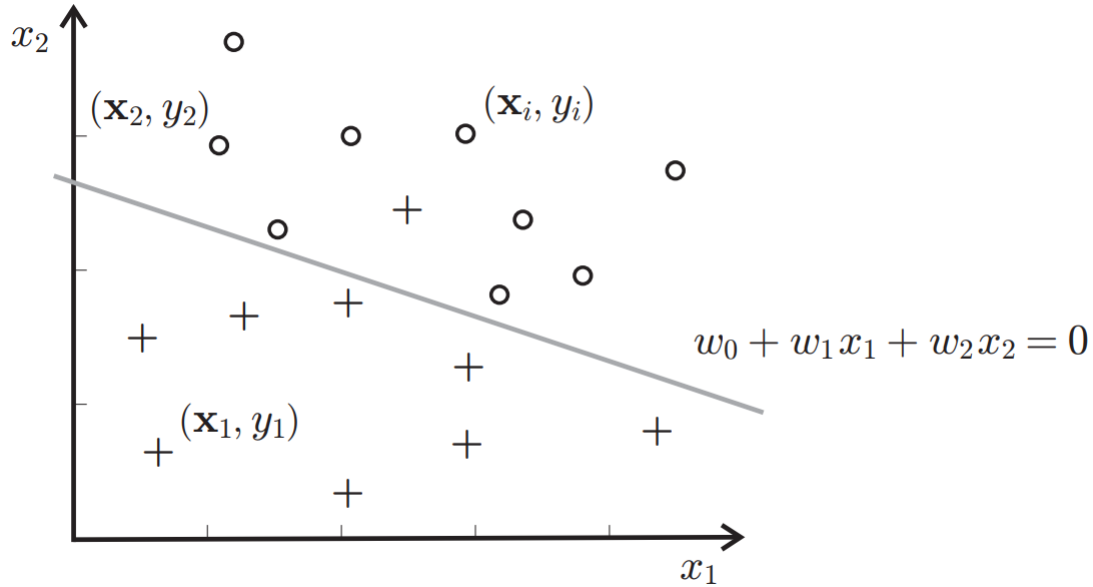
فصل ۸

طبقه‌بندی‌های خطی

فرض کنید ما علاقه‌مند به ساخت یک طبقه‌بندی کننده خطی $f : \mathbb{R}^d \rightarrow \{-1, +1\}$ هستیم. طبقه‌بندی خطی برای یافتن رابطه بین ورودی‌ها و خروجی‌ها با ساخت یک تابع خطی (نقطه، خط، صفحه یا ابر صفحه) که \mathbb{R}^d را به دو نیم فاصله تقسیم می‌کند. دو نیمه فضا به ترتیب به عنوان مناطق تصمیم‌گیری برای مثال‌های مثبت و منفی عمل می‌کنند. با توجه به مجموعه داده‌ای $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ متشکل از مثال‌های مثبت و منفی، روش‌های زیادی وجود دارد که طبقه‌بندی کننده‌های خطی را می‌توان ساخت. به عنوان مثال، یک الگوریتم آموزشی ممکن است به طور صریح برای موقعیت‌یابی سطح تصمیم به منظور جداسازی مثال‌های مثبت و منفی بر اساس برخی معیارهای مرتبط با مشکل کار کند. به عنوان مثال، ممکن است سعی کند کسر نمونه‌ها را در سمت نادرست سطح تصمیم به حداقل برساند. از طرف دیگر، هدف الگوریتم آموزشی ممکن است تخمین مستقیم توزیع پسینی $p(y|x)$ باشد، در این صورت الگوریتم به احتمال زیاد بر اصول تخمین پارامتر رسمی تکیه می‌کند. به عنوان مثال، ممکن است احتمال را به حداکثر برساند. نمونه‌ای از یک طبقه‌بندی کننده با سطح تصمیم‌گیری خطی در شکل ۸.۱ نشان داده شده است

برای ساده‌سازی فرمالیسم در بخش‌های بعدی، یک جزء $x_0 = 1$ به هر ورودی اضافه می‌کنیم (x_1, \dots, x_d) . این فضای ورودی را به $\mathcal{X} = \mathbb{R}^{d+1}$ گسترش می‌دهد، اما خوشبختانه ما را به یک نماد ساده‌شده نیز هدایت می‌کند که در آن مرز تصمیم‌گیری در \mathbb{R}^d می‌تواند به صورت $\mathbf{w}^T \mathbf{x} = 0$ نوشته شود، که در آن $\mathbf{w} = (w_0, w_1, \dots, w_d)$ مجموعه‌ای از وزن‌ها است و $\mathbf{x} = (x_0 = 1, x_1, \dots, x_d)$ هر عنصر فضای ورودی است. با این وجود، باید به یاد داشته باشیم که ورودی‌های واقعی d -بعدی هستند.

قبلاً در اظهارات مقدماتی، ما یک طبقه‌بندی کننده را به عنوان تابع $f : \mathcal{X} \rightarrow \mathcal{Y}$ ارائه کردیم و مسئله یادگیری را به $p(y|x)$ تقریبی تبدیل کردیم. در شرایطی که طبقه‌بندی کننده‌های خطی، انعطاف‌پذیری ما محدود است، زیرا روش ما باید احتمالات پسینی $p(y|x)$ را بیاموزد و در عین حال سطح تصمیم‌گیری خطی در \mathbb{R}^d داشته باشد. با این حال، اگر $p(y|x)$ به عنوان تابع یکنواخت $\mathbf{w}^T \mathbf{x}$ مدل‌سازی شود، می‌توان به این امر دست یافت. به عنوان مثال، $\tanh(\mathbf{w}^T \mathbf{x})$ یا $(1 + e^{-\mathbf{w}^T \mathbf{x}})^{-1}$. البته، یک مدل آموزش دیده برای یادگیری احتمالات پسین $p(y|x)$ را می‌توان به عنوان یک پیش‌بینی کننده "نرم" یا یک تابع امتیاز دهی $s : \mathcal{X} \rightarrow [0, 1]$ مشاهده کرد. سپس، تبدیل از s به f یک کاربرد مستقیم از حداکثر اصل پسینی است: خروجی پیش‌بینی شده مثبت است اگر $s(x) \geq 0.5$ باشد و اگر $s(x) < 0.5$ باشد، منفی است. به طور کلی، تابع امتیازدهی می‌تواند هر نقشه برداری باشد $s : \mathcal{X} \rightarrow \mathbb{R}$ ، با آستانه اعمال شده بر اساس هر مقدار خاص τ .



شکل ۸/۱: مجموعه‌ای از داده‌ها در \mathbb{R}^2 شامل نه مثال مثبت و نه مثال منفی است. خط خاکستری نشان دهنده سطح تصمیم‌گیری خطی در \mathbb{R}^2 است. سطح تصمیم به طور کامل نکات مثبت را از منفی جدا نمی‌کند.

۸-۱ رگرسیون لجستیک^۱

اجازه دهید طبقه‌بندی باینری را در \mathbb{R}^d در نظر بگیریم، که در آن $\mathcal{X} = \mathbb{R}^{d+1}$ و $Y = \{0, 1\}$. رگرسیون لجستیک یک مدل خطی تعمیم یافته است، که در آن توزیع بر روی Y داده شده \mathbf{x} یک توزیع برنولی است و تابع انتقال تابع سیگموئید است که تابع لجستیک نیز نامیده می‌شود.

$$\sigma(t) = (1 + e^{-t})^{-1}$$

در شکل ۸،۲ ترسیم شده است. در اصطلاحات مشابه GLM ها، تابع انتقال سیگموئید است و تابع پیوند - معکوس تابع انتقال - تابع logit است $\text{logit}(x) = \ln \frac{x}{1-x}$ ، با

$$1. E[y|\mathbf{x}] = \sigma(\mathbf{w}^T \mathbf{x})$$

$$2. p(y|\mathbf{x}) = \text{Bernoulli}(\alpha) \text{ with } \alpha = E[y|\mathbf{x}]$$

توزیع برنولی با α تابعی از \mathbf{x} است

$$p(y|\mathbf{x}) = \begin{cases} \left(\frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \right)^y & \text{for } y = 1 \\ \left(1 - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \right)^{1-y} & \text{for } y = 0 \end{cases}$$

$$= \sigma(\mathbf{x}^T \mathbf{w})^y (1 - \sigma(\mathbf{x}^T \mathbf{w}))^{1-y}$$

¹ Logistic regression

که در آن $\omega = (\omega_0, \omega_1, \dots, \omega_d)$ مجموعه‌ای از ضرایب مجهول است که می‌خواهیم بازیابی کنیم (یا یاد بگیریم). توجه کنید که پیش‌بینی ما $\sigma(\omega^T \mathbf{x})$ است و برآورده می‌شود

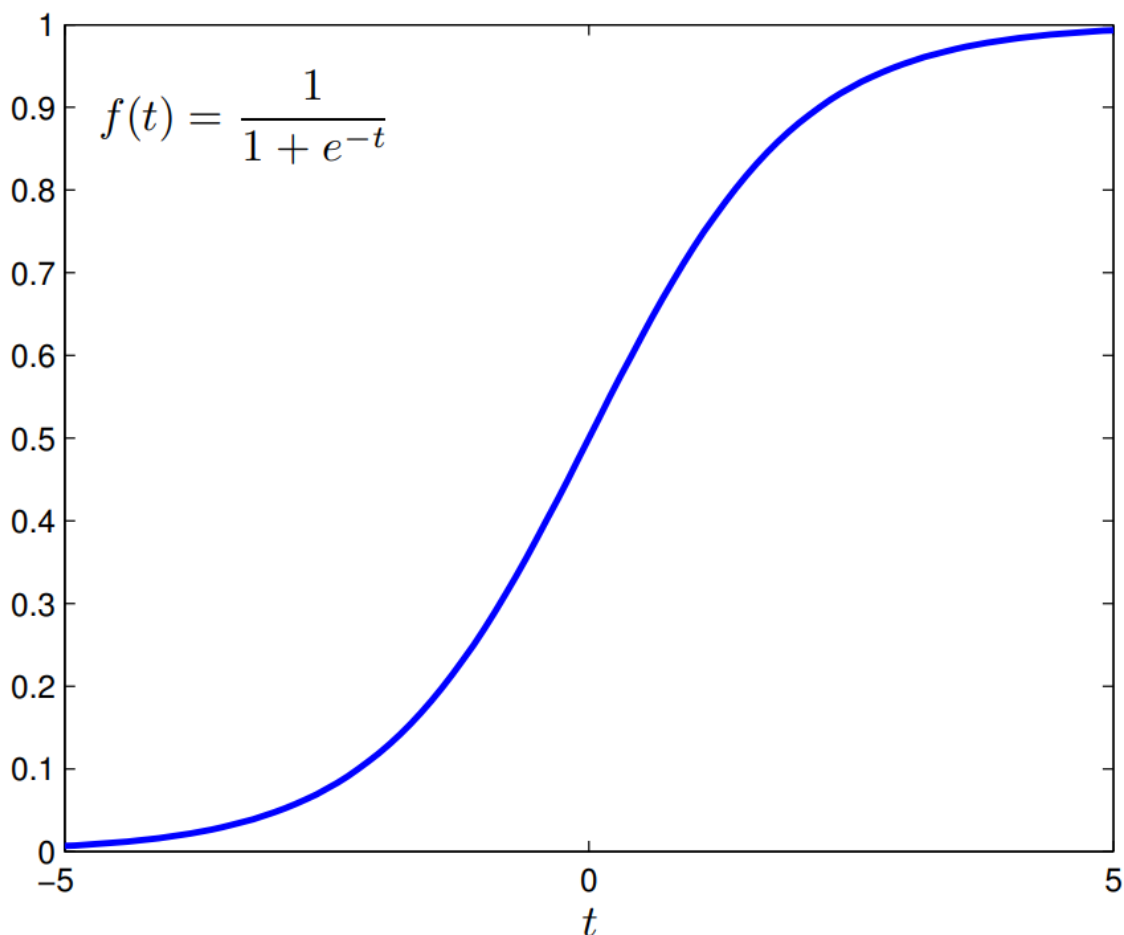
$$p(y = 1|\mathbf{x}) = \sigma(\omega^T \mathbf{x})$$

بنابراین، مانند بسیاری از رویکردهای طبقه‌بندی باینری، هدف ما پیش‌بینی احتمال 1 بودن کلاس است. با توجه به این احتمال، می‌توانیم $p(y = 0|\mathbf{x}) = 1 - p(y = 1|\mathbf{x})$ را استنباط کنیم.

۸-۱-۱ پیش‌بینی برچسب‌های کلاس

تابعی که توسط رگرسیون لجستیک آموخته می‌شود، به جای یک پیش‌بینی صریح 0 یا 1، یک احتمال را بر می‌گرداند. بنابراین، ما باید این تخمین احتمال را بگیریم و آن را به یک پیش‌بینی مناسب کلاس تبدیل کنیم. برای یک نقطه داده‌ای که قبلاً دیده نشده بود \mathbf{x} و مجموعه‌ای از ضرایب آموخته شده \mathbf{w} ، به سادگی احتمال پسین را به این صورت محاسبه می‌کنیم.

$$P(Y = 1|\mathbf{x}, \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$



شکل ۸/۲: تابع سیگموئید در بازه $[-5, 5]$.

اگر $P(Y = 1|x, w) \geq 0.5$ نتیجه می‌گیریم که نقطه داده x باید مثبت باشد ($\hat{y} = 1$). در غیر این صورت، اگر $P(Y = 1|x, w^*) < 0.5$ ، نقطه داده را منفی می‌گذاریم ($\hat{y} = 0$). پیش‌بینی کننده یک بردار $(d + 1)$ بعدی $x = (x_0 = 1, x_1, \dots, x_d)$ را به صفر یا یک نگاشت می‌کند.

با این حال، همانطور که در فصل ۱۰ بحث می‌کنیم، این آستانه نباید 0.5 باشد. در برخی موارد، ممکن است فرد بیشتر به عدم شناسایی مثبت اهمیت دهد (به عنوان مثال، ناتوانی در شناسایی یک بیماری). در چنین حالتی، ممکن است اشتباه کردن در آستانه کوچکتر ایمن‌تر باشد، به طوری که نمونه‌های بیشتری به عنوان مثبت برچسب گذاری شوند. علاوه بر این، مقادیر احتمال خود می‌توانند منعکس کننده باشند: حتی اگر هر دو طبقه‌بندی کننده دقت خوبی داشته باشند، ترجیح داده می‌شود طبقه‌بندی کننده‌ای داشته باشیم که به طور مداوم احتمالات نزدیک به 0.9 و 0.1 را تولید کند، به جای احتمالات با اطمینان کمتر که در اطراف 0.5 قرار دارند. دلیل این امر این است که انتظار می‌رود اغتشاش‌های کوچک تأثیر بیشتری بر طبقه‌بندی کننده دوم داشته باشند، که می‌تواند به طور ناگهانی برچسب گذاری را در یک نمونه تغییر دهد. برای تعیین کمیت این جنبه از پیش‌بینی، معیارهای دیگری مانند منحنی عملیاتی گزارش شده ترجیح داده می‌شوند که در بخش ۱۰.۴ در مورد آن بحث می‌کنیم. در حال حاضر، این آستانه ساده‌تر را فرض می‌کنیم و ارزیابی پیشرفته‌تر الگوریتم‌های طبقه‌بندی را به بعد واگذار می‌کنیم.

توجه داشته باشید که طبقه‌بندی کننده رگرسیون لجستیک یک طبقه‌بندی کننده خطی است، علیرغم اینکه سیگموئید غیرخطی است. این به این دلیل است که $P(Y = 1|x, w) \geq 0.5$ فقط زمانی که $w^T x \geq 0$ باشد. عبارت $w^T x = 0$ معادله یک ابر صفحه را نشان می‌دهد که مثال‌های مثبت و منفی را از هم جدا می‌کند.

۸-۱-۲ برآورد حداکثر احتمال برای رگرسیون لجستیک

از آنجایی که رگرسیون لجستیک یک GLM است، می‌توانیم از الگوریتم گرادین کاهشی عمومی مشتق شده در بخش ۷.۳ با انتقال $f = \sigma$ ، با گرادین زیر از احتمال ورود به سیستم منفی در هر نمونه استفاده کنیم.

$$-\frac{\partial \text{ll}_i(w)}{\partial w_j} = (\sigma(x_i^T w) - y_i)x_{ij}$$

حتی اگر به روزرسانی نهایی را می‌دانیم، به عنوان یک تمرین، به صراحت راه حل حداکثر احتمال را استخراج خواهیم کرد. مانند قبل، فرض کنید که مجموعه داده $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ یک $i.i.d.$ است. نمونه از یک توزیع احتمال ثابت اما ناشناخته $p(x, y) = p(y|x)p(x)$ داده‌ها با رسم تصادفی یک نقطه x مطابق $p(x)$ تولید می‌شوند و سپس برچسب کلاس Y را مطابق توزیع برنولی در $(\lambda, 1)$ تنظیم می‌کنند. احتمال ورود منفی برای این مجموعه داده $\text{ll}(w) = \sum_{i=1}^n -\text{ll}_i(w)$ است که در آن

$$\begin{aligned} \text{ll}_i(w) &= \log p(y_i | x) \\ &= y_i \log \sigma(w^T x_i) + (1 - y_i) \log (1 - \sigma(w^T x_i)) \\ &= \left(y_i \log \left(\frac{1}{1 + e^{-w^T x_i}} \right) + (1 - y_i) \log \left(1 - \frac{1}{1 + e^{-w^T x_i}} \right) \right) \end{aligned}$$

منفی $\log - \text{likelihood}$ با این ll_i معمولاً به عنوان آنتروپی متقاطع^۱ نامیده می‌شود

^۱ cross-entropy

از اینجا، می‌توانید مشتق هر جزء را در این مجموع، با استفاده از قانون زنجیره‌ای برای سیگموئید، بگیرید. برای جزء اول، با $p_i = \sigma(\theta_i)$

$$\begin{aligned}\frac{\partial y_i \log \sigma(\mathbf{w}^T \mathbf{x}_i)}{\partial w_j} &= y_i \frac{\partial \log \sigma(\theta_i)}{\partial \theta_i} \frac{\partial \theta_i}{\partial w_j} \\ &= \frac{\partial \log p_i}{\partial p_i} \frac{\partial p_i}{\partial \theta_i} \frac{\partial \theta_i}{\partial w_j} \\ &= y_i \frac{1}{p_i} \frac{\partial p_i}{\partial \theta_i} \frac{\partial \theta_i}{\partial w_j} \\ &= y_i \frac{1}{p_i} \sigma(\theta_i)(1 - \sigma(\theta_i)) x_{ij} \\ &= y_i(1 - \sigma(\theta_i)) x_{ij}\end{aligned}$$

زیرا

$$\frac{\partial \sigma(\theta_i)}{\partial \theta_i} = \sigma(\theta_i)(1 - \sigma(\theta_i))$$

شما می‌توانید این مرحله را برای خودتان تأیید کنید، اما به صراحت تعریف σ را وصل کنید. برای جزء دوم، با دنبال کردن مراحل مشابه، دریافت می‌کنیم

$$\frac{\partial (1 - y_i) \log (1 - \sigma(\mathbf{w}^T \mathbf{x}_i))}{\partial w_j} = (y_i - 1) \sigma(\theta_i) x_{ij}$$

با جمع این موارد و گرفتن منفی، به گرادیان $(p_i - y_i) x_{ij}$ می‌رسیم.

برای تمرین بیشتر گرفتن شیب از این اهداف، می‌توانیم قبل از گرفتن گرادیان، هدف را کمی تغییر دهیم. این به مسیر دیگری برای استخراج قانون به‌روزرسانی برای رگرسیون لجستیک منجر می‌شود که اکنون از آن عبور می‌کنیم. ابتدا توجه کنید که

$$\left(1 - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}}\right) = \frac{e^{-\mathbf{w}^T \mathbf{x}_i}}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}}$$

که میدهد

$$\begin{aligned}ll(\mathbf{w}) &= \sum_{i=1}^n \left(-y_i \cdot \log(1 + e^{-\mathbf{w}^T \mathbf{x}_i}) + (1 - y_i) \cdot \log(e^{-\mathbf{w}^T \mathbf{x}_i}) - (1 - y_i) \cdot \log(1 + e^{-\mathbf{w}^T \mathbf{x}_i}) \right) \\ &= \sum_{i=1}^n \left((y_i - 1) \mathbf{w}^T \mathbf{x}_i + \log \left(\frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}} \right) \right)\end{aligned}$$

باز هم، بر خلاف رگرسیون خطی، واضح است که هیچ راه‌حل بسته‌ای برای $\nabla ll(\mathbf{w}) = 0$ وجود ندارد. بنابراین، ما باید با روش‌های بهینه‌سازی تکراری پیش برویم. ما \mathbf{w}_0 را معمولاً به یک بردار تصادفی یا به طور بالقوه با راه‌حل رگرسیون خطی که نقطه اولیه بسیار بهتری را ارائه می‌دهد مقداردهی اولیه می‌کنیم. از آنجایی که هدف محدب است، مقداردهی اولیه تنها بر تعداد مراحل تأثیر می‌گذارد، اما نباید مانع از همگرایی شیب نزولی به مینیمم مطلق شود. به روزرسانی گرادیان کاهشی تصادفی است

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t (\sigma(\mathbf{x}_i^T \mathbf{w}_t) - y_i) \mathbf{x}_i$$

و به روزرسانی گرادیان کاهش دهنده است

$$\begin{aligned} \mathbf{w}_{t+1} &= \mathbf{w}_t - \eta_t \sum_{i=1}^n (\sigma(\mathbf{x}_i^T \mathbf{w}_t) - y_i) \mathbf{x}_i \\ &= \mathbf{w}_t - \eta_t \mathbf{X}^T (\sigma(\mathbf{X} \mathbf{w}_t) - \mathbf{y}) \end{aligned}$$

جایی که ما تعریف σ را هنگام اعمال بر یک بردار اضافه بارگذاری می‌کنیم به این معنی که به طور جداگانه برای هر عنصر در آن بردار اعمال می‌شود: $\sigma(v) = [\sigma(v_1), \dots, \sigma(v_n)]$. $\mathbf{H}_{-ll(w)} = -\mathbf{X}^T \mathbf{D} \mathbf{X}$ هسین است که در آن \mathbf{D} یک ماتریس مورب $n \times n$ با $p_i(1 - p_i)$ در قطر است (نگاه کنید به (۷,۳))

تابع درستنمایی شرطی وزنی

در شرایط خاص، ممکن است توجیهی باشد که به شما اهمیت یکسانی برای هر نقطه داده، داده شود. این تابع احتمال شرطی را تغییر می‌دهد

$$l(\mathbf{w}) = \prod_{i=1}^n p_i^{c_i y_i} \cdot (1 - p_i)^{c_i (1 - y_i)}$$

که در آن $0 \leq c_i \leq 1$ هزینه برای نقطه داده i است. با در نظر گرفتن $\mathbf{C} = \text{diag}(c_1, c_2, \dots, c_n)$ اکنون می‌توانیم گرادیان $\log - \text{relihood}$ منفی را به این صورت بیان کنیم

$$-\nabla l(\mathbf{w}) = \mathbf{X}^T \mathbf{C} (\mathbf{p} - \mathbf{y})$$

و هسین را به عنوان

$$\mathbf{H}_{-ll(w)} = -\mathbf{X}^T \mathbf{C} \mathbf{P} (\mathbf{I} - \mathbf{P}) \mathbf{X}$$

مشاهده این نکته جالب است که هسین به صورت نیمه قطعی مثبت باقی می‌ماند. بنابراین، انتظار می‌رود قانون به روزرسانی به مینیمم مطلق همگرا شود.

۸-۱-۳ مسائل مربوط به مینیمم کردن فاصله اقلیدسی

یک سوال طبیعی این است که چرا ما این مسیر را برای طبقه‌بندی خطی دنبال کردیم. به جای اینکه صریحاً فرض کنیم $\sigma(\mathbf{x}^T \mathbf{w}) = E[Y|\mathbf{x}]$ برای $P(Y = 1|\mathbf{x}, \mathbf{w})$ یک توزیع برنولی است و راه‌حل حداکثر درستنمایی را برای $P(Y = 1|\mathbf{x}, \mathbf{w})$ محاسبه کنیم. می‌توانستیم به سادگی تصمیم بگیریم از $\sigma(\mathbf{x}^T \mathbf{w})$ برای پیش‌بینی اهداف $y \in \{0, 1\}$ استفاده کنیم و سپس سعی کنیم تفاوت آنها را با استفاده از هزینه مورد علاقه خود (هزینه مجذور) به حداقل برسانیم. متأسفانه، این مشخصات مسئله تصادفی‌تر منجر به بهینه‌سازی غیر محدب می‌شود. در واقع، نتیجه‌ای وجود دارد که استفاده از خطای اقلیدسی برای انتقال سیگموئید به طور تصاعدی مینیمم‌های محلی زیادی در تعداد ویژگی‌ها به دست می‌دهد [۱]. برای علاقه به این مسیر جایگزین، نشان خواهیم داد که این جهت به یک بهینه‌سازی غیر محدب منجر می‌شود

اجازه دهید تابع خطا با فاصله اقلیدسی اکنون به این صورت نوشته شود

$$\text{Err}(\mathbf{w}) = \sum_{i=1}^n (y_i - p_i)^2 \quad \triangleright p_i = \sigma(\mathbf{x}_i^T \mathbf{w}), e_i = y_i - p_i$$

کمینه سازی $Err(\mathbf{w})$ به طور فرمول بندی شده به این صورت بیان می شود

$$\begin{aligned}\mathbf{w}^* &= \arg \min_{\mathbf{w}} \{Err(\mathbf{w})\} \\ &= \arg \min_{\mathbf{w}} \left\{ \sum_{i=1}^n (y_i - p_i)^2 \right\}\end{aligned}$$

مشابه فرآیند حداکثر احتمال، هدف ما محاسبه بردار گرادیان و هسین تابع خطا خواهد بود. مشتقات جزئی تابع خطا را می توان به صورت زیر محاسبه کرد

$$\begin{aligned}\frac{\partial Err(\mathbf{w})}{\partial w_j} &= \frac{\partial}{\partial w_j} \sum_{i=1}^n e_i^2 \\ &= \sum_{i=1}^n 2 \cdot e_i \cdot \frac{\partial e_i}{\partial w_j} \\ &= 2 \cdot \sum_{i=1}^n \left(y_i - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}} \right) \cdot \frac{1}{(1 + e^{-\mathbf{w}^T \mathbf{x}_i})^2} \cdot e^{-\mathbf{w}^T \mathbf{x}_i} \cdot (-x_{ij}) \\ &= 2 \cdot \sum_{i=1}^n x_{ij} \cdot \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}} \cdot \left(1 - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}} \right) \cdot \left(y_i - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}} \right) \\ &= -2 \mathbf{f}_j^T \mathbf{P} (\mathbf{I} - \mathbf{P}) (\mathbf{y} - \mathbf{p})\end{aligned}$$

این بردار گرادیان را به شکل زیر ارائه می دهد

$$\nabla Err(\mathbf{w}) = -2 \mathbf{X}^T \mathbf{P} (\mathbf{I} - \mathbf{P}) (\mathbf{y} - \mathbf{p})$$

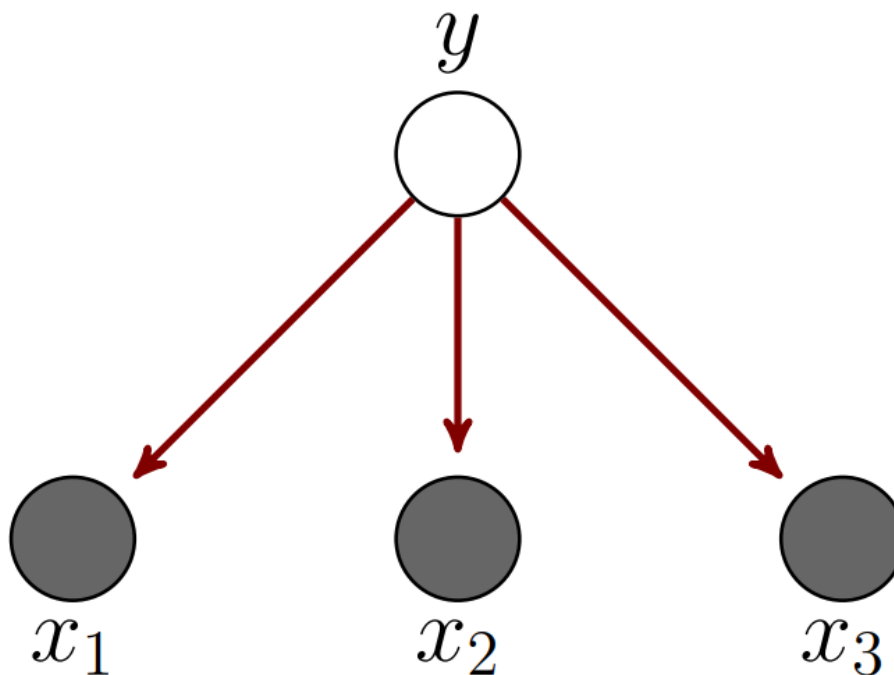
ماتریس $\mathbf{J} = \mathbf{P} (\mathbf{I} - \mathbf{P}) \mathbf{X}$ به عنوان ژاکوبین¹ نامیده می شود. به طور کلی، ژاکوبین یک ماتریس $n \times d$ است که به صورت محاسبه می شود

$$\begin{bmatrix} \frac{\partial e_1}{\partial w_1} & \frac{\partial e_1}{\partial w_2} & \dots & \frac{\partial e_1}{\partial w_d} \\ \vdots & \ddots & & \\ \frac{\partial e_n}{\partial w_1} & & & \frac{\partial e_n}{\partial w_d} \end{bmatrix}$$

دومین مشتق جزئی تابع خطا را می توان به این صورت پیدا کرد

$$\begin{aligned}\frac{\partial^2 Err(\mathbf{w})}{\partial w_j \partial w_d} &= 2 \cdot \sum_{i=1}^n \frac{\partial e_i}{\partial w_d} \cdot \frac{\partial e_i}{\partial w_j} + e_i \cdot \frac{\partial^2 e_i}{\partial w_j \partial w_d} \\ &= 2 \cdot \sum_{i=1}^n x_{ij} \cdot \left(p_i^2 (1 - p_i)^2 + p_i \cdot (1 - p_i) \cdot (2p_i - 1) \cdot (y_i - p_i) \right) \cdot x_{ik}\end{aligned}$$

¹ Jacobian



شکل ۸-۳: مدل گرافیکی Naive Bayes، با سه ویژگی

بنابراین، هسین را می‌توان به این صورت محاسبه کرد

$$\mathbf{H}_{\text{Err}(\mathbf{w})} = 2\mathbf{X}^T(\mathbf{I} - \mathbf{P})^T\mathbf{P}^T\mathbf{P}(\mathbf{I} - \mathbf{P})\mathbf{X} + 2\mathbf{X}^T(\mathbf{I} - \mathbf{P})^T\mathbf{P}^T\mathbf{E}(2\mathbf{P} - \mathbf{I})\mathbf{X} = 2\mathbf{J}^T\mathbf{J} + 2\mathbf{J}^T\mathbf{E}(2\mathbf{P} - \mathbf{I})\mathbf{X}$$

که در آن $\mathbf{P} = \text{diag}\{\mathbf{p}\}$, $\mathbf{E} = \text{diag}\{\mathbf{e}\}$ یک ماتریس مورب حاوی عناصر $E_{ii} = e_i = y_i - p_i$ و \mathbf{I} یک ماتریس همانی است.

اکنون می‌بینیم که هسین تضمینی برای مثبت بودن نیمه قطعی نیست. این بدان معنی است که $\text{Err}(\mathbf{w})$ محدب نیست، یعنی باید مینیمم‌های متعدد با مقادیر مختلف تابع هدف داشته باشد. یافتن یک بهینه مطلق بستگی به این دارد که راه‌حل اولیه $\mathbf{w}^{(0)}$ چقدر مطلوب است و چگونه مرحله به روزرسانی وزن می‌تواند از مینیمم‌های محلی برای یافتن راه‌حل‌های بهتر فرار کند. با این حال، به حداقل رساندن این تابع غیر محدب، بسیار مشکل‌سازتر از آنتروپی متقابل محدب خواهد بود.

۸-۲ طبقه‌بندی کننده بیز ساده

طبقه‌بندی بیز ساده یک رویکرد مولد برای پیش‌بینی است. تا کنون، ما در مورد رویکردهای افتراقی^۱ (رگرسیون خطی، رگرسیون لجستیک) بحث کرده‌ایم که سعی در یادگیری $p(y|\mathbf{x})$ دارند. برای یک تنظیم مولد، $p(\mathbf{x}, y) = p(\mathbf{x}|y)p(y)$ را یاد می‌گیریم. همانطور که می‌توانید تصور کنید، این می‌تواند یک کار دشوارتر باشد، زیرا ما همچنین نیاز داریم که توزیع را روی

¹ discriminative

خود ویژگی‌ها بیاموزیم. برای بیز ساده ، ما به طور قابل توجهی به سادگی این توزیع مشترک را با ایجاد یک فرض قوی یاد می‌گیریم: ویژگی‌ها با توجه به برچسب مستقل هستند. این فرض با مدل گرافیکی در شکل ۸،۳ نشان داده شده است.

همانند طبقه‌بندی‌کننده‌های متمایز، مانند رگرسیون لجستیک، قانون تصمیم‌گیری برای برچسب‌گذاری یک نقطه به عنوان کلاس c (یعنی $y = 1$) است.

$$\operatorname{argmax}_{c \in \{1, \dots, k\}} p(y = c | \mathbf{x})$$

که در آن $p(y = 1 | \mathbf{x}) \propto p(\mathbf{x} | y = 1)p(y = 1)$ برای دو کلاس، اگر $p(y = 1 | \mathbf{x}) \geq p(y = 0 | \mathbf{x})$ باشد، این به انتخاب کلاس 1 مربوط می‌شود.

برای شروع، تنظیمات ساده‌تری را با ویژگی‌های باینری در نظر می‌گیریم و سپس به ویژگی‌های پیوسته می‌پردازیم. توجه داشته باشید که بیز ساده یک طبقه‌بندی خطی برای ویژگی‌های باینری است. با این حال، به طور کلی‌تر، لزوماً یک طبقه‌بندی خطی نیست. توجه داشته باشید که یک طبقه‌بندی خطی طبقه‌بندی‌کننده‌ای است که در آن دو کلاس با یک صفحه خطی از هم جدا می‌شوند، یعنی مرز تصمیم‌گیری بر اساس ترکیب خطی ویژگی‌ها است.

۸-۲-۱ ویژگی‌های باینری و طبقه‌بندی خطی

اجازه دهید $\mathcal{D} = \{(x_i, y)\}_{i=1}^n$ یک مجموعه داده ورودی باشد، که در آن $\mathcal{X} = \{0, 1\}^d$ و $\mathcal{Y} = \{0, 1\}$. تحت فرض بیز ساده ، ویژگی‌ها با توجه به برچسب کلاس مستقل هستند. بنابراین، ما می‌توانیم بنویسیم

$$p(\mathbf{x} | y) = \prod_{i=1}^d p(x_i | y)$$

یک انتخاب مناسب برای این احتمالات تک متغیره ساده‌تر، توزیع برنولی است، زیرا هر x_j باینری است،

$$p(x_j | y = c) = p_{j,c}^{x_j} (1 - p_{j,c})^{1-x_j}$$

پارامترهای توزیع برنولی $p_{j,c} = p(x_j = 1 | y = c)$ با یک پارامتر مختلف $p_{j,c}$ برای هر کلاس c و برای هر ویژگی است. با محاسبه به راحتی می‌توانیم این پارامتر را از داده‌ها یاد بگیریم

$$p_{j,c} = \frac{\text{number of times } x_j = 1 \text{ for class } c}{\text{number of datapoints labeled as class } c}$$

به طور مشابه، ما می‌توانیم با استفاده از $p_c = p(y = c)$ قبلی یاد بگیریم

$$p_c = p(y = c) = \frac{\text{number of datapoint labeled as class } c}{\text{total number of datapoints}}$$

توجه داشته باشید که این رویکرد می‌تواند برای بیشتر از دو کلاس نیز انجام شود. پیش‌بینی در نقطه جدید \mathbf{x} پس از آن اینگونه است

$$\max_{c \in \mathcal{Y}} p(y = c | \mathbf{x}) = \max_{c \in \mathcal{Y}} p(\mathbf{x} | y = c)p(y = c)$$

$$\begin{aligned}
&= \max_{c \in \mathcal{Y}} \prod_{j=1}^d p(x_j|y=c) p(y=c) \\
&= \max_{c \in \mathcal{Y}} \prod_{j=1}^d p_{j,c} p_c
\end{aligned}$$

تمرین ۱۰: راه حل بالا بصری است و از استخراج راه حل حداکثر احتمال به طور مشابه به دست می آید. با فرض اینکه شما n نقطه داده دارید و توزیع انتخاب شده $p(x_i|y)$ برنولی است همانطور که در بالا توضیح داده شد، حداکثر پارامترهای احتمال را برای $p_{j,c}, p_c$ for $j = 1, \dots, d, c = 0, 1$ استخراج کنید. برای ساده تر کردن کارها، از گزارش احتمال استفاده کنید

مرز طبقه بندی خطی با ویژگی ها و اهداف باینری

جالب اینجاست که طبقه بندی کننده بیز ساده با ویژگی های باینری و دو کلاس یک طبقه بندی کننده خطی است. این تا حدودی تعجب آور است، زیرا این رویکرد مولد بسیار متفاوت از آنچه قبلا انجام دادیم به نظر می رسد. برای اینکه ببینید چرا چنین است، توجه کنید که طبقه بندی کننده در چه زمانی تصمیم مثبت خواهد گرفت

$$p(y = 1|x) \geq p(y = 0|x)$$

آن موقع است که

$$p(x|y = 1)p(y = 1) \geq p(x|y = 0)p(y = 0)$$

ما این نماد را با استفاده از $p(y = 0) = p(0)$, $p(x|y = 0) = p(x|0)$ کوتاه می کنیم. با استفاده از فرض بیز ساده، اکنون داریم

$$p(1) \prod_{j=1}^d p(x_j|1) \geq p(0) \prod_{j=1}^d p(x_j|0)$$

که پس از اعمال لگاریتم تبدیل می شود به

$$\log p(1) + \sum_{j=1}^d \log p(x_j|1) \geq \log p(0) + \sum_{j=1}^d \log p(x_j|0)$$

اجازه دهید اکنون احتمالات شرطی کلاس $p(x_j|y)$ را بررسی کنیم، زمانی که $y \in \{0, 1\}$ است. به یاد بیاورید که هر ویژگی برنولی توزیع شده است، یعنی.

$$p(x_j|1) = p_{j,1}^{x_j} (1 - p_{j,1})^{1-x_j}$$

و

$$p(x_j|0) = p_{j,0}^{x_j} (1 - p_{j,0})^{1-x_j}$$

که در آن پارامترهای $p_{j,c}$ از مجموعه آموزشی تخمین زده می شوند. با گرفتن $p(y = c) = p_c$ داریم

$$\sum_{j=1}^d x_j \log \frac{p_{j,1}(1-p_{j,0})}{(1-p_{j,1})p_{j,0}} + \sum_{j=1}^d \log \frac{1-p_{j,1}}{1-p_{j,0}} + \log \frac{p_1}{p_0} \geq 0$$

می‌توانیم عبارت قبلی را به این صورت بنویسیم

$$w_0 + \sum_{j=1}^d w_j x_j \geq 0$$

که

$$w_0 = \log \frac{p_1}{p_0} + \sum_{j=1}^d \log \frac{1-p_{j,1}}{1-p_{j,0}}$$

$$w_j = \log \frac{p_{j,1}(1-p_{j,0})}{(1-p_{j,1})p_{j,0}} \quad j \in \{1, 2, \dots, d\}$$

بنابراین، در مورد ویژگی‌های باینری، بیز ساده یک طبقه‌بندی خطی است.

۸-۲-۲ بیز ساده پیوسته

برای ویژگی‌های پیوسته، توزیع برنولی دیگر برای $p(x_j|y)$ مناسب نیست و باید توزیع شرطی متفاوتی $p(\mathbf{x}|y)$ انتخاب کنیم. یک انتخاب متداول توزیع گاوسی است، اکنون با میانگین و واریانس متفاوت برای هر ویژگی و کلاس، $\mu_{j,c}, \sigma_{j,c}^2$:

$$p(x_j|y=c) = (2\pi\sigma_{j,c}^2)^{-1/2} \exp\left(-\frac{(x_j - \mu_{j,c})^2}{2\sigma_{j,c}^2}\right)$$

از آنجایی که \mathcal{Y} هنوز گسسته است، می‌توانیم $p(y)$ را با استفاده از شمارش‌های قبلی تقریب کنیم. پارامترهای میانگین و واریانس حداکثر درست‌نمایی با میانگین نمونه و کوواریانس نمونه برای هر کلاس به طور جداگانه مطابقت دارد. این شامل محاسبه میانگین و واریانس ویژگی j در نقاط داده برچسب‌گذاری شده با کلاس c است:

$$\mu_{j,c} = \frac{\sum_{i=1}^n 1(y_i = c) x_j}{\text{number of datapoints labeled as class } c}$$

$$\sigma_{j,c}^2 = \frac{\sum_{i=1}^n 1(y_i = c) (x_j - \mu_{j,c})^2}{\text{number of datapoints labeled as class } c}$$

تمرین ۱۱: فرمول حداکثر احتمال را برای یک مدل بیز ساده گاوسی استخراج کنید و بررسی کنید که آیا راه‌حل در واقع با میانگین نمونه و واریانس هر ویژگی و کلاس به طور جداگانه مطابق بالا مطابقت دارد.

۳-۸ رگرسیون لجستیک چند جمله ای

حال اجازه دهید طبقه‌بندی چندطبقه‌ای متمایز را در نظر بگیریم، که در آن $\mathcal{X} = \mathbb{R}^d$ و $\mathcal{Y} = \{1, 2, \dots, k\}$. این تنظیم به طور طبیعی در یادگیری ماشینی ایجاد می‌شود، جایی که اغلب بیش از دو دسته وجود دارد. به عنوان مثال، اگر بخواهیم گروه خونی (O و A ، B ، AB) یک فرد را پیش‌بینی کنیم، چهار کلاس داریم. در اینجا ما در مورد طبقه‌بندی چند کلاسه بحث می‌کنیم که در آن فقط می‌خواهیم یک نقطه داده را با یک کلاس از k برچسب گذاری کنیم. در تنظیمات دیگر، ممکن است بخواهید یک نقطه داده را با چندین کلاس برچسب گذاری کنید. در پایان این بخش به اختصار به این موضوع اشاره شده است.

ما می‌توانیم با استفاده از ایده چندجمله‌ای و تابع پیوند مربوطه، مانند دیگر مدل‌های خطی تعمیم‌یافته، به این تنظیمات تعمیم دهیم. توزیع چند جمله‌ای عضوی از خانواده نمایی است. ما می‌توانیم بنویسیم

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{y_1! \dots y_k!} p(y_1 = 1|\mathbf{x})^{y_1} \dots p(y_k = 1|\mathbf{x})^{y_k} \quad (8.3)$$

که در آن صورت معمولی $n! = 1$ زیرا $n = \sum_{j=1}^k y_j = 1$ زیرا ما فقط می‌توانیم یک مقدار کلاس داشته باشیم. مانند رگرسیون لجستیک، می‌توانیم $p(y_j = 1|\mathbf{x}) = \sigma(\mathbf{x}^T \mathbf{w}_j)$ را پارامتر کنیم. با این حال، ما همچنین باید اطمینان حاصل کنیم که $\sum_{j=1}^k p(y_j = 1|\mathbf{x}) = 1$ برای انجام این کار، حول کلاس نهایی، $p(y_k = 1|\mathbf{x}) = 1 - \sum_{j=1}^{k-1} p(y_j = 1|\mathbf{x})$ «پیوت»^۱ می‌کنیم. و فقط به طور صریح $\mathbf{w}_1, \dots, \mathbf{w}_{k-1}$ را یاد می‌گیرد. توجه داشته باشید که این مدل‌ها به طور مستقل یاد نمی‌گیرند، زیرا آنها با احتمال برای آخرین کلاس گره خورده‌اند. پارامترها را می‌توان نمایش داد. به عنوان یک ماتریس $\mathbf{W} \in \mathbb{R}^{d \times k}$ که در آن $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_k]$ از k بردار وزن با $\mathbf{w}_k = 0$ تشکیل شده است. خواهیم دید که چرا $\mathbf{w}_k = 0$ را ثابت می‌کنیم.

انتقال (معکوس پیوند) برای این تنظیم، انتقال سافت مکس است

$$\begin{aligned} \text{softmax}(\mathbf{x}^T \mathbf{W}) &= \left[\frac{\exp(\mathbf{x}^T \mathbf{w}_1)}{\sum_{j=1}^k \exp(\mathbf{x}^T \mathbf{w}_j)}, \dots, \frac{\exp(\mathbf{x}^T \mathbf{w}_k)}{\sum_{j=1}^k \exp(\mathbf{x}^T \mathbf{w}_j)} \right] \\ &= \left[\frac{\exp(\mathbf{x}^T \mathbf{w}_1)}{1^T \exp(\mathbf{x}^T \mathbf{W})}, \dots, \frac{\exp(\mathbf{x}^T \mathbf{w}_k)}{1^T \exp(\mathbf{x}^T \mathbf{W})} \right] \end{aligned}$$

و پیش‌بینی $\hat{y} = \text{softmax}(\mathbf{x}) \in [0, 1]^k$ است که در هر ورودی احتمال برچسب‌گذاری آن کلاس را می‌دهد، جایی که $\hat{\mathbf{y}}^T \mathbf{1} = 1$ نشان می‌دهد که مجموع احتمالات برابر با 1 است. توجه داشته باشید که این مدل شامل تنظیمات باینری برای رگرسیون لجستیک است، زیرا $\sigma(\mathbf{x}^T \mathbf{w}) = \frac{\exp(\mathbf{x}^T \mathbf{w})}{1 + \exp(\mathbf{x}^T \mathbf{w})}$ و $\sigma(\mathbf{x}^T \mathbf{w}) = (1 + \exp(-\mathbf{x}^T \mathbf{w}))^{-1}$. وزن رگرسیون لجستیک چند جمله‌ای با دو کلاس $\mathbf{W} = [\mathbf{w}, 0]$ است.

$$\begin{aligned} p(y = 0|\mathbf{x}) &= \frac{\exp(\mathbf{x}^T \mathbf{w})}{1^T \exp(\mathbf{x}^T \mathbf{W})} \\ &= \frac{\exp(\mathbf{x}^T \mathbf{w})}{\exp(\mathbf{x}^T \mathbf{w}) + \exp(\mathbf{x}^T \mathbf{0})} \end{aligned}$$

¹ pivot

$$= \frac{\exp(\mathbf{x}^T \mathbf{w})}{\exp(\mathbf{x}^T \mathbf{w}) + 1}$$

$$= \sigma(\mathbf{x}^T \mathbf{w})$$

به طور مشابه، برای $k > 2$ ، با ثابت کردن $\mathbf{w}_k = 0$ ، وزن‌های دیگر $\mathbf{w}_1, \dots, \mathbf{w}_{k-1}$ برای اطمینان از اینکه $p(y = k|x) =$

$$\sum_{j=1}^k p(y = j|x) = 1 \text{ یاد گرفته می‌شود و آن } \frac{\exp(\mathbf{x}^T \mathbf{w}_k)}{1^T \exp(\mathbf{x}^T \mathbf{w})} = \frac{1}{1 + \sum_{j=1}^{k-1} \exp(\mathbf{x}^T \mathbf{w}_j)}$$

با پارامترهای مدل با پارامتر \mathbf{W} و انتقال softmax، می‌توانیم فرمول حداکثر درست‌نمایی را تعیین کنیم. با وارد کردن پارامتر به معادله (۸،۳)، با گرفتن log منفی آن احتمال و حذف ثابت‌ها، به از خطای زیر برای نمونه‌های $(x_1, y_1), \dots, (x_n, y_n)$

می‌رسیم

$$\min_{\mathbf{W} \in \mathbb{R}^{d \times k}} \sum_{i=1}^n \log(1^T \exp(\mathbf{x}_i^T \mathbf{W})) - \mathbf{x}_i^T \mathbf{W} \mathbf{y}_i$$

با گرادیان

$$\nabla \sum_{i=1}^n (\log(1^T \exp(\mathbf{x}_i^T \mathbf{W})) - \mathbf{x}_i^T \mathbf{W} \mathbf{y}_i) = \sum_{i=1}^n \frac{\exp(\mathbf{x}_i^T \mathbf{W})^T \mathbf{x}_i^T}{1^T \exp(\mathbf{x}_i^T \mathbf{W})} - \mathbf{x}_i \mathbf{y}_i^T$$

مانند قبل برای این گرادیان راه‌حل بسته‌ای نداریم و از روش‌های تکراری برای حل \mathbf{W} استفاده می‌کنیم. توجه داشته باشید که در اینجا برخلاف روش‌های قبلی، روی بخشی از متغیر محدودیت داریم. با این حال، این فقط برای راحتی نوشته شده است. ما $\mathbf{W}_{:k}$ را بهینه نمی‌کنیم، زیرا روی صفر ثابت شده است. می‌توان این کمینه‌سازی و گرادیان را بازنویسی کرد تا فقط برای $\mathbf{W}_{:(1:k-1)}$ اعمال شود. این مربوط به مقداردهی اولیه $\mathbf{W}_{:k} = 0$ ، و سپس تنها با استفاده از اولین ستون‌های $k-1$ گرادیان در به روزرسانی به $\mathbf{W}_{:(1:k-1)}$ است.

پیش‌بینی نهایی $\text{softmax}(\mathbf{x}^T \mathbf{W}) \in [0, 1]$ احتمال حضور در یک کلاس را نشان می‌دهد. همانند رگرسیون لجستیک، برای انتخاب یک کلاس، بالاترین مقدار احتمال انتخاب می‌شود. برای مثال، با $k = 4$ ، ممکن است $[0.1 \ 0.2 \ 0.6 \ 0.1]$ را پیش‌بینی کنیم و بنابراین تصمیم بگیریم که نقطه را در کلاس 3 واحد طبقه‌بندی کنیم.

نکته در مورد همپوشانی کلاس‌ها: اگر می‌خواهید چندین کلاس را برای نقطه داده \mathbf{x} پیش‌بینی کنید، یک استراتژی رایج این است که پیش‌بینی‌های باینری جداگانه برای هر کلاس را یاد بگیرید. هر پیش‌بینی‌کننده به‌طور جداگانه مورد پرسش قرار می‌گیرد، و یک نقطه داده هر کلاس را با 0 یا 1 برچسب‌گذاری می‌کند، که احتمالاً بیش از یک کلاس دارای 1 است. در بالا، موردی را بررسی کردیم که در آن نقطه داده منحصرأ در یکی از کلاس‌های ارائه شده قرار داشت، با تنظیم $n = 1$ در چند جمله‌ای.

فصل ۹

نمایش‌هایی برای یادگیری ماشینی

در ابتدا، ممکن است به نظر برسد که کاربرد رگرسیون خطی و طبقه‌بندی برای مسائل زندگی واقعی بسیار محدود است. و مشخص نیست که آیا درست است (بیشتر اوقات) فرض کنیم که متغیر هدف ترکیبی خطی از ویژگی‌ها است یا خیر. خوشبختانه، کاربرد رگرسیون خطی گسترده‌تر از آن چیزی است که در ابتدا تصور می‌شد. ایده اصلی این است که قبل از مرحله برازش، یک تبدیل غیر خطی به ماتریس داده X اعمال شود، که برازش غیر خطی را ممکن می‌کند. به دست آوردن چنین نمایش ویژگی مفیدی یک مشکل اساسی در یادگیری ماشین است.

ما ابتدا نمایش‌های ثابت را برای رگرسیون خطی بررسی می‌کنیم: شبکه‌های برازش منحنی چند جمله‌ای و تابع پایه شعاعی (RBF). سپس، در مورد بازنمایی‌های یادگیری بحث خواهیم کرد.

۹-۱ شبکه‌های تابع پایه شعاعی^۱ و نمایش هسته

ایده شبکه‌های تابع پایه شعاعی (RBF) یک تعمیم طبیعی از برازش منحنی چند اسمی و رویکردهای بخش قبل است. با توجه به مجموعه داده $D = \{(x_i, y_i)\}_{i=1}^n$ ، با انتخاب p نقطه به عنوان "مراکز" در فضای ورودی X شروع می‌کنیم. ما آن مراکز را به صورت c_1, c_2, \dots, c_p نشان می‌دهیم. معمولاً، اینها را می‌توان از D انتخاب کرد یا با استفاده از برخی تکنیک‌های خوشه‌بندی (مانند الگوریتم EM، میانگین K) محاسبه کرد.

هنگامی که خوشه‌ها با استفاده از مدل مخلوط گاوسی تعیین می‌شوند، توابع پایه را می‌توان با این عنوان انتخاب کرد

$$\varphi_j(\mathbf{x}) = e^{-\frac{1}{2}(\mathbf{x}-\mathbf{c}_j)^T \Sigma_j^{-1}(\mathbf{x}-\mathbf{c}_j)}$$

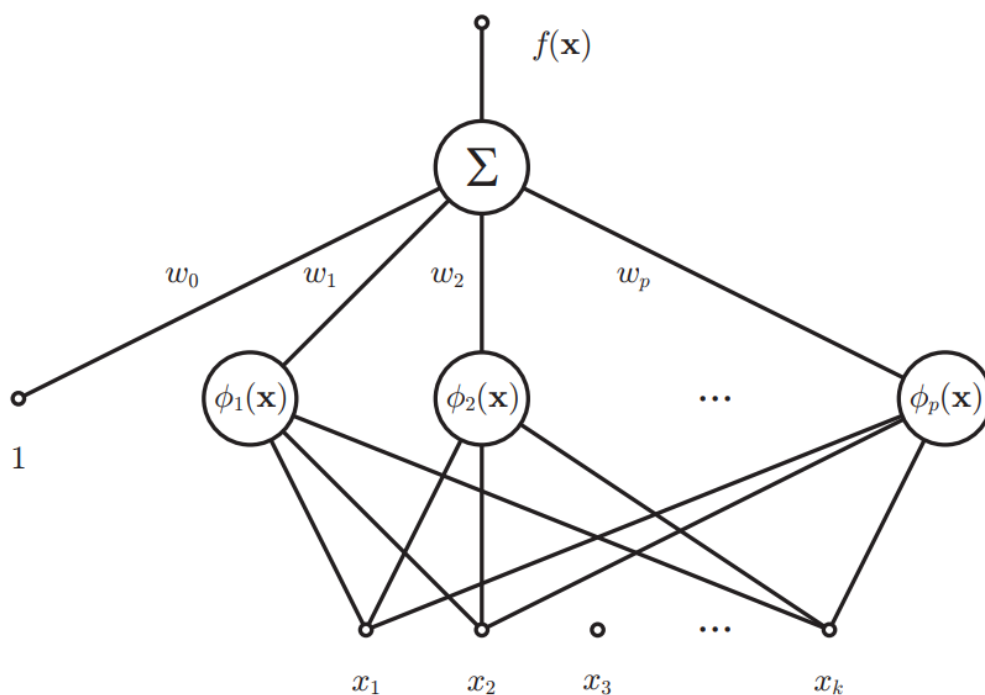
که در آن مراکز خوشه و ماتریس کوواریانس در طول خوشه‌بندی یافت می‌شوند. وقتی از $K - means$ یا خوشه‌بندی دیگر استفاده می‌شود، می‌توانیم از این فرمول استفاده کنیم

$$\varphi_j(\mathbf{x}) = e^{-\frac{\|\mathbf{x}-\mathbf{c}_j\|^2}{2\sigma_j^2}}$$

¹ Radial basis function

جایی که σ_j را می‌توان به طور جداگانه بهینه کرد. به عنوان مثال، با استفاده از مجموعه اعتبارسنجی در زمینه تبدیل‌های چند بعدی از \mathbf{x} به Φ ، توابع پایه را می‌توان به عنوان توابع هسته نیز نام برد، یعنی ماتریس $\varphi_j(\mathbf{x}) = k_j(\mathbf{x}, \mathbf{c}_j)$

$$\Phi = \begin{bmatrix} \varphi_0(\mathbf{x}_1) & \varphi_1(\mathbf{x}_1) & \cdots & \varphi_p(\mathbf{x}_1) \\ \varphi_0(\mathbf{x}_2) & \varphi_1(\mathbf{x}_2) & & \\ \vdots & & \ddots & \\ \varphi_0(\mathbf{x}_n) & & & \varphi_p(\mathbf{x}_n) \end{bmatrix}$$



شکل ۹/۱: شبکه تابع پایه شعاعی.

اکنون به عنوان یک ماتریس داده جدید استفاده می‌شود. برای یک ورودی داده شده \mathbf{x} ، پیش‌بینی هدف y به این صورت محاسبه می‌شود

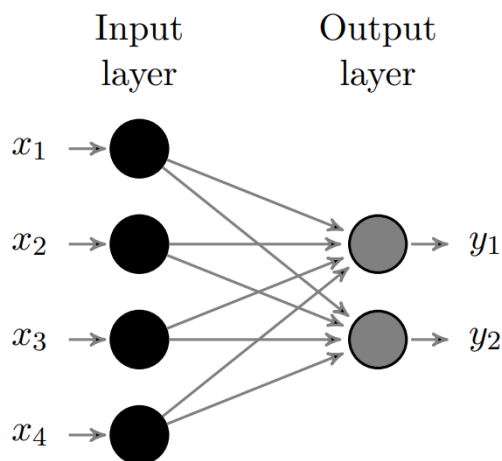
$$\begin{aligned} f(\mathbf{x}) &= w_0 + \sum_{j=1}^p w_j \varphi_j(\mathbf{x}) \\ &= \sum_{j=0}^p w_j \varphi_j(\mathbf{x}) \end{aligned}$$

جایی که $\varphi_0(\mathbf{x}) = 1$ و \mathbf{w} یافت می‌شود. می‌توان ثابت کرد که با تعداد کافی از توابع پایه شعاعی می‌توانیم هر تابع را به طور دقیق تقریب کنیم. همانطور که در شکل ۹،۱ مشاهده می‌شود، می‌توانیم RBFها را به عنوان شبکه‌های عصبی در نظر بگیریم.

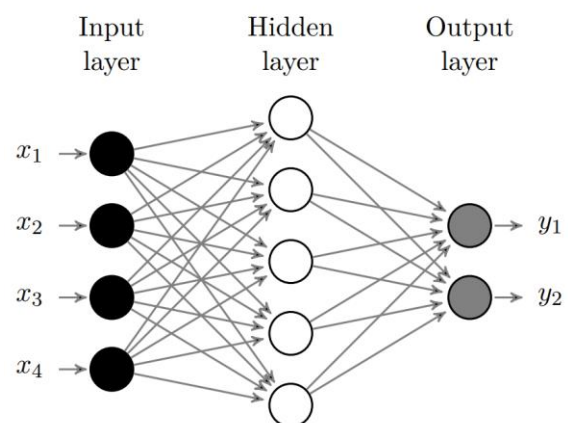
شبکه‌های RBF و نمایش‌های هسته بسیار مرتبط هستند. تمایز اصلی این است که نمایش‌های هسته از هر تابع هسته برای اندازه‌گیری شباهت $k(x, c_j) = \varphi_j(x)$ استفاده می‌کنند، که در آن توابع پایه شعاعی نمونه‌ای از یک هسته هستند. علاوه بر این، اگر یک هسته RBF انتخاب شود، برای بازنمایی هسته معمولی، مراکز از مجموعه داده آموزشی انتخاب می‌شوند. برای شبکه‌های RBF، انتخاب مراکز به طور کلی به عنوان یک مرحله مهم باقی مانده است، جایی که می‌توان آنها را از مجموعه آموزشی انتخاب کرد، اما همچنین می‌تواند به روش‌های دیگر انتخاب شود.

۹-۲ بازنمایی‌های یادگیری

رویکردهای زیادی برای بازنمایی یادگیری وجود دارد. دو رویکرد غالب، تکنیک‌های فاکتورسازی ماتریس (نیمه نظارت‌شده) و شبکه‌های عصبی هستند. شبکه‌های عصبی بر اساس مدل‌های خطی تعمیم‌یافته‌ای که مورد بحث قرار گرفتیم، ساخته می‌شوند و چندین مدل خطی تعمیم‌یافته را با هم ترکیب می‌کنند. تکنیک‌های فاکتورسازی ماتریسی (به عنوان مثال، کاهش ابعاد، کدگذاری پراکنده) معمولاً داده‌های ورودی را در یک فرهنگ لغت و یک نمایش جدید (پایه) فاکتور می‌کنند. ما ابتدا شبکه‌های عصبی را مورد بحث قرار می‌دهیم و سپس در مورد بسیاری از تکنیک‌های یادگیری بدون نظارت و نیمه نظارت که توسط فاکتورسازی‌های ماتریسی احاطه شده‌اند بحث خواهیم کرد.



شکل ۹/۲: مدل خطی تعمیم یافته، مانند رگرسیون لجستیک.



شکل ۹/۳: شبکه عصبی دو لایه استاندارد

۹-۲-۱ شبکه‌های عصبی

شبکه‌های عصبی شکلی از یادگیری بازنمایی نظارت شده هستند. مانند قبل، هدف یادگیری تابعی از ورودی‌ها، f ، برای تولید یک پیش‌بینی از هدف است: $f(x)$. افزودن لایه‌های پنهان با توابع فعال‌سازی غیرخطی، یادگیری توابع غیرخطی f را امکان‌پذیر می‌سازد. برای برخی از شهود، می‌توان در نظر گرفت که اولین لایه‌های پنهان لایه بازنمایی را تشکیل می‌دهند، با یادگیری در آخرین لایه مربوط به قسمت پیش‌بینی نظارت شده است. شکل ۹،۲ مدل گرافیکی مدل‌های خطی تعمیم‌یافته‌ای را که در فصل‌های قبل مورد بحث قرار دادیم، نشان می‌دهد. جایی که وزن‌ها و انتقال مربوطه را می‌توان روی فلش‌ها در نظر گرفت (زیرا آنها متغیرهای تصادفی نیستند). شکل ۹،۳ یک شبکه عصبی با یک لایه پنهان را نشان می‌دهد. این شبکه عصبی دو لایه نامیده می‌شود، زیرا دو لایه وزن دارد

در شکل، شبکه عصبی یک بردار ویژگی ۴ بعدی $x = [x_1, x_2, x_3, x_4]$ (یعنی $d = 4$) را وارد کرده و یک پیش‌بینی دو بعدی $y = [y_1, y_2]$ (یعنی $m = 2$)، لایه پنهان شامل یک نقشه برداری از x به یک نمایش جدید است که ۵ بعدی است (یعنی $k_1 = 5$ طبق نماد زیر). برای شبکه عصبی، اجازه دهید هر گره در این نمایش پنهان با $k \in \{1, \dots, 5\}$ ، نمایه شود. هر h_k از تبدیل وزن خطی x تشکیل شده است، مانند انتقال سیگموئید: $\{h_k = \sigma(\sum_{j=1}^d (x_j w_{kj})) = \sigma(xw_k)\}$ که در آن $w_k \in \mathbb{R}^d$ وزن‌های روی اولین لایه است که برای تولید گره در نمایش پنهان شماره k استفاده می‌شود.

مثال ۱۸: برای مثال ساده، $d = 1$ (یعنی یک مشاهده ورودی)، $m = 1$ (یعنی یک خروجی)، $k_1 = 2$ (یعنی لایه پنهان ۲ بعدی) و یک انتقال سیگموئید را در نظر بگیرید تا اولین مورد را بدست آورید. لایه پنهان فرض کنید یک نمونه به ما داده شده است (x, y) . سپس مشاهده ورودی x تبدیل به

$$h = [h_1, h_2], \quad \text{with } h_1 = \sigma(xw_{1(2)}) \text{ and } h_2 = \sigma(xw_{2(2)}) \quad \text{for } w_{1(2)}, w_{2(2)} \in \mathbb{R}.$$

برای جلوگیری از انتقال نماد، از $x \in \mathbb{R}^{1 \times d}$ برای ارائه یک ردیف از ماتریس داده $x \in \mathbb{R}^{n \times d}$ و بردار ردیف $h \in \mathbb{R}^{1 \times k_1}$ استفاده کردیم. برای تمایز بین وزن‌های لایه اول و آخر از نماد بالانویس استفاده می‌کنیم. ممکن است غیرقابل درک به نظر برسد که چرا ما برچسب $w^{(2)}$ را برای لایه ورودی و $w^{(1)}$ برای لایه خروجی می‌زنیم، اما در زیر خواهید دید که نشان‌گذاری شروع نمایه‌سازی از لایه خروجی را ساده‌تر می‌کند.

هنگامی که h را داریم، می‌توانیم وانمود کنیم که h نمایش ورودی جدید است و ادامه می‌دهیم و یک مدل خطی (تعمیم‌یافته) در این لایه آخر یاد می‌گیریم. بیایید دو حالت را در نظر بگیریم: $y \in \mathbb{R}$ و $y \in \{0, 1\}$. اگر $y \in \mathbb{R}$ ، از رگرسیون خطی برای آخرین لایه استفاده می‌کنیم و بنابراین وزن‌های $w^{(2)} \in \mathbb{R}^2$ را طوری یاد می‌گیریم که $hw^{(2)}$ خروجی واقعی y را تقریب می‌زند. اگر $y \in \{0, 1\}$ ، از رگرسیون لجستیک برای آخرین لایه استفاده می‌کنیم و بنابراین وزن‌های $w^{(2)} \in \mathbb{R}^2$ را طوری یاد می‌گیریم که $\sigma(hw^{(2)})$ خروجی واقعی y را تقریب می‌زند.

حال حالت کلی‌تر را با هر m, k_1, d در نظر می‌گیریم. برای ارائه شهودی برای این تنظیم کلی تر، ما با یک لایه پنهان، برای تابع انتقال سیگموئید و از دست دادن خروجی آنتروپی متقابل شروع می‌کنیم. برای رگرسیون لجستیک $W \in \mathbb{R}^{d \times m}$ را با $y \in \mathbb{R}^m$ را پیش‌بینی می‌کنیم، زیرا تعمیم‌های بعدی را واضح‌تر می‌کند و نماد وزن‌های هر لایه را یکنواخت‌تر می‌کند. وقتی یک لایه پنهان اضافه می‌کنیم، دو ماتریس پارامتر $W^{(1)} \in \mathbb{R}^{k_1 \times m}$ و $W^{(2)} \in \mathbb{R}^{d \times k_1}$ داریم که k_1 بعد لایه پنهان است.

$$h = \sigma(xW^{(2)}) = \begin{bmatrix} \sigma(xw_{:1}^{(2)}) \\ \sigma(xw_{:2}^{(2)}) \\ \vdots \\ \sigma(xw_{:k_1}^{(2)}) \end{bmatrix} \in \mathbb{R}^{k_1}$$

که در آن تابع سیگموئید برای هر ورودی در $xW^{(2)}$ و اعمال می‌شود. این لایه پنهان مجموعه جدیدی از ویژگی‌ها است و مجدداً بهینه‌سازی رگرسیون لجستیک معمولی را برای یادگیری وزن‌ها در h انجام خواهید داد:

$$p(y = 1|x) = \sigma(hW^{(1)}) = \sigma(\sigma(xW^{(2)})W^{(1)})$$

با مدل احتمالی و پارامتر مشخص شده، اکنون باید یک الگوریتم برای بدست آوردن آن پارامترها استخراج کنیم. مانند قبل، ما یک رویکرد ماکسیمم احتمال را اتخاذ کرده و به‌روزرسانی‌های شیب نزولی را استخراج می‌کنیم. به نظر می‌رسد این ترکیب

انتقال‌ها مسائل را پیچیده می‌کند، اما همچنان می‌توانیم شیب $w.r.t$ را بگیریم. ما اکنون پارامترهای بیشتری داریم: $W^{(2)} \in \mathbb{R}^{1 \times k_1} \mathbb{R}^{k_1 \times d}$ ، وقتی گرادیان $w.r.t$ را داریم. هر ماتریس پارامتر، طبق معمول به سادگی یک قدم در جهت منفی گرادیان برداریم. گرادیان برای این پارامترها اطلاعات را به اشتراک می‌گذارند. برای کارایی محاسباتی، گرادیان ابتدا برای $W^{(1)}$ محاسبه می‌شود، و اطلاعات گرادیان تکراری برای محاسبه گرادیان برای $W^{(2)}$ ارسال می‌شود. این الگوریتم معمولاً به نام انتشار برگشتی نامیده می‌شود که در ادامه به توضیح آن می‌پردازیم.

به طور کلی، ما می‌توانیم گرادیان را برای هر تعداد لایه پنهان محاسبه کنیم. هر تابع انتقال قابل تفکیک f_1, \dots, f_H را، نشان دهید. با f_1 به عنوان انتقال خروجی، و k_1, \dots, k_{H-1} مرتب شده است. به عنوان ابعاد پنهان با لایه‌های پنهان $H - 1$. سپس خروجی از شبکه عصبی است

$$f_1(f_2(\dots f_{H-1}(f_H(xW^{(H)}))W^{(H-1)}))\dots)W^{(1)})$$

$$W^{(1)} \in \mathbb{R}^{k_1 \times m}, W^{(2)} \in \mathbb{R}^{k_2 \times k_1}, \dots, W^{(H)} \in \mathbb{R}^{d \times k_{H-1}}$$
 که

الگوریتم پس انتشار

ما با استخراج پس انتشار برای دو لایه شروع خواهیم کرد. گسترش چندین لایه با توجه به این اشتقاق واضح‌تر خواهد بود. با توجه به اندازه شبکه، ما اغلب با گرادیان کاهشی تصادفی یاد خواهیم گرفت. بنابراین، ابتدا این گرادیان را با فرض اینکه فقط یک نمونه (x, y) داریم محاسبه می‌کنیم.

الگوریتم انتشار برگشتی به سادگی گرادیان کاهشی بر روی یک هدف غیر محدب، با ترتیب دقیق محاسبات برای جلوگیری از تکرار محاسبات است. به ویژه، ابتدا به جلو منتشر می‌شود و متغیر $h = f_2(xW^{(2)}) \in \mathbb{R}^{1 \times k}$ و سپس $\hat{y} = f_1(hw^{(1)}) = f_1(f_2(xW^{(2)})W^{(1)})$ را محاسبه می‌کند. سپس خطای بین پیش‌بینی \hat{y} و برچسب واقعی را محاسبه می‌کنیم. گرادیان این خطا (از دست دادن) را $w.r.t$ می‌گیریم. به پارامترهای ما؛ در این حالت، برای محاسبه کارآمد، بهترین ترتیب، محاسبه گرادیان $w.r.t$ است. ابتدا به آخرین پارامتر $W^{(1)}$ و سپس $W^{(2)}$ برسید. دلیل استفاده از اصطلاح *back-propagation* به همین دلیل است، زیرا خطا ابتدا از آخرین لایه به عقب منتشر می‌شود.

سپس انتخاب‌ها شامل انتخاب انتقال‌ها در هر لایه، تعداد گره‌های پنهان و از دست دادن لایه آخر است. تطبیق محدب $L(\cdot, y)$ به $p(y|x)$ انتخاب شده و تابع انتقال متناظر برای آخرین لایه شبکه عصبی بستگی دارد، درست مانند مدل‌های خطی تعمیم یافته. برای سهولت علامت‌گذاری، این تابع خطا را به صورت تعریف می‌کنیم

$$Err(W^{(1)}, W^{(2)}) \stackrel{\text{def}}{=} \sum_{k=1}^m L(f_1(f_2(xW^{(2)})W_{:,k}^{(1)}), y_k)$$

برای یک نمونه (x, y) . به عنوان مثال، برای $p(y = 1|x)$ گاوسی و انتقال هویت¹ f_2 ، $Err(W^{(1)}, W^{(2)}) = (f_2(xW^{(2)})W^{(1)} - y)^2$ را بدست می‌آوریم. اگر $p(y = 1|x)$ یک توزیع برنولی باشد، آنگاه از دست دادن رگرسیون لجستیک (آنتروپی متقاطع) را انتخاب می‌کنیم.

مانند قبل، گرادیان‌های اتلاف $w.r.t$ را محاسبه می‌کنیم. پارامترهای ما ابتدا مشتق جزئی $w.r.t$ را می‌گیریم. پارامترهای $W^{(1)}$ (با فرض اینکه $W^{(2)}$ ثابت باشد).

¹ identity transfer

$$\begin{aligned}\frac{\partial \text{Err}(\mathbf{W}^{(1)}, \mathbf{W}^{(2)})}{\partial \mathbf{W}_{jk}^{(1)}} &= \frac{\partial L(f_1(f_2(\mathbf{x}\mathbf{W}^{(2)})\mathbf{W}^{(1)}), \mathbf{y})}{\partial \mathbf{W}_{jk}^{(1)}} \quad \triangleright \hat{y}_k = f_1(\mathbf{h}\mathbf{W}_{:k}^{(1)}) \\ &= \left(\frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \right) \frac{\partial \hat{\mathbf{y}}_k}{\partial \mathbf{W}_{jk}^{(1)}}\end{aligned}$$

که در آن فقط \hat{y}_k تحت تأثیر $\mathbf{W}_{jk}^{(1)}$ در از دست دادن قرار می‌گیرد، و بنابراین گرادیان برای بقیه صفر است. در ادامه،

$$\begin{aligned}\frac{\partial \text{Err}(\mathbf{W}^{(1)}, \mathbf{W}^{(2)})}{\partial \mathbf{W}_{jk}^{(1)}} &= \left(\frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \right) \frac{\partial f_1(\theta_k^{(1)})}{\partial \theta_k^{(1)}} \frac{\partial \theta_k^{(1)}}{\partial \mathbf{W}_{jk}^{(1)}} \quad \triangleright \theta_k^{(1)} = \mathbf{h}\mathbf{W}_{:k}^{(1)} \\ &= \left(\frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \right) \frac{\partial f_1(\theta_k^{(1)})}{\partial \theta_k^{(1)}} \mathbf{h}_j\end{aligned}$$

در این مرحله این معادلات انتزاعی هستند. اما محاسبه آنها برای هزینه و انتقالاتی که ما بررسی کردیم ساده است. به عنوان

مثال، برای $L(\hat{y}_k, y_k) = \frac{1}{2}(\hat{y}_k - y_k)^2$ و f_2 هویت، به دست می‌آوریم

$$\begin{aligned}\frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} &= (\hat{\mathbf{y}}_k - \mathbf{y}_k) \\ \frac{\partial f_1(\theta_k^{(1)})}{\partial \theta_k^{(1)}} &= 1\end{aligned}$$

که می‌دهد

$$\frac{\partial \text{Err}(\mathbf{W}^{(1)}, \mathbf{W}^{(2)})}{\partial \mathbf{W}_{jk}^{(1)}} = \left(\frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \right) \frac{\partial f_1(\theta_k^{(1)})}{\partial \theta_k^{(1)}} \mathbf{h}_j = (\hat{\mathbf{y}}_k - \mathbf{y}_k) \mathbf{h}_j$$

به روز رسانی گرادیان طبق معمول با $\mathbf{W}^{(1)} = \mathbf{W}^{(1)} - \alpha(\hat{\mathbf{y}} - \mathbf{y})\mathbf{h}^T$ برای مقداری α اندازه گام است.

سپس، گرادیان جزئی را با توجه به $\mathbf{W}^{(2)}$ محاسبه می‌کنیم. اما اکنون، کل متغیر خروجی $\mathbf{y} \in \mathbb{R}^{1 \times m}$ تحت تأثیر انتخاب $\mathbf{W}_{ij}^{(2)}$ برای همه $i \in 1, \dots, k_2$ و $j \in 1, \dots, k_1$ بنابرین، باید مشتق جزئی w.r.t را برای هر \mathbf{y} بگیریم.

$$\begin{aligned}\frac{\partial \text{Err}(\mathbf{W}^{(1)}, \mathbf{W}^{(2)})}{\partial \mathbf{W}_{ij}^{(2)}} &= \frac{\partial \sum_{k=1}^m L(f_1(f_2(\mathbf{x}\mathbf{W}^{(2)})\mathbf{W}_{:k}^{(1)}), \mathbf{y}_k)}{\partial \mathbf{W}_{ij}^{(2)}} \\ \sum_{k=1}^m \frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \frac{\partial \hat{\mathbf{y}}_k}{\partial \mathbf{W}_{ij}^{(2)}} &\quad \triangleright \hat{y}_k = f_1(\mathbf{h}\mathbf{W}_{:k}^{(1)}) = f_1(\theta_k^{(1)}) \\ &= \sum_{k=1}^m \frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \frac{\partial f_1(\theta_k^{(1)})}{\partial \theta_k^{(1)}} \frac{\partial \theta_k^{(1)}}{\partial \mathbf{W}_{ij}^{(2)}}\end{aligned}$$

در ادامه

$$\frac{\partial \theta_k^{(1)}}{\partial \mathbf{W}_{ij}^{(2)}} = \frac{\partial \mathbf{h} \mathbf{W}_{:k}^{(1)}}{\partial \mathbf{W}_{ij}^{(2)}} = \frac{\partial \sum_{l=1}^k \mathbf{h}_l \mathbf{W}_{lk}^{(1)}}{\partial \mathbf{W}^{(2)}}$$

$$\frac{\partial \sum_{l=1}^k f_2(\mathbf{x} \mathbf{W}_{:l}^{(2)}) \mathbf{W}_{lk}^{(1)}}{\partial \mathbf{W}_{ij}^{(2)}}$$

$$\sum_{l=1}^k \mathbf{W}_{lk}^{(1)} \frac{\partial f_2(\mathbf{x} \mathbf{W}_{:l}^{(2)})}{\partial \mathbf{W}_{ij}^{(2)}} \\ = \mathbf{W}_{jk}^{(1)} \frac{\partial f_2(\mathbf{x} \mathbf{W}_{:j}^{(2)})}{\partial \mathbf{W}_{ij}^{(2)}}$$

زیرا $0 = \frac{\partial f_2(\mathbf{x} \mathbf{W}_{:l}^{(2)})}{\partial \mathbf{W}_{ij}^{(2)}}$ برای $j \neq l$. حالا قانون زنجیره‌ای را ادامه دهید

$$\frac{\partial f_2(\mathbf{x} \mathbf{W}_{:j}^{(2)})}{\partial \mathbf{W}_{ij}^{(2)}} = \frac{\partial f_2(\theta_j^{(2)})}{\partial \theta_j^{(2)}} \frac{\partial \theta_j^{(2)}}{\partial \mathbf{W}_{ij}^{(2)}} \quad \triangleright \theta_j^{(2)} = \mathbf{x} \mathbf{W}_{:j}^{(2)} \\ = \frac{\partial f_2(\theta_j^{(2)})}{\partial \theta_j^{(2)}} \mathbf{x}_i$$

با کنار هم قرار دادن اینها، می‌گیریم

$$\frac{\partial \text{Err}(\mathbf{W}^{(1)}, \mathbf{W}^{(2)})}{\partial \mathbf{W}_{ij}^{(2)}} = \sum_{k=1}^m \frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \frac{\partial f_1(\theta_k^{(1)})}{\partial \theta_k^{(1)}} \frac{\partial \theta_k^{(1)}}{\partial \mathbf{W}_{ij}^{(2)}}$$

توجه داشته باشید که مقداری از گرادیان مانند $\mathbf{W}^{(1)}$ است، یعنی.

$$\delta_k^{(1)} = \frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \frac{\partial f_1(\theta_k^{(1)})}{\partial \theta_k^{(1)}}$$

محاسبه این مؤلفه‌ها فقط باید یک بار برای $\mathbf{W}^{(1)}$ انجام شود و این اطلاعات دوباره منتشر شود تا گرادیان $\mathbf{W}^{(2)}$ بدست آید.

تفاوت در گرادیان $\frac{\partial \theta^{(1)}}{\partial \mathbf{W}^{(2)}}$ است، زیرا \mathbf{h} به $\mathbf{W}^{(2)}$ متکی است. برای $\mathbf{W}^{(1)}$ ، $\mathbf{h} = f_2(\mathbf{x}_i \mathbf{W}^{(2)})$ یک ثابت است، و بنابراین بر شیب $\mathbf{W}^{(1)}$ تأثیر نمی‌گذارد. گرادیان نهایی است

$$\frac{\partial \text{Err}(\mathbf{W}^{(1)}, \mathbf{W}^{(2)})}{\partial \mathbf{W}_{ij}^{(2)}} = \left(\sum_{k=1}^m \delta_k^{(1)} \mathbf{W}_{jk}^{(1)} \right) \frac{\partial f_2(\theta_j^{(2)})}{\partial \theta_j^{(2)}} \mathbf{x}_i \\ = (\mathbf{W}_{j:}^{(1)} \delta^{(1)}) \frac{\partial f_2(\theta_j^{(2)})}{\partial \theta_j^{(2)}} \mathbf{x}_i$$

اگر لایه دیگری قبل از $\mathbf{W}^{(2)}$ اضافه شود، اطلاعات منتشر شده به عقب است

$$\delta_j^{(2)} = \left(\mathbf{w}_{j:}^{(1)} \delta^{(1)} \right) \frac{\partial f_2 \left(\theta_j^{(2)} \right)}{\partial \theta_j^{(2)}}$$

و x_i با $h_i^{(2)}$ جایگزین می‌شود. گرادیان برای $W_{ij}^{(3)}$ است

$$\left(\mathbf{w}_{j:}^{(2)} \delta^{(2)} \right) \frac{\partial f_3 \left(\theta_j^{(3)} \right)}{\partial \theta_j^{(3)}} \mathbf{x}_i$$

مثال ۱۹: فرض کنید $p(y = 1|x)$ یک توزیع برنولی باشد، با f_1 و f_2 هر دو تابع سیگموئید. هزینه آنتروپی متقاطع است. می‌توانیم قانون به‌روزرسانی دو لایه را با این تنظیمات با وصل کردن در بالا استخراج کنیم.

$$L(\hat{y}, y) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y}) \quad \triangleright \text{cross - entropy}$$

$$\frac{\partial L(\hat{y}, y)}{\partial \hat{y}} = -\frac{y}{\hat{y}} + \frac{1 - y}{1 - \hat{y}}$$

$$f_2 \left(\mathbf{x} \mathbf{W}_{:j}^{(2)} \right) = \sigma \left(\mathbf{x} \mathbf{W}_{:j}^{(2)} \right) = \frac{1}{1 + \exp \left(-\mathbf{x} \mathbf{W}_{:j}^{(2)} \right)}$$

$$f_1 \left(\mathbf{h} \mathbf{W}_{:k}^{(1)} \right) = \sigma \left(\mathbf{h} \mathbf{W}_{:k}^{(1)} \right) = \frac{1}{1 + \exp \left(-\mathbf{h} \mathbf{W}_{:k}^{(1)} \right)}$$

$$\partial \sigma(\theta) = \sigma(\theta)(1 - \sigma(\theta))$$

اکنون می‌توانیم به‌روزرسانی پس انتشار را با انتشار به جلو محاسبه کنیم

$$\mathbf{h} = \sigma(\mathbf{x} \mathbf{W}^{(2)})$$

$$\hat{y} = \sigma(\mathbf{h} \mathbf{W}^{(1)})$$

$$\begin{array}{c} d \\ \boxed{\mathbf{X}} \\ n \end{array} \approx \begin{array}{c} k \\ \boxed{\mathbf{H}} \\ n \end{array} \times \begin{array}{c} d \\ \boxed{\mathbf{D}} \\ k \end{array}$$

شکل ۹/۴: فاکتورسازی ماتریسی ماتریس داده $X \in \mathbb{R}^{n \times d}$.

و سپس انتشار گرادیان به عقب

$$\delta_k^{(1)} = \frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \frac{\partial f_1 \left(\theta_k^{(1)} \right)}{\partial \theta_k^{(1)}}$$

$$\begin{aligned}
&= \left(-\frac{y_k}{\hat{y}_k} + \frac{1 - y_k}{1 - \hat{y}_k} \right) \hat{y}_k (1 - \hat{y}_k) \\
&= -y_k (1 - \hat{y}_k) + (1 - y_k) \hat{y}_k \\
&= \hat{y}_k - y_k \\
&\frac{\partial}{\partial \mathbf{W}_{jk}^{(1)}} = \delta_k^{(1)} \mathbf{h}_j \\
&\delta_j^{(2)} = \left(\mathbf{W}_{j:}^{(1)} \delta^{(1)} \right) \mathbf{h}_j (1 - \mathbf{h}_j) \\
&\frac{\partial}{\partial \mathbf{W}_{ij}^{(2)}} = \delta_j^{(2)} \mathbf{x}_i
\end{aligned}$$

به روزرسانی به سادگی شامل گام برداشتن در جهت این شیب‌ها است، همانطور که برای گرادیان کاهشی معمول است. ما با برخی از $\mathbf{W}^{(1)}$ و $\mathbf{W}^{(2)}$ اولیه شروع می‌کنیم (مثلاً با مقادیر تصادفی پر شده است)، و سپس قوانین گرادیان کاهشی را با این گرادیان‌ها اعمال می‌کنیم.

۹-۲-۲ یادگیری بدون نظارت^۱ و فاکتورسازی ماتریس^۲

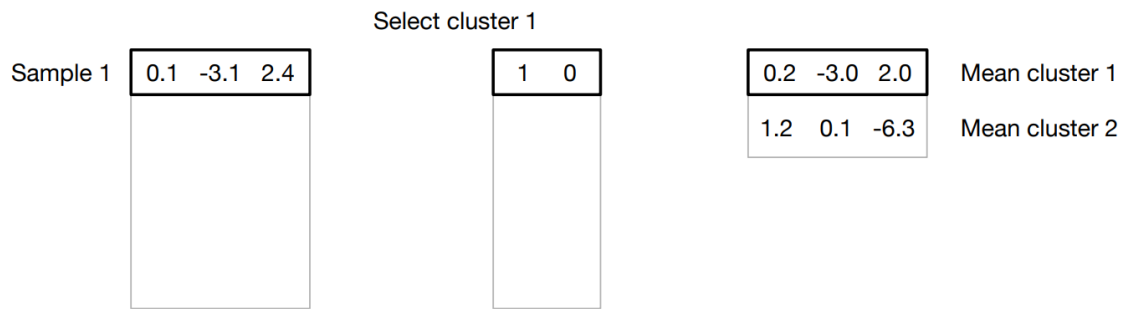
استراتژی دیگر برای به دست آوردن یک نمایش جدید از طریق فاکتورسازی ماتریس است. ماتریس داده \mathbf{X} به یک فرهنگ لغت \mathbf{D} و یک پایه یا نمایش جدید \mathbf{H} تبدیل می‌شود (شکل ۹،۴ را ببینید). در واقع، بسیاری از الگوریتم‌های یادگیری بدون نظارت (مانند کاهش ابعاد، کدگذاری پراکنده) و الگوریتم‌های یادگیری نیمه‌نظارت‌شده (مانند یادگیری فرهنگ لغت نظارت‌شده) در واقع می‌توانند به عنوان عامل‌بندی ماتریسی فرموله شوند. ما به عنوان مثال به خوشه‌بندی k - means و تجزیه و تحلیل اجزای اصلی نگاه خواهیم کرد. الگوریتم‌های باقی مانده به سادگی در جدول زیر خلاصه شده‌اند. این رویکرد کلی برای به دست آوردن یک نمایش جدید با استفاده از فاکتورسازی، یادگیری فرهنگ لغت^۳ نامیده می‌شود.

خوشه‌بندی k - means یک مشکل یادگیری بدون نظارت برای گروه‌بندی نقاط داده در k خوشه‌ها با به مینیمم رساندن فاصله تا میانگین هر خوشه است. این مشکل معمولاً به عنوان یک رویکرد یادگیری بازنمایی در نظر گرفته نمی‌شود، زیرا شماره خوشه معمولاً به عنوان یک نمایش استفاده نمی‌شود. با این حال، با این حال، ما با k - means شروع می‌کنیم زیرا این یک مثال شهودی از این است که چگونه این الگوریتم‌های یادگیری بدون نظارت را می‌توان به عنوان ماتریس در نظر گرفت.

¹ Unsupervised learning

² matrix factorization

³ dictionary learning



شکل ۹/۵: K - means خوشه‌بندی به عنوان فاکتورسازی ماتریس برای ماتریس داده $X \in \mathbb{R}^{n \times d}$.

فاکتورسازی علاوه بر این، رویکرد خوشه‌بندی را می‌توان به عنوان یک رویکرد یادگیری بازنمایی در نظر گرفت، زیرا این یک گسسته‌سازی آموخته‌شده از فضا است. ما این دیدگاه k - means را پس از بحث در مورد آن به عنوان فاکتورسازی ماتریسی مورد بحث قرار خواهیم داد.

تصور کنید که شما دو خوشه دارید ($k = 2$)، با ابعاد داده $\mathbf{d} = 3$. فرض کنید \mathbf{d}_1 میانگین برای خوشه 1 و \mathbf{d}_2 میانگین برای خوشه 2 باشد. هدف این است که فاصله ℓ_2 هر نقطه داده \mathbf{x} را به مینیمم برسانید. به مرکز خوشه آن

$$\left\| \mathbf{x} - \sum_{i=1}^2 1(\mathbf{x} \text{ in cluster } i) \mathbf{d}_i \right\|_2^2 = \|\mathbf{x} - \mathbf{hD}\|_2^2$$

بردارهای خوشه‌های مختلف \mathbf{h} برای هر \mathbf{x} آموخته می‌شوند، اما فرهنگ لغت میانگین‌ها بین تمام نقاط داده مشترک است. بهینه‌سازی مشخص‌شده باید فرهنگ لغت \mathbf{D} را انتخاب کند که کمترین فاصله را تا نقاط در مجموعه داده آموزشی فراهم می‌کند.

که در آن $\mathbf{h} = [1 \ 0]$ یا $\mathbf{h} = [0 \ 1]$ و $\mathbf{D} = [\mathbf{d}_1; \mathbf{d}_2]$. یک مثال در شکل ۹،۵ نشان داده شده است. برای یک نقطه $\mathbf{x} = [0.1 \ -3.1 \ 2.4]$ ، $\mathbf{h} = [1 \ 0]$ به این معنی که در خوشه 1 با میانگین $\mathbf{d}_1 = [0.2 \ -3.0 \ 2.0]$ قرار می‌گیرد. قرار دادن \mathbf{x} در خوشه 2 که میانگین آن متفاوت‌تر است، خطای بیشتری دارد: $\mathbf{d}_2 = [1.2 \ 0.1 \ -6.3]$.

$$\min_{\substack{\mathbf{H} \in \{0,1\}^{n \times k}, \mathbf{1H}=\mathbf{1} \\ \mathbf{D} \in \mathbb{R}^{k \times d}}} \|\mathbf{X} - \mathbf{HD}\|_F^2$$

تجزیه و تحلیل مولفه‌های اصلی^۱ (PCA) یک تکنیک کاهش ابعاد استاندارد است که در آن داده‌های ورودی $\mathbf{x} \in \mathbb{R}^{1 \times d}$ به ابعاد $\mathbf{h} \in \mathbb{R}^{1 \times k}$ با فضای مولفه‌های اصلی پیش‌بینی می‌شود. این مولفه‌های اصلی جهت ماکسیمم واریانس در داده‌ها هستند. برای به دست آوردن این k مولفه‌های اصلی $\mathbf{D} \in \mathbb{R}^{k \times d}$ ، روش حل رایج به دست آوردن تجزیه ارزش منفرد ماتریس داده $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \in \mathbb{R}^{n \times d}$ است که

$$\mathbf{D} = \mathbf{V}_k^T \in \mathbb{R}^{k \times d}$$

$$\mathbf{H} = \mathbf{U}_k \mathbf{\Sigma}_k \in \mathbb{R}^{n \times k}$$

¹ Principal components analysis

جایی که $\Sigma_k \in \mathbb{R}^{k \times k}$ از بزرگترین مقادیر k مفرد (به ترتیب نزولی) تشکیل شده است و $U_k \in \mathbb{R}^{n \times k}$ و $V_k \in \mathbb{R}^{k \times d}$ بردارهای مفرد متناظر هستند، یعنی $U_k = U_{:,1:k}$ و $V_k = V_{:,1:k}$. نمایش جدید برای X (با استفاده از PCA) این H است. توجه داشته باشید که PCA ویژگی‌های زیر را انتخاب نمی‌کند، بلکه ویژگی‌های جدیدی ایجاد می‌کند: h تولید شده زیرمجموعه‌ای از \mathcal{X} اصلی نیست.

این تکنیک کاهش ابعاد را می‌توان به عنوان یک فاکتورسازی ماتریسی نیز فرموله کرد. بهینه‌سازی مربوطه نشان داده شده است

$$\min_{D \in \mathbb{R}^{k \times d}, H \in \mathbb{R}^{n \times k}} \|X - HD\|_F^2$$

یک راه ساده برای فهمیدن چرایی این امر، یادآوری قضیه معروف *Eckart – Young Mirsky* است که ماتریس K رتبه ای \hat{X} که بهترین تقریب X را بر حسب نرم فروبنیوس مینیمال^۱ دارد، $\hat{X} = U_k \Sigma_k V_k^T$ است.

مانند خوشه‌بندی k - means، ممکن است به سختی بفهمیم که چرا h تولید شده توسط PCA می‌تواند به عنوان یک نمایش مفید باشد. در واقع، PCA اغلب برای تجسم استفاده می‌شود، و بنابراین همیشه برای یادگیری بازنمایی استفاده نمی‌شود. برای تجسم، طرح ریزی اغلب به دو یا سه بعدی تهاجمی است. با این حال، به طور کلی، طرح ریزی به ابعاد پایین این خاصیت را دارد که نویز را حذف می‌کند و فقط معنی دارترین جهت‌ها را حفظ می‌کند. بنابراین، این پیش‌بینی با کاهش تعداد ویژگی‌ها و ترویج تعمیم، با جلوگیری از تطابق بیش از حد با نویز به سرعت یادگیری کمک می‌کند.

کدگذاری پراکنده رویکرد متفاوتی دارد، جایی که داده‌های ورودی به یک نمایش پراکنده گسترش می‌یابد. کدگذاری پراکنده از نظر بیولوژیکی [۱۵] بر اساس فعالیت‌های پراکنده برای حافظه در مغز پستانداران است. تفسیر دیگر این است که کدگذاری پراکنده به طور موثر فضا را گسسته می‌کند، مانند خوشه‌بندی k - means، اما با خوشه‌های همپوشانی و مقدار مربوط به اینکه چقدر یک نقطه به آن خوشه تعلق دارد.

یک استراتژی معمول برای به دست آوردن نمایش‌های پراکنده استفاده از تنظیم کننده پراکنده در نمایش آموخته شده h است. این مربوط به بهینه‌سازی است

$$\min_{D \in \mathbb{R}^{k \times d}, H \in \mathbb{R}^{n \times k}} \|X - HD\|_F^2 + \lambda \sum_{i=1}^k \|H_{:,i}\|_1 + \lambda \sum_{i=1}^k \|D_{:,i}\|_2^2$$

همانطور که در بخش ۵.۴.۲ مورد بحث قرار گرفت، تنظیم کننده ℓ_1 ورودی‌های صفر شده را گسترش می‌دهد، و بنابراین H را با ماکسیمم صفرهای ممکن ترجیح می‌دهد. یک تنظیم کننده نیز به D اضافه می‌شود تا اطمینان حاصل شود که D خیلی بزرگ نمی‌شود. در غیر این صورت، تمام وزن در DH به D منتقل می‌شود.

به طور کلی، انواع مختلفی از الگوریتم‌های یادگیری بدون نظارت وجود دارد که در واقع با فاکتورسازی ماتریس داده‌ها مطابقت دارند. جزئیات بیشتر در ضمیمه D داده شده است. ما همچنین جزئیاتی را در مورد نحوه یادگیری این فاکتورگیری‌ها در پیوست ارائه می‌دهیم. همانند الگوریتم‌های قبلی، آنها به سادگی بر روی متغیرهای (ماتریس) نزولی هستند. تنها تمایز در اینجا این است که استفاده از نزول مختصات بلوکی به جای الگوریتم گرادیان کاهشی استانداردتر رایج است. این تمایز جزئی است و استفاده از گرادیان کاهشی استاندارد کاملاً معتبر است.

¹ minimal Frobenius

فصل ۱۰

ارزیابی الگوریتم‌های یادگیری

اکثر این کتاب بر روی استخراج الگوریتم و به دست آوردن مدل‌ها متمرکز شده است، اما ما هنوز باید به نحوه ارزیابی این مدل‌ها بپردازیم. فرمالیسم حداکثر احتمال برای استخراج الگوریتم‌های یادگیری برخی از نتایج سازگاری را ارائه می‌دهد، که در حد نمونه‌ها می‌توانیم نقطه همگرایی یک برآوردگر را مورد بحث قرار دهیم. با این حال، در عمل، ما می‌خواهیم الگوریتم‌ها را بر اساس یک نمونه محدود ارزیابی کنیم. محیطی را تصور کنید که در آن دو مدل را یاد می‌گیرید، مثلاً با استفاده از رگرسیون لجستیک با دو پارامتر تنظیم متفاوت. کدام یک از این دو مدل "بهتر" است؟ حتی بهتر است بگوییم چه معنایی دارد؟ آیا می‌خواهید بگویید که مدل برای این مشکل (تنظیم داده‌ها) بهتر است یا برای مشکلات متعدد؟ آیا ما سعی می‌کنیم الگوریتم‌ها یا مدل‌های به دست آمده را با هم مقایسه کنیم. از یک نمونه خاص از یک الگوریتم؟ چگونه می‌توانیم مطمئن باشیم که عملکرد اندازه گیری شده دقیقاً عملکردی را که انتظار داریم در داده‌های جدید ببینیم منعکس می‌کند؟ سوالاتی در مورد ویژگی‌های آن هدف و ویژگی‌های تجربی مدل‌های آموخته شده.

در این فصل، ابزارهای نظری و تجربی را برای ارزیابی بهتر ویژگی‌های الگوریتم‌های یادگیری ارائه می‌کنیم. ما با برخی از نتایج نظری نمونه محدود اولیه شروع می‌کنیم، که پیچیدگی کلاس مدل را به تعداد نمونه‌های مورد نیاز برای به دست آوردن یک تخمین منطقی از خطای مورد انتظار (خطای تعمیم) مرتبط می‌کند. این بخش همچنین ایده‌های بهینه سازی در یک کلاس تابع و اهداف ما برای به دست آوردن بهترین مدل از نظر خطای تعمیم را معرفی می‌کند. حوزه ای که با این نوع خصوصیات نظری سروکار دارد، نظریه یادگیری آماری نامیده می‌شود. ما یک نتیجه را با استفاده از نابرابری‌های تمرکز و پیچیدگی *Rademacher* برای توصیف پیچیدگی کلاس مدل مورد بحث قرار خواهیم داد. برای اطلاعات بیشتر، می‌توانید این آموزش را در مورد موضوع [۷] در نظر بگیرید.

سپس، نحوه مقایسه تجربی الگوریتم‌ها را مورد بحث قرار خواهیم داد. در اکثر تنظیمات دنیای واقعی، شما بین الگوریتم‌ها بر اساس عملکرد آنها در داده‌های موجود، یکی را انتخاب خواهید کرد. شما می‌خواهید این انتخاب منعکس کننده میزان عملکرد آن الگوریتم‌ها بر روی داده‌های جدید باشد. برای رسیدن به این هدف، ما در مورد چگونگی تقسیم داده‌ها و نحوه استفاده از آزمون‌های معناداری آماری برای ارائه سطحی از اطمینان به اینکه یک الگوریتم یا مدل از دیگری بهتر است، تحت برخی معیارهای خاص بحث خواهیم کرد. ما به ندرت قادر به نتیجه گیری قوی بر اساس آزمایش خواهیم بود، اما می‌توانیم شواهدی در مورد ویژگی‌های الگوریتم ایجاد کنیم.

این ابزارها مسلماً حیاتی ترین جنبه‌های استفاده صحیح از الگوریتم‌های یادگیری ماشین در عمل هستند. می‌توان یک مدل پیچیده را یاد گرفت، اما بدون درک نحوه عملکرد آن در عمل روی داده‌های جدید، استفاده واقعی از این مدل‌ها امکان‌پذیر

نیست. این که آیا یک الگوریتم برای مقاصد علمی استفاده می‌شود یا در سیستم‌های واقعی به کار می‌رود، داشتن درک درستی از ویژگی‌های آن چه از نظر تئوری و چه از لحاظ تجربی، کلیدی برای به دست آوردن نتایج مورد انتظار است. این فصل فقط شروع به خراش دادن سطح این ابزارها می‌کند، با این هدف که علاقه شما را برانگیزد و شما را به سمت مطالب بیشتری برای یادگیری در مورد ارزیابی هدایت کند.

۱-۱۰ مقدمه‌ای کوتاه بر مرزهای تعمیم

هدف ما در سراسر این کتاب این بوده است که تابعی را بر اساس مجموعه‌ای از مثال‌ها به دست آوریم که به دقت پیش‌بینی می‌کند: خطای کم مورد انتظار را در فضای نمونه‌های ممکن ایجاد می‌کند. با این حال، ما نمی‌توانیم خطای مورد انتظار را اندازه گیری کنیم. از نظر آماری می‌دانیم که با یک نمونه کافی می‌توانیم یک انتظار را تقریبی کنیم. در اینجا، ما این را با دقت بیشتری برای توابع آموخته شده تعیین می‌کنیم.

هدف ما به‌طور دقیق‌تر انتخاب تابعی از یک تابع کلاس H برای به حداقل رساندن یک تابع هزینه است: $\ell: \mathbb{R} \times \mathbb{R} \rightarrow [0, \infty)$ در انتظار روی همه جفت‌ها (\mathbf{x}, y)

$$\min_{f \in \mathcal{H}} \mathbb{E}[\ell(f(\mathbf{X}), \mathbf{Y})]$$

به عنوان مثال، در رگرسیون خطی، $\mathcal{H} = \{f: \mathbb{R}^d \rightarrow \mathbb{R} \mid f(\mathbf{x}) = \mathbf{x}^T \mathbf{w}, \text{ for any } \mathbf{w} \in \mathbb{R}^d\}$. این فضای توابع \mathcal{H} همه توابع خطی ممکن ورودی $\mathbf{x} \in \mathbb{R}^d$ را برای تولید یک خروجی اسکالر نشان می‌دهد. هدف ما در رگرسیون خطی، به حداقل رساندن یک پروکسی برای خطای واقعی مورد انتظار، یعنی خطای نمونه بود: $\frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i), y_i)$. یک سوال طبیعی مطرح می‌شود: آیا این خطای نمونه تخمین دقیقی از خطای مورد انتظار واقعی ارائه می‌دهد؟ و در مورد عملکرد تعمیم واقعی، یعنی خطای واقعی مورد انتظار چه چیزی به ما می‌گوید؟

بیا با یک مثال ساده با استفاده از رگرسیون خطی شروع کنیم. یک تابع محدود کلاس \mathcal{H} را فرض کنید، که در آن $\mathcal{H} = \{f: \mathbb{R}^d \rightarrow \mathbb{R} \mid f(\mathbf{x}) = \mathbf{x}^T \mathbf{w}, \text{ for any } \mathbf{w} \in \mathbb{R}^d \text{ such that } \|\mathbf{w}\|_2 \leq \mathbf{B}_w\}$. فرض کنید ویژگی‌های ورودی از یک فضای محدود می‌آیند، به طوری که برای همه \mathbf{x} ، $\|\mathbf{x}\|_2 \leq \mathbf{B}_x$ برای برخی اسکالر محدود $\mathbf{B}_x > 0$ ، و همچنین خروجی $y \in [-B_y, B_y]$ برای برخی از $B_y > 0$. فرض کنید از اتلاف $\ell(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$ استفاده می‌کنیم که (به صورت محلی) *Lipschitz* پیوسته است. برای منطقه محدود ما، با ثابت *Lipschitz* $c = B_y + B_x B_w$ این به این دلیل است که $|\hat{y}| \leq B_x B_w$ و

$$\left| \frac{d\ell(\hat{y}, y)}{d\hat{y}} \right| = |\hat{y} - y| \leq |\hat{y}| + |y| \leq B_y + B_x B_w$$

بعلاوه، چون $y \in [-B_y, B_y]$ ، می‌دانیم که هزینه به صورت محدود می‌شود

$$\ell(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2 \leq \frac{1}{2} (B_y^2 + B_x^2 B_w^2)$$

برای خطای تقریبی

$$\widehat{\text{Err}} d(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i), \mathbf{y}_i)$$

و خطای واقعی

$$\text{Err}(f) = \mathbb{E}[\ell(f(\mathbf{X}), \mathbf{Y})] = \int_{\mathbf{x} \times \mathbf{y}} p(\mathbf{x}, \mathbf{y}) \ell(f(\mathbf{x}), \mathbf{y}) d\mathbf{x} d\mathbf{y}$$

با استفاده از معادله ۱۰،۲ زیر، دریافت می‌کنیم که با احتمال $1 - \delta$ ، برای $\delta \in (0, 1]$

$$\text{Err}(f) \leq \widehat{\text{Err}} d(f) + \frac{2cB_x B_w}{\sqrt{n}} + \frac{1}{2} (B_y^2 + B_x^2 B_w^2) \quad (10.1)$$

با افزایش نمونه‌های n ، دو عبارت دوم ناپدید می‌شوند و خطای نمونه به خطای واقعی مورد انتظار نزدیک می‌شود. این کران میزان ناپدید شدن این اختلاف را نشان می‌دهد. برای اطمینان بیشتر δ - کوچک که $\ln(1/\delta)$ را بزرگتر می‌کند - به نمونه‌های بیشتری نیاز است تا ترم سوم کوچک باشد. این جمله سوم با استفاده از نابرابری‌های غلظت به دست می‌آید، که ما را قادر می‌سازد نرخ را بیان کنیم که میانگین نمونه به مقدار مورد انتظارش نزدیک می‌شود. برای مقادیر احتمالاً بزرگ ویژگی‌ها یا وزن‌های آموخته شده، عبارت دوم می‌تواند بزرگ باشد و دوباره به نمونه‌های بیشتری نیاز دارد. عبارت دوم ویژگی‌های کلاس تابع ما را منعکس می‌کند: یک کلاس ساده‌تر، با وزن‌های محدود کوچک، می‌تواند تخمین دقیق‌تری از هزینه در تعداد کمتری از نمونه‌ها داشته باشد. به طور کلی تر، این اندازه گیری پیچیدگی، پیچیدگی Rademacher نامیده می‌شود. 1 برای توابع خطی بالا، با نرم‌های 2 محدود برای \mathbf{x}, \mathbf{w} ، پیچیدگی Rademacher به صورت $R_n(\mathcal{H}) \leq B_x B_y / \sqrt{n}$ محدود می‌شود (به [۱۰، معادله ۳] مراجعه کنید).

در چند بخش بعدی، یک نتیجه تعمیم برای توابع کلی تر و همچنین پیش زمینه مورد نیاز برای تعیین آن نتیجه ارائه می‌دهیم.

۱-۱-۱۰ نابرابری‌های تمرکز

ما استفاده از نابرابری‌های تمرکز را با یک مثال رایج بررسی خواهیم کرد: نابرابری Hoeffding. برای تعمیم کران زیر، یک تعمیم به نام نابرابری McDiarmid استفاده می‌شود

برای $i.i.d.$ متغیرهای تصادفی X_1, \dots, X_n ، به طوری که $0 \leq X_i \leq 1$ ، اجازه دهید $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ میانگین نمونه باشد. سپس نابرابری هوفدینگ بیان می‌کند که برای هر کدام

$$\Pr(\bar{X} - \mathbb{E}[\bar{X}] \geq \epsilon) \leq \exp(-2n\epsilon^2)$$

ما با تنظیم این مقدار احتمال روی δ شروع می‌کنیم، به طوری که می‌توانیم با احتمال δ ، $\Pr(\bar{X} - \mathbb{E}[\bar{X}] \geq \delta)$ function(δ) بگوییم. ما می‌توانیم برای هر حساب δ حل کنیم تا بدست آوریم

$$\delta = \exp(-2n\epsilon^2) \Rightarrow \epsilon = \pm \sqrt{\frac{\ln(1/\delta)}{2n}}$$

ما می‌توانیم برای محدود کردن \bar{X} به نزدیک $\mathbb{E}[\bar{X}]$ از بالا و پایین ϵ را روی $\sqrt{\frac{\ln(1/\delta)}{2n}}$ یا $-\sqrt{\frac{\ln(1/\delta)}{2n}}$ تنظیم کنیم، ما آن را با احتمال $1 - \delta$ دریافت می‌کنیم، $|\bar{X} - \mathbb{E}[\bar{X}]| \leq |\epsilon| = \sqrt{\frac{\ln(1/\delta)}{2n}}$

این نابرابری غلظت، مفروضات کمی در مورد متغیرهای تصادفی ایجاد می‌کند و به هیچ فرض توزیعی نیاز ندارد. در نتیجه، نرخ همگرایی به میانگین واقعی تنها $1/\sqrt{n}$ است. نرخ‌های سریع تری را می‌توان با فرضیات بیشتر به دست آورد.

۲-۱-۱۰ پیچیدگی یک کلاس تابع

پیچیدگی *Rademacher* یک کلاس تابع، توانایی بیش از حد برازش توابع را در یک نمونه خاص مشخص می‌کند. کلاس توابعی که معمولاً پیچیده‌تر هستند، به احتمال زیاد می‌توانند نویز تصادفی را تطبیق دهند و بنابراین پیچیدگی *Rademacher* بالاتری دارند. پیچیدگی تجربی *Rademacher*، برای نمونه z_1, \dots, z_n - جایی که معمولاً $z_i = (x_i, y_i)$ را در نظر می‌گیریم - به صورت تعریف می‌شود

$$\hat{R}_n(\mathcal{H}) = \mathbb{E} \left[\max_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(\mathbf{x}_i) \right]$$

جایی که انتظار به پایان رسیده است *i. i. d.* متغیرهای تصادفی $\sigma_1, \dots, \sigma_n$ به طور یکنواخت از $\{-1, 1\}$ انتخاب شده است. این انتخاب نشان می‌دهد که چگونه کلاس تابع می‌تواند با این نویز تصادفی ارتباط داشته باشد. به عنوان مثال اگر $f(x)$ مقدار 1 یا -1 را پیش‌بینی می‌کند، مانند طبقه‌بندی باینری، در نظر بگیرید. اگر تابعی در کلاس توابع وجود داشته باشد که بتواند کاملاً با علامت σ_i نمونه برداری تصادفی مطابقت داشته باشد، آن تابع بالاترین مقدار $\sum_{i=1}^n \sigma_i f(\mathbf{x}_i)$ را تولید می‌کند. پیچیدگی تجربی *Rademacher* برای یک کلاس تابع زیاد است، اگر برای هر σ_i نمونه‌برداری تصادفی، چنین تابعی در کلاس تابع وجود داشته باشد (می‌تواند برای هر $\sigma_1, \dots, \sigma_n$ یک تابع متفاوت باشد). پیچیدگی *Rademacher* پیچیدگی تجربی *Rademacher* مورد انتظار است، بیش از همه نمونه‌های احتمالاً از n نمونه

برای کلاس‌های تابع با پیچیدگی *Rademacher* بالا، خطا در مجموعه آموزشی بعید است که منعکس کننده خطای تعمیم باشد، تا زمانی که تعداد نمونه کافی وجود داشته باشد. این در تعمیم محدود در بخش ۳ - ۱ - ۱۰ منعکس شده است.

اتصال به بعد VC: پیچیدگی یک کلاس تابع را می‌توان با بعد VC نیز مشخص کرد. ایده بعد VC برای مشخص کردن تعداد نقاطی که می‌توانند توسط یک کلاس تابع جدا شوند (یا خرد شوند). توابع ساده ابعاد VC پایینی دارند، زیرا به اندازه کافی پیچیده نیستند تا نقاط زیادی را از هم جدا کنند. توابع پیچیده تر، که مرزهای پیچیده را فعال می‌کنند، بعد VC بالاتری دارند. به عنوان مثال، برای توابع به شکل $f((x_1, x_2)) = \text{sign}(x_1 w_1 + x_2 w_2 + w_0)$ بعد VC برابر 3 است. به طور کلی، برای $x \in \mathbb{R}^d$ ، بعد VC برابر $d + 1$ است. بعد VC ایده‌ای مشابه با پیچیدگی *Rademacher* است، اما به طبقه‌بندی کننده‌های باینری محدود می‌شود. به همین دلیل، ما به طور مستقیم پیچیدگی *Rademacher* را مورد بحث قرار می‌دهیم، که برای طبقه‌بندی کننده‌های باینری می‌تواند برحسب بعد VC محدود شود. با لمای Sauer، ما معمولاً می‌توانیم پیچیدگی

$$\text{Rademacher یک کلاس فرضیه را با } \sqrt{\frac{2VC - \text{dimension} \ln n}{n}} \text{ محدود کنیم.}$$

۳-۱-۱۰ مرزهای تعمیم

کران تعمیم برای یک کلاس از مدل‌ها را می‌توان با ترکیب نابرابری‌های غلظت به انحراف کران از میانگین برای نمونه‌های کمتر، و استفاده از پیچیدگی $Rademacher$ برای محدود کردن تفاوت بین خطای نمونه و خطای واقعی مورد انتظار در همه توابع در کلاس تابع ما علاوه بر این باید مجموعه زبان‌ها را محدود کنیم. ما فرض می‌کنیم که هزینه‌های $Lipschitz$ با ثابت c هستند، به این معنی که آنها خیلی سریع در یک منطقه تغییر نمی‌کنند، با c نشان دهنده نرخ تغییر است. علاوه بر این، ما همچنین فرض می‌کنیم که هزینه با b محدود می‌شود، یعنی به مقادیر $[-b, b]$ می‌رسد. همانطور که در بالا، اگر z_1, \dots, z_n یک $i.i.d$ باشد، سپس با احتمال $1 - \delta$ ، برای هر $f \in \mathcal{H}$

$$\mathbb{E}[\ell(f(\mathbf{X}), \mathbf{Y})] \leq \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i), y_i) + 2c R_n(\mathcal{H}) + b \sqrt{\frac{\ln(1/\delta)}{2n}} \quad (10.2)$$

برای بیان قضیه و اثبات دقیق‌تر، به [۳، قضیه ۷] و [۱۰، قضیه ۱] مراجعه کنید.

۲-۱۰ مقایسه الگوریتم‌های یادگیری

برای ارزیابی تجربی الگوریتم‌ها، می‌توانیم تنظیماتی را با یک یا چند الگوریتم در یک یا چند مجموعه داده در نظر بگیریم. بسته به تنظیمات، ارزیابی‌های متفاوتی به کار گرفته می‌شود. برای یک نمای کلی خوب از ارزیابی الگوریتم‌های یادگیری ماشین، [۹] را ببینید.

در حال حاضر، اجازه دهید با یک مورد ساده شروع کنیم، که در آن دو الگوریتم را با هم مقایسه کرده و از تست دو جمله‌ای استفاده می‌کنیم. فرض کنید مجموعه‌ای از مسائل یادگیری D_1, D_2, \dots, D_m و مایل به مقایسه الگوریتم‌های یادگیری a_1 و a_2 هستند. ما می‌توانیم چنین مقایسه‌ای را با استفاده از آزمون شمارش به صورت زیر انجام دهیم: برای هر مجموعه داده، هر دو الگوریتم بر حسب معیار عملکرد انتخابی ارزیابی می‌شوند و الگوریتمی با دقت عملکرد بالاتر برنده می‌شود، در حالی که به دیگری هزینه داده می‌شود. (در صورت عملکرد دقیقاً یکسان، ما می‌توانیم یک برد / باخت به صورت تصادفی ارائه دهیم).

	D_1	D_2	D_3	D_4		D_{m-1}	D_m
a_1	1	0	1	1	...	0	1
a_2	0	1	0	0		1	0

جدول ۱۰/۱: یک آزمون شمارش که در آن الگوریتم‌های یادگیری a_1 و a_2 بر روی مجموعه‌ای از m مجموعه داده‌های مستقل مقایسه می‌شوند. الگوریتمی با عملکرد بهتر در یک مجموعه داده خاص، برد (1) را جمع‌آوری می‌کند، در حالی که الگوریتم دیگر هزینه (0) را جمع‌آوری می‌کند.

ما اکنون علاقه مند به ارائه شواهد آماری هستیم که می‌گویید الگوریتم a_1 بهتر از الگوریتم a_2 است. فرض کنید a_1 تعداد k برد از m داشته باشد و الگوریتم a_2 دارای $m - k$ برد باشد، همانطور که در جدول ۱۰/۱ نشان داده شده است. ما می‌خواهیم فرضیه صفر H_0 را ارزیابی کنیم که الگوریتم‌های a_1 و a_2 عملکرد یکسانی دارند با ارائه یک فرضیه جایگزین H_1 مبنی بر اینکه الگوریتم a_1 بهتر از a_2 است. به اختصار،

$$H_0: \text{quality}(a_1) = \text{quality}(a_2)$$

$$H_1: \text{quality}(a_1) > \text{quality}(a_2)$$

اگر فرضیه صفر درست باشد، برد/ باخت در هر مجموعه داده به همان اندازه محتمل خواهد بود و با تغییرات جزئی تعیین می‌شود. بنابراین، احتمال برد در هر مجموعه داده تقریباً برابر با $p = 1/2$ خواهد بود. اکنون، می‌توانیم با استفاده از توزیع دوجمله‌ای، این احتمال را بیان کنیم که الگوریتم a_1 دارای k برنده یا بیشتر را تحت فرضیه صفر می‌برد.

$$P = \sum_{i=k}^m \binom{m}{i} p^i (1-p)^{m-i}$$

و از آن به عنوان P -value یاد کنید. این مقدار احتمال پیروزی k ، به اضافه احتمال برد $k+1$ ، تا احتمال m برد، تحت فرضیه صفر است. یک رویکرد معمولی در این موارد، ایجاد یک مقدار معنی‌داری است، مثلاً $\alpha = 0.05$ و رد فرضیه صفر اگر $P \leq \alpha$ باشد. اگر مقدار P بزرگتر از α باشد، می‌گوییم شواهد کافی برای رد H_0 وجود ندارد. برای مقادیر P به اندازه کافی پایین، ممکن است به این نتیجه برسیم که شواهد کافی وجود دارد که الگوریتم a_1 بهتر از الگوریتم a_2 است.

انتخاب آستانه اهمیت α تا حدودی دلخواه است. به طور معمول، 5% یک مقدار معقول است، اما مقادیر پایین‌تر نشان می‌دهد که وضعیت خاص k برد از m بسیار بعید است، که می‌توانیم شواهدی برای رد H_0 بسیار قوی در نظر بگیریم. توانایی رد فرضیه صفر تا حدی اطمینان می‌دهد که نتیجه تصادفی رخ نداده است.

به طور کلی‌تر، می‌توانیم آزمون‌های معناداری آماری دیگری را بر اساس توزیع معیارهای عملکرد در نظر بگیریم. در مثال بالا، توزیع دوجمله‌ای مناسب بود. اگر در عوض خطاهای واقعی مجموعه داده‌ها را در نظر بگیریم، آنگاه جفت مقادیر واقعی خواهیم داشت. در این مورد، اگر هر دو خطا به طور نرمال توزیع شده باشند و واریانس مشابهی داشته باشند، یک انتخاب رایج، آزمون t زوجی است. آزمون t زوجی تفاوت‌های نمونه بین الگوریتم‌ها را می‌گیرد (خط ۳ در جدول ۱۰، ۲). d_1, \dots, d_m از آنجایی که مجدداً فرضیه صفر ما این است که الگوریتم‌ها به یک اندازه عمل می‌کنند، در فرضیه صفر میانگین این تفاوت‌ها ۰ است. اگر تفاوت‌ها به طور نرمال توزیع شوند، برای میانگین نمونه $\bar{d} = \frac{1}{m} \sum_{i=1}^m d_i$ و انحراف استاندارد نمونه $S_d = \sqrt{\frac{1}{m-1} \sum_{i=1}^m (d_i - \bar{d})^2}$ ، متغیر تصادفی $t = \frac{\bar{d}-0}{S_d/\sqrt{m}}$ بر اساس توزیع t دانش‌آموز توزیع می‌شود. توزیع t - Student تقریباً مانند یک توزیع نرمال است، با پارامتر درجه آزادی $m-1$ که باعث می‌شود با بزرگتر شدن m ، توزیع بیشتر شبیه توزیع نرمال به نظر برسد.

اکنون می‌توانیم در مورد احتمال این متغیر تصادفی T نسبت به آمار محاسبه شده بپرسیم. اگر ما فقط به دانستن اینکه آیا الگوریتم ۱ بهتر از الگوریتم ۲ است اهمیت می‌دهیم، یک آزمایش یک دنباله انجام می‌دهیم. اگر احتمال اینکه T بزرگتر از t باشد، یعنی $p = Pr(T > t)$ ، کوچک باشد، شواهدی به دست می‌آوریم که نشان می‌دهد الگوریتم ۱ بهتر از الگوریتم ۲ است. ما می‌توانیم ترتیب تفاوت را مبادله کنیم. اگر $p = Pr(T > -t)$ کوچک باشد، شواهدی به دست می‌آوریم که نشان می‌دهد الگوریتم ۲ بهتر از الگوریتم ۱ است. در عوض یک تست دو طرفه می‌پرسد که آیا این دو الگوریتم متفاوت هستند یا خیر. در این مورد، از $p = Pr(T > |t|)$ استفاده می‌شود.

	D_1	D_2	D_3	D_4		D_{m-1}	D_m
a_1	0.11	0.08	0.15	0.12	...	0.07	0.09
a_2	0.10	0.09	0.11	0.12	...	0.10	0.09
d	0.01	-0.01	0.04	0.0	...	-0.03	0.0

جدول ۱۰/۲: جدولی از خطاها برای دو الگوریتم یادگیری a_1 و a_2 بر روی مجموعه ای از m مجموعه داده های مستقل مقایسه شده است. سطر آخر شامل تفاوت‌هایی است که برای آزمون t زوجی استفاده می‌شود.

۳-۱۰ به دست آوردن نمونه‌های خطا

یک مرحله کلیدی در مقایسه الگوریتم‌ها، به دست آوردن معیارهای معتبر عملکرد برای مقایسه است. تا اینجا ما فرض کردیم که اینها داده شده است. یک روش برای به دست آوردن نمونه‌های بی طرفانه از خطا، نگه داشتن یک مجموعه تست نگهدارنده است. تصور کنید m نمونه در رزرو تنظیم شده است، که الگوریتم‌ها روی آن‌ها آموزش داده نشده‌اند و تا زمانی که آماده ارزیابی نباشیم نمی‌توانیم به آن‌ها نگاه کنیم. ما می‌توانیم دو مدل را در مجموعه آموزشی آموزش دهیم، و سپس m نمونه‌های جفتی از خطا را بدست آوریم. سپس می‌توانیم از آزمون t زوجی استفاده کنیم تا ادعا کنیم که آیا این دو مدل از نظر آماری به طور معنی‌داری برای مسئله متفاوت هستند یا خیر.

یک رویکرد جایگزین برای به دست آوردن تخمین‌های خطا، استفاده از تکنیک‌های نمونه‌گیری مجدد از کل مجموعه داده است. دو تکنیک نمونه‌گیری مجدد رایج عبارتند از اعتبارسنجی متقاطع k -fold و نمونه‌برداری مجدد $bootstrap$. در مرحله اول، داده‌ها به k مجموعه‌های ناهمگون (فولد) تقسیم می‌شوند. این مدل روی $k - 1$ چین آموزش داده می‌شود و روی چین دیگر آزمایش می‌شود. این k بار تکرار می‌شود که در آن هر فولد به عنوان تای تست عمل می‌کند. این رویکرد محیط یادگیری رایج را شبیه سازی می‌کند که در آن مجموعه‌های آموزشی و آزمون از هم جدا هستند. تخمین‌های عملکرد k به دست آمده عمدتاً مستقل هستند، با برخی وابستگی‌ها به دلیل وابستگی‌های بین مجموعه‌های آموزشی در سراسر k اجراها معرفی شده‌اند. برخی از تعصبات اضافی ارائه شده از این واقعیت است که ما مدل را در کل مجموعه آموزشی اجرا نمی‌کنیم، بلکه تخمینی از خطای الگوریتم آموزش داده شده بر روی $(n/k) - n$ دریافت می‌کنیم. برای هر مدل نهایی که پس از انجام این ارزیابی‌ها وارد تولید می‌شود، احتمالاً کل مجموعه n نمونه را آموزش خواهیم داد.

نمونه مجدد بوت استرپ با داده‌ها از ایده پشت بوت استرپ استفاده می‌کند: داده‌ها مدل معقولی از داده‌ها را تشکیل می‌دهند. با نمونه برداری از داده‌ها، مانند نمونه برداری از توزیعی است که داده‌ها را تولید کرده است. برای ایجاد تقسیم‌های آموزشی/آزمون، داده‌ها با جایگزینی برای ایجاد مجموعه آموزشی نمونه‌برداری می‌شوند و نمونه‌های استفاده نشده باقی‌مانده برای آزمایش استفاده می‌شوند. اگر k نمونه مجدد بدست آید، دوباره k معیارهای عملکرد را بدست می‌آوریم و می‌توانیم میانگین نمونه عملکرد را در تقسیمات مختلف بدست آوریم و از آزمون معناداری آماری استفاده کنیم.

برای درک بهتر خواص این دو رویکرد، به توضیح کامل و قابل دسترس در [۸، فصل ۵] مراجعه کنید.

۴-۱۰ معیارهای عملکرد برای مدل‌های طبقه بندی

در طبقه بندی، انواع معیارهای عملکردی وجود دارد که اهمیت نسبی پیش بینی‌های نادرست را برای هر یک از کلاس‌ها منعکس می‌کند. به عنوان مثال، پیش‌بینی اینکه بیمار واقعاً بیمار است (منفی کاذب) می‌تواند مضرتر باشد، که در نتیجه تصمیم به انجام ندادن تشخیص بیشتر می‌شود و در نتیجه باعث عوارض جدی ناشی از عدم درمان بیماری می‌شود. هنگام آموزش و ارزیابی الگوریتم‌های طبقه‌بندی، این اولویت‌ها باید کدگذاری شوند. جدول ۱۰،۳ برخی از اصطلاحات را برای بحث در مورد عملکرد مدل‌های طبقه بندی خلاصه می‌کند.

Name	Symbol	Definition
Classification error	$error$	$error = \frac{fp+fn}{tp+fp+tn+fn}$
Classification accuracy	$accuracy$	$accuracy = 1 - error$
True positive rate	tpr	$tpr = \frac{tp}{tp+fn}$
False negative rate	fnr	$fnr = \frac{fn}{tp+fn}$
True negative rate	tnr	$tnr = \frac{tn}{tn+fp}$
False positive rate	fpr	$fpr = \frac{fp}{tn+fp}$
Precision	pr	$pr = \frac{tp}{tp+fp}$
Recall	rc	$rc = \frac{tp}{tp+fn}$

جدول ۱۰/۳: برخی از معیارهای طبقه بندی.

- [1] P Auer, M Herbster, and Manfred K Warmuth. Exponentially many local minima for single neurons. In *Advances in Neural Information Processing Systems*, 1996.
- [2] A Banerjee, S Merugu, I S Dhillon, and J Ghosh. Clustering with Bregman divergences. *Journal of Machine Learning Research*, 2005.
- [3] Peter L Bartlett and Shahar Mendelson. Rademacher and Gaussian Complexities: Risk Bounds and Structural Results. *Journal of Machine Learning Research*, 2002.
- [4] Amir Beck and Marc Teboulle. A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems. *SIAM J. Imaging Sciences*, 2009.
- [5] Léon Bottou and Yann Le Cun. On-line learning for very large data sets. *Applied Stochastic Models in Business and Industry*, 2005.
- [6] Léon Bottou. Online learning and stochastic approximations. *Online Learning and Neural Networks*, 1998.
- [7] Olivier Bousquet, Stéphane Boucheron, and Gábor Lugosi. Introduction to Statistical Learning Theory. In *Advanced Lectures on Machine Learning*. Springer Berlin Heidelberg, 2004.
- [8] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An Introduction to Statistical Learning*. Springer New York, 2013.
- [9] Nathalie Japkowicz and Mohak Shah. *Evaluating Learning Algorithms - A Classification Perspective*. Cambridge University Press, 2011.
- [10] Sham M Kakade, Karthik Sridharan, and Ambuj Tewari. On the Complexity of Linear Prediction: Risk Bounds, Margin Bounds, and Regularization. In *Advances in Neural Information Processing Systems*, 2008.
- [11] Diederik P Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *International Conference on Learning Representations*, 2015.
- [12] Lei Le and Martha White. Global optimization of factor models using alternating minimization. *arXiv.org*, 2016.
- [13] Julien Mairal, Francis Bach, Jean Ponce, Guillermo Sapiro, and Andrew Zisserman. Supervised dictionary learning. In *Advances in Neural Information Processing Systems*, 2009.
- [14] Jorge Nocedal. Updating quasi-Newton matrices with limited storage. *Mathematics of Computation*, 1980.
- [15] Bruno A Olshausen and David J Field. Sparse coding with an overcomplete basis set: a strategy employed by V1? *Vision Research*, 1997.
- [16] K B Petersen. *The matrix cookbook*. Technical University of Denmark, 2004.
- [17] Dinah Shender and John Lafferty. Computation-Risk Tradeoffs for Covariance-Thresholded Regression. *International Conference on Machine Learning*, 2013.
- [18] Ajit P Singh and Geoffrey J Gordon. A unified view of matrix factorization models. In *European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases*, 2008.

- [19] Larry Wasserman. All of Statistics: A Concise Course in Statistical Inference. Springer, 2004.
- [20] Martha White. Regularized factor models. PhD thesis, University of Alberta, 2014.
- [21] Matthew D Zeiler. ADADELTA: An Adaptive Learning Rate Method. arXiv.org, 2012.

پیوست A

مطالب اضافی برای نظریه احتمال

A.1 بدیهیات احتمال

ما می‌توانیم ویژگی‌های دیگری را از تعریف اولیه مجموعه رویدادهای قابل اندازه‌گیری (جبر سیگما) و توزیع‌های احتمال استخراج کنیم. تعریف میدان سیگما مستلزم آن است که \mathcal{E} تحت هر دو تعداد محدود و قابل شمارش نامتناهی از عملیات مجموعه پایه (اتحاد، تقاطع، متمم و اختلاف مجموعه) بسته شود. اتحاد عملیات و تکمیل در تعریف است. برای تقاطع، می‌توانیم از قوانین دی مورگان¹ استفاده کنیم: $\bigcap A_i = (\bigcup A_i^c)^c$ و $\bigcup A_i = (\bigcap A_i^c)^c$. هر تقاطع مجموعه‌ها در \mathcal{E} باید دوباره در \mathcal{E} باشد زیرا \mathcal{E} تحت اتحاد و مکمل بسته است. بنابراین، یک میدان سیگما نیز در زیر تقاطع بسته است. به طور مشابه برای اختلاف مجموعه، می‌توانیم $A_1 - A_2 = (A_1 \cap A_2)^c \cap A_1$ را بنویسیم، که سپس به $A_1 - A_2 \in \mathcal{E}$ اشاره می‌کند زیرا $A_1 \cap A_2 \in \mathcal{F} \Rightarrow (A_1 \cap A_2)^c \in \mathcal{E} \Rightarrow (A_1 \cap A_2)^c \cap A_1 \in \mathcal{F}$ که همه شرایط فوق نشان می‌دهد که $\emptyset \in \mathcal{E}$ و $\Omega \in \mathcal{E}$ که در آن مجموعه خالی است.

برای توزیع احتمال $P: \mathcal{E} \rightarrow [0, 1]$ ، ما نیاز داریم

$$1. P(\Omega) = 1$$

$$2. A_1, A_2, \dots \in \mathcal{E}, A_i \cap A_j = \emptyset \forall i, j \Rightarrow P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i)$$

تاپل (Ω, \mathcal{E}, P) فضای احتمال نامیده می‌شود. به نظر می‌رسد بصری است که شرط دوم را می‌توان با اتحادیه‌ای از مجموعه‌های محدود (نیاز ساده‌تر افزایش به جای σ افزایشی) جایگزین کرد. با این حال، برای فیلدهای سیگما، بسته شدن تحت اتحادیه‌های محدود ممکن است منجر به بسته شدن تحت اتحادیه‌های نامحدود نشود.

زیبایی این بدیهیات در فشردگی و ظرافت آنها نهفته است. بسیاری از عبارات مفید را می‌توان از بدیهیات احتمال استخراج کرد. برای مثال، واضح است که $P(\emptyset) = 0$ یا $P(A^c) = 1 - P(A)$. به طور مشابه، بسته شدن تحت اتحادیه‌های نامتناهی مجموعه‌های غیرمتناسب (σ افزایشی) به بسته شدن محدود (افزودنی) دلالت دارد، زیرا مجموعه‌های باقیمانده را می‌توان به مجموعه تهی $\emptyset: \forall A_1, A_2 \in \mathcal{E}$ تنظیم کرد با: $A_1 \cap A_2 = \emptyset$ مجموعه $A_i = \emptyset$ در $i > 2$ برای بدست آوردن $P(A_1 \cup A_2) = P(\bigcup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i) = P(A_1) + P(A_2)$ فرمول دیگری که اهمیت ویژه‌ای دارد را می‌توان با

¹ De Morgan's laws

در نظر گرفتن پارتیشن^۱ فضای نمونه بدست آورد. یعنی مجموعه‌ای از k مجموعه‌های غیر همپوشانی $\{B_i\}_{i=1}^k$ به طوری که $\Omega = \bigcup_{i=1}^k B_i$. یعنی اگر A هر مجموعه‌ای در Ω باشد و اگر $\{B_i\}_{i=1}^k$ پارتیشنی از Ω باشد نتیجه می‌شود که

$$\begin{aligned} P(A) &= P(A \cap \Omega) \\ &= P\left(A \cap \left(\bigcup_{i=1}^k B_i\right)\right) \\ &= P\left(\left(\bigcup_{i=1}^k A \cap B_i\right)\right) \quad (A.1) \\ &= \sum_{i=1}^k p(A \cap B_i) \end{aligned}$$

که در آن خط آخر از بدیهیات احتمال پیروی می‌کند. ما به این عبارت به عنوان قانون جمع اشاره خواهیم کرد. عبارت مهم دیگری که در اینجا بدون مشتق نشان داده شده است این است که $P(A \cup B) = P(A) + P(B) - P(A \cap B)$.

A.2 چند pmf مفید دیگر

ما چند نمونه دیگر از pmf ها را در اینجا ارائه می‌کنیم، عمدتاً به این منظور که مثال‌های ملموس اضافی می‌توانند برای درک مفید باشند. ما بیشتر از این pmf ها در این یادداشت ها استفاده نخواهیم کرد.

توزیع دو جمله‌ای برای توصیف دنباله‌ای از n آزمایش برنولی مستقل و توزیع شده یکسان (*i.i.d.*) استفاده می‌شود. در هر مقدار k در فضای نمونه، توزیع این احتمال را می‌دهد که موفقیت دقیقاً k بار از n آزمایش اتفاق افتاده است، که البته $0 \leq k \leq n$ به طور فرمول‌بندی، $\Omega = \{0, 1, \dots, n\}$ برای $\forall k \in \Omega$ دو جمله‌ای pmf به این صورت تعریف می‌شود

$$p(k) = \binom{n}{k} \alpha^k (1 - \alpha)^{n-k}$$

که در آن $\alpha \in (0, 1)$ ، مانند قبل، پارامتری است که احتمال موفقیت در یک آزمایش را نشان می‌دهد. در اینجا، ضریب دو جمله‌ای

$$\binom{n}{k} = \frac{n!}{k! (n - k)!}$$

همه راه‌هایی را برمی‌شمارد که از طریق آن می‌توان k عنصر را از فهرستی از n عنصر انتخاب کرد (به عنوان مثال، ۳ روش مختلف وجود دارد که از طریق آنها می‌توان $k = 2$ عنصر را از یک گروه $n = 3$ عنصر انتخاب کرد). ما به توزیع دو جمله‌ای با پارامترهای n و α به عنوان $Binomial(n, \alpha)$ اشاره خواهیم کرد. آزمایشی که منجر به توزیع دو جمله‌ای می‌شود را می‌توان به موقعیتی با بیش از دو نتیجه ممکن تعمیم داد. این آزمایش منجر به یک تابع جرم احتمال چند بعدی (یک بعد در هر نتیجه ممکن) به نام توزیع چند جمله‌ای می‌شود.

¹ partition

توزیع هندسی همچنین برای مدل سازی دنباله ای از آزمایش های مستقل برنولی با احتمال موفقیت α استفاده می شود. در هر نقطه $\Omega \in k$ ، این احتمال را می دهد که اولین موفقیت دقیقاً در آزمایش k -ام اتفاق می افتد. در اینجا، $\Omega = \{1, 2, \dots\}$ و برای $\forall k \in \Omega$

$$p(k) = (1 - \alpha)^{k-1} \alpha$$

که $\alpha \in (0, 1)$ یک پارامتر است. توزیع هندسی، $\text{Geometric}(\alpha)$ ، بر روی یک فضای نمونه بی نهایت تعریف شده است. یعنی $\Omega = N$

برای توزیع فراهندسی، جمعیت محدودی از عناصر N از دو نوع (مثلاً موفقیت و شکست) را در نظر بگیرید که K از یک نوع (مثلاً موفقیت) هستند. این آزمایش شامل ترسیم n عنصر، بدون جایگزینی، از این جمعیت است، به طوری که عناصر باقی مانده در جمعیت از نظر انتخاب شدن در قرعه کشی بعدی، یکسان هستند. احتمال استخراج k موفقیت از n آزمایش را می توان به صورت توصیف کرد

$$p(k) = \frac{\binom{K}{k} \cdot \binom{N-K}{n-k}}{\binom{N}{n}}$$

که در آن $0 \leq n \leq N$ و $0 \leq k \leq n$. توزیع ابر هندسی ارتباط نزدیکی با توزیع دو جمله ای دارد که در آن عناصر با جایگزینی ترسیم می شوند ($\alpha = K/N$). در آنجا، احتمال ترسیم موفقیت در آزمایشات بعدی تغییر نمی کند. ما به توزیع فوق هندسی به عنوان $\text{Hypergeometric}(n, N, K)$ اشاره خواهیم کرد.

A.3 چند pdf مفید دیگر

همانطور که در بالا ذکر شد، ما چند فایل pdf دیگر را برای ارائه مثال های عینی ارائه می دهیم، اگرچه به صراحت از این pdf ها در یادداشت ها استفاده نخواهیم کرد.

توزیع \log نرمال یک تغییر توزیع نرمال است. در اینجا، برای $\Omega = (0, \infty)$ چگالی \log نرمال را می توان به این صورت بیان کرد

$$p(\omega) = \frac{1}{\omega \sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(\ln\omega - \mu)^2}$$

که در آن $\mu \in \mathbb{R}$ و $\sigma > 0$ پارامترها هستند. ما به این توزیع به عنوان $\text{Lognormal}(\mu, \sigma^2)$ یا $\ln N(\mu, \sigma^2)$ اشاره خواهیم کرد.

توزیع $Gumbel$ متعلق به کلاس توزیع های ارزش شدید است. تابع چگالی احتمال آن بر روی $\Omega = \mathbb{R}$ به این صورت تعریف می شود

$$p(\omega) = \frac{1}{\beta} e^{-\frac{\omega-\alpha}{\beta}} e^{-e^{-\frac{\omega-\alpha}{\beta}}}$$

که $\alpha \in \mathbb{R}$ پارامتر مکان و $\beta > 0$ پارامتر مقیاس است. ما به این توزیع به عنوان $\text{Gumbel}(\alpha, \beta)$ اشاره خواهیم کرد.

توزیع پارتو^۱ برای مدل سازی رویدادهایی با موارد نادر مقادیر شدید مفید است. تابع چگالی احتمال آن بر روی $\Omega = [\omega_{\min}, \infty)$ به صورت تعریف شده است

$$p(\omega) = \frac{\alpha \omega_{\min}}{\omega^{\alpha+1}}$$

که $\alpha > 0$ یک پارامتر و $\omega_{\min} > 0$ مینیمم مقدار مجاز برای ω است. ما به توزیع پارتو به عنوان $\text{Pareto}(\alpha, \omega_{\min})$ اشاره خواهیم کرد. هنگامی که $\alpha \in (0, 2]$ به یک ویژگی بدون مقیاس منجر می شود.

A.4 متغیرهای تصادفی

در بسیاری از موقعیت ها، ما می خواهیم از مدل سازی احتمالی در مجموعه ها (به عنوان مثال، گروهی از افراد) استفاده کنیم که در آن عناصر می توانند با توصیفگرهای مختلف مرتبط شوند. به عنوان مثال، یک فرد ممکن است با سن، قد، شهروندی، ضریب هوشی یا وضعیت تأهل او مرتبط باشد و ممکن است ما به رویدادهای مرتبط با این توصیف کننده ها علاقه مند باشیم. در موقعیت های دیگر، ممکن است به تبدیل فضاهای نمونه مانند فضاهایی که مربوط به دیجیتالی کردن سیگنال آنالوگ از یک میکروفون به مجموعه ای از اعداد صحیح بر اساس مجموعه ای از آستانه های ولتاژ هستند، علاقه مند باشیم. مکانیسم یک متغیر تصادفی پرداختن به همه چنین موقعیت هایی را به روشی ساده، دقیق و یکپارچه تسهیل می کند.

متغیر تصادفی متغیری است که از دیدگاه ناظر، مقادیر را به صورت غیر قطعی، با ترجیحات کلی متفاوت برای نتایج متفاوت می گیرد. اما از نظر ریاضی، تابعی است که یک فضای نمونه را به فضای دیگر نگاشت می کند، با چند اخطار فنی که بعداً معرفی خواهیم کرد. اجازه دهید نیاز به متغیرهای تصادفی را تحریک کنیم. یک فضای احتمال (Ω, E, P) را در نظر بگیرید، که در آن Ω مجموعه ای از افراد است و اجازه دهید احتمال خوشحالی یک فرد تصادفی انتخاب شده $\omega \in \Omega$ را بررسی کنیم (ممکن است فرض کنیم یک روش تشخیصی برای ارزیابی وضعیت هر فرد داریم). ما با تعریف یک رویداد A شروع می کنیم

$$A = \{\omega \in \Omega : \text{Status}(\omega) = \text{happy}\}$$

و به سادگی احتمال این رویداد را محاسبه کنید. این یک رویکرد کاملاً قانونی است، اما می توان آن را با استفاده از مکانیسم متغیر تصادفی بسیار ساده کرد. ابتدا توجه می کنیم که از نظر فنی، روش تشخیصی ما با یک تابع مطابقت دارد: $\Omega \rightarrow S$ که فضای نمونه Ω را به یک فضای نمونه باینری جدید نگاشت می کند $S = \{\text{happy}, \text{not happy}\}$. جالبتر اینکه رویکرد ما همچنین توزیع احتمال P را به یک توزیع احتمال P_{status} جدید که بر روی برخی از جبر سیگما S تعریف شده است، ترسیم می کند. مثلاً E_{status} (برای اینکه نگاشت مطابق انتظار عمل کند، E_{status} باید مجموعه توان S باشد). اکنون می توانیم ببینیم که می توانیم وضعیت $P(\{\text{happy}\})$ را از روی احتمال رویداد A محاسبه کنیم. به عنوان مثال، $P_{\text{status}}(\text{happy}) = P(A)$. این یک نماد به هم ریخته است، بنابراین ممکن است بخواهیم آن را با استفاده از $P(\text{Status} = \text{happy})$ ، که در آن وضعیت یک "متغیر تصادفی" است، ساده کنیم.

از حروف بزرگ x, y, \dots برای نشان دادن متغیرهای تصادفی (مانند وضعیت) و حروف کوچک x, y, \dots برای نشان دادن عناصر (مانند "happy") استفاده خواهیم کرد فضاهای جدید X, Y, \dots به طور کلی، احتمالات را به صورت $P(X = x)$ می نویسیم، که یک ریلیکسیشن^۲ نمادین از $P(\{\omega : X(\omega) = x\})$ یا $P(X \leq x)$ برای $P(\{\omega : X(\omega) \leq x\})$ است. زمانی

¹ Pareto distribution

² relaxation

که هم دامنه X پیوسته باشد. همچنین زمانی که نیاز به توضیح بیشتر در مورد متغیر تصادفی داشته باشیم، به توابع جرم یا چگالی احتمال مربوطه به صورت $p(x)$ یا $p_X(x)$ اشاره خواهیم کرد. این در واقع زمانی اتفاق می‌افتد که X یک مقدار خاص را بگیرد. مثلاً برای $x = 1$ ، $p_X(1)$ را خواهیم نوشت. قبل از اینکه به تعریف رسمی متغیرهای تصادفی بپردازیم، به دو مثال گویا نگاه خواهیم کرد.

مثال ۲۰: [پرتاب‌های متوالی یک سکه منصفانه.] فرآیندی از سه پرتاب سکه و دو متغیر تصادفی X و Y را در نظر بگیرید که در فضای نمونه تعریف شده اند. X را به عنوان تعداد سرها در اولین پرتاب و Y را به عنوان تعداد سرهای روی هر سه پرتاب تعریف می‌کنیم. هدف ما یافتن فضاهای احتمالی است که پس از تبدیل‌ها ایجاد می‌شوند.

ابتدا $\Omega = \{HHH, HHT, HTH, HTT, THH, THT, TTH, TTT\}$ و

ω	HHH	HHT	HTH	HTT	THH	THT	TTH	TTT
$X(\omega)$	1	1	1	1	0	0	0	0
$Y(\omega)$	3	2	2	1	2	1	1	0

بگذارید فقط روی متغیر Y تمرکز کنیم. واضح است که $Y : \Omega \rightarrow \{0, 1, 2, 3\}$ اما ما همچنین باید \mathcal{E}_Y و P_Y را پیدا کنیم. برای محاسبه P_Y ، یک روش ساده یافتن pmf $p(y)$ آن است. برای مثال، اجازه دهید $P_Y(\{2\}) = p_Y(2)$ را به صورت محاسبه کنیم

$$\begin{aligned} P_Y(\{2\}) &= P(Y = 2) \\ &= P(\{\omega : Y(\omega) = 2\}) \\ &= P(\{HHT, HTH, THH\}) \\ &= \frac{3}{8} \end{aligned}$$

به دلیل توزیع یکنواخت در فضای اصلی (Ω, \mathcal{E}, P) ، به روشی مشابه، می‌توانیم محاسبه کنیم که $P(Y = 0) = P(Y = 3) = 1/8$ و $P(Y = 1) = 3/8$. در این مثال، ما در نظر گرفتیم که $\mathcal{E} = P(\Omega)$ و $\mathcal{E}_Y = P(Y)$ به عنوان نکته پایانی، اشاره می‌کنیم که تمام تصادفی بودن در فضای احتمال اصلی (P, \mathcal{E}, Ω) تعریف می‌شود و فضای احتمال جدید (Y, \mathcal{E}_Y, P_Y) به سادگی آن را از طریق یک تبدیل قطعی به ارث می‌برد.

مثال ۲۱: [کوانتیزاسیون] (Ω, \mathcal{E}, P) را در نظر بگیرید که در آن $\Omega = [0, 1]$ ، $\mathcal{E} = B(\Omega)$ و P توسط یک pdf یکنواخت القا می‌شود. X را تعریف کنید: $X : \Omega \rightarrow \{0, 1\}$ به عنوان

$$X(\omega) = \begin{cases} 0 & \omega \leq 0.5 \\ 1 & \omega > 0.5 \end{cases}$$

و فضای احتمال تبدیل شده را پیدا کنید.

از نظر فنی، فضای نمونه را به $X = \{0, 1\}$ تغییر داده‌ایم. برای فضای رویداد $\mathcal{E}_X = P(X) = \emptyset, 0, 1, 0, 1$ ما می‌خواهیم توزیع احتمال جدید P_X را درک کنیم. ما داریم

$$\begin{aligned} p_X(0) &= P_X(\{0\}) \\ &= P(X = 0) \end{aligned}$$

$$= P(\{\omega: \omega \in [0,0.5]\})$$

$$= \frac{1}{2}$$

9

$$p_X(1) = P_X(\{1\})$$

$$= P(X = 1)$$

$$= P(\{\omega : \omega \in (0.5, 1]\})$$

$$= \frac{1}{2}$$

از اینجا به راحتی می‌توانیم ببینیم که $P_X(\{0, 1\}) = 1$ و $P_X(\emptyset) = 0$ ، و بنابراین P_X در واقع یک توزیع احتمال است. دوباره، P_X به طور طبیعی با استفاده از P تعریف می‌شود. بنابراین، فضای احتمال (Ω, \mathcal{E}, P) را به (X, \mathcal{E}_X, P_X) تبدیل کرده‌ایم.

A.4.1 تعریف رسمی متغیر تصادفی

اکنون به طور رسمی یک متغیر تصادفی تعریف می‌کنیم. با در نظر گرفتن یک فضای احتمال (Ω, \mathcal{E}, P) ، یک متغیر تصادفی $X: \Omega \rightarrow \mathcal{X}$ تابع X است به طوری که برای هر $A \in \mathcal{B}(\mathcal{X})$ این گونه است که $\{\omega : X(\omega) \in A\} \in \mathcal{E}$ است. نتیجه می‌شود که

$$P_X(A) = P(\{\omega: X(\omega) \in A\})$$

ذکر این نکته ضروری است که به طور پیش فرض، فضای رویداد یک متغیر تصادفی را فیلد $Borel X$ تعریف کرده‌ایم. این راحت است زیرا یک میدان بول از یک مجموعه قابل شمارش Ω مجموعه توان آن است. بنابراین، ما در حال کار با بزرگترین فضاهای رویداد ممکن برای متغیرهای تصادفی گسسته و پیوسته هستیم.

اکنون یک متغیر تصادفی گسسته X را در نظر بگیرید که روی (Ω, \mathcal{E}, P) تعریف شده است. همانطور که از مثال‌های قبلی می‌بینیم، توزیع احتمال X را می‌توان به صورت پیدا کرد

$$p(x) = P_{X(x)}$$

$$= P(\{\omega: X(\omega) = x\})$$

برای $\forall x \in \mathcal{X}$. احتمال یک رویداد A را می‌توان به صورت پیدا کرد

$$P_X(A) = P(\{\omega: X(\omega) \in A\})$$

$$= \sum_{x \in A} p(x)$$

برای $\forall A \subseteq \mathcal{X}$

مورد متغیرهای تصادفی پیوسته پیچیده‌تر است، اما به رویکردی شبیه به متغیرهای تصادفی گسسته کاهش می‌یابد. در اینجا ابتدا یک تابع توزیع تجمعی (cdf) را به صورت تعریف می‌کنیم

$$\begin{aligned}
F_X(t) &= P_X(\{x: x \leq t\}) \\
&= P(\{\omega : X(\omega) \leq t\}) \\
&= P(X \leq t)
\end{aligned}$$

که در آن $P(X \leq t)$ ، مانند قبل، سوء استفاده جزئی از نماد را نشان می‌دهد. اگر تابع توزیع تجمعی قابل تمایز باشد، تابع چگالی احتمال یک متغیر تصادفی پیوسته به صورت تعریف می‌شود.

$$p(x) = \left. \frac{dF_X(t)}{dt} \right|_{t=x}$$

متناوباً، اگر $p(x)$ وجود داشته باشد، آنگاه

$$F_X(t) = \int_{-\infty}^t p(x) dx$$

برای هر $t \in \mathbb{R}$ تمرکز ما منحصراً بر روی متغیرهای تصادفی است که توابع چگالی احتمال خود را دارند. با این حال، برای یک دید کلی تر، ما همیشه باید "اگر وجود دارد" را در هنگام مراجعه به فایل‌های pdf در نظر داشته باشیم.

احتمال اینکه یک متغیر تصادفی مقداری از بازه (a, b) بگیرد اکنون می‌تواند به صورت محاسبه شود.

$$\begin{aligned}
P_X((a, b]) &= P(a < X \leq b) \\
&= \int_a^b p(x) dx \\
&= F_X(b) - F_X(a)
\end{aligned}$$

که از ویژگی‌های یکپارچه سازی به دست می‌آید.

حالا فرض کنید که متغیر تصادفی X فضای احتمال (Ω, \mathcal{E}, P) را به $(X, B(\mathcal{X}), P_X)$ تبدیل می‌کند. برای توصیف فضای احتمال حاصل، معمولاً از توابع جرم و چگالی احتمال القاکننده P_X استفاده می‌کنیم. به عنوان مثال، اگر P_X توسط توزیع گاوسی با پارامترهای μ و σ^2 القا شود، از

$$X: \mathcal{N}(\mu, \sigma^2) \quad \text{or} \quad X \sim \mathcal{N}(\mu, \sigma^2)$$

هر دو نماد نشان می‌دهند که تابع چگالی احتمال برای متغیر تصادفی X است

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

چگالی گاوسی به طور ضمنی تعریف می‌کند که $X = \mathbb{R}$. اما این نقطه سطحی است زیرا ما همیشه می‌توانیم دامنه یک تابع چگالی را به \mathbb{R} گسترش دهیم و هر جا که تابع اصلی تعریف نشده است، $p(x) = 0$ را تنظیم کنیم.

گروهی از d متغیرهای تصادفی $\{X_i\}_{i=1}^d$ که در فضای احتمال یکسان (Ω, \mathcal{E}, P) تعریف شده‌اند، بردار تصادفی یا متغیر تصادفی چند متغیره (چند بعدی) نامیده می‌شوند. ما قبلاً نمونه‌ای از یک بردار تصادفی ارائه شده توسط متغیرهای تصادفی (X, Y) را در مثال ۲۰ دیده ایم. یعنی $\{X_i: i \in \mathcal{T}\}$ که در آن \mathcal{T} یک مجموعه شاخص است که معمولاً به عنوان مجموعه‌ای از شاخص‌های زمانی تفسیر می‌شود. در مورد شاخص‌های زمان گسسته (به عنوان مثال، $\mathcal{T} = N$) فرآیند تصادفی یک فرآیند

تصادفی زمان گسسته نامیده می‌شود. در غیر این صورت (به عنوان مثال، $\mathcal{T} = \mathbb{R}$) به آن یک فرآیند تصادفی زمان پیوسته می‌گویند. مدل‌های زیادی در یادگیری ماشینی وجود دارد که با متغیرهای تصادفی مرتبط زمانی سروکار دارند (مانند مدل‌های خودبازگشتی برای سری‌های زمانی، زنجیره‌های مارکوف، مدل‌های مارکوف پنهان، شبکه‌های بیزی پویا). زبان متغیرهای تصادفی، از طریق فرآیندهای تصادفی، به خوبی امکان فرمول‌بندی این مدل‌ها را فراهم می‌کند. با این حال، بیشتر این یادداشته‌ها با تنظیمات ساده‌تر که فقط به متغیرهای تصادفی چند متغیره ($i.i.d.$) نیاز دارند، سروکار دارند.

مثال ۲۲: [سه پرتاب یک سکه منصفانه برای احتمالات مشترک]. دو متغیر تصادفی از مثال ۲۰ را در نظر بگیرید و فضاهای احتمال، توزیع مشترک و حاشیه آنها را محاسبه کنید. $Recall X$ تعداد سرها در اولین پرتاب و Y تعداد سرهای روی هر سه پرتاب است.

تابع جرم احتمال مشترک $p(x, y) = P(X = x, Y = y)$ در زیر نشان داده شده است

		Y			
		0	1	2	3
X	0	$1/8$	$1/4$	$1/8$	0
	1	0	$1/8$	$1/4$	$1/8$

اما اجازه دهید لحظه‌ای به عقب برگردیم و نشان دهیم که چگونه می‌توانیم آن را محاسبه کنیم. اجازه دهید دو مجموعه $A = \{HHT, HTH, THH\}$ و $B = \{HHH, HHT, HTH, HTT\}$ را در نظر بگیریم که به ترتیب مربوط به رویدادهایی است که اولین پرتاب سر بود و دقیقاً دو سر روی سه پرتاب وجود داشت. حال، اجازه دهید به احتمال تقاطع A و B نگاه کنیم

$$\begin{aligned} P(A \cap B) &= P(\{HHT, HTH\}) \\ &= \frac{1}{4} \end{aligned}$$

می‌توانیم احتمال عبارت منطقی $X = 1 \wedge Y = 2$ را به عنوان نمایش دهیم

$$\begin{aligned} p_{XY}(1, 2) &= P(X = 1, Y = 2) \\ &= P(A \cap B) \\ &= P(HHT, HTH) \\ &= \frac{1}{4} \end{aligned}$$

توزیع احتمال حاشیه‌ای را می‌توان به روشی ساده پیدا کرد

$$p(x) = \sum_{y \in \mathcal{Y}} p(x, y)$$

که در آن $\mathcal{Y} = \{0, 1, 2, 3\}$ بدین ترتیب

$$p_X(0) = \sum_{y \in \mathcal{Y}} p_{XY}(0, y)$$

برای پایان خاطرنشان می‌کنیم که در حالت $1 - |\mathcal{X}| \cdot |\mathcal{Y}|$ داریم پارامتر آزاد (زیرا مجموع باید برابر ۱ باشد) تا توزیع مشترک $p(x, y)$ را به طور کامل توصیف کند. به طور مجانبی، این مربوط به رشد تصاعدی تعداد ورودی‌های جدول با تعداد متغیرهای تصادفی (d) است. برای مثال، اگر $|X_i| = 2$ برای $\forall X_i$ ، $d - 1$ عنصر آزاد در توزیع احتمال مشترک وجود دارد. برآورد چنین توزیع‌هایی از داده‌ها غیرقابل حل است و یکی از اشکال لعنت ابعاد ابعادی است.

A.4.2. مثالی از استقلال مشروط

در دو مثال ساده از شکل A.1 نشان می‌دهیم که استقلال و استقلال مشروط بر یکدیگر دلالت ندارند. این مثال بیشتر جبری است و تمرین مفیدی برای نشان دادن این ویژگی است، اما شهود زیادی در مورد اینکه چرا این اتفاق می‌افتد ارائه نمی‌دهد. قوانین جداسازی d ارائه شده در بخش A.6 بیشتر توضیح می‌دهد که چرا می‌توانید این وابستگی‌های مختلف را داشته باشید. ما قبلاً در مثال ۷ مثالی از X و Y ارائه کرده‌ایم که به طور مشروط مستقل هستند، اما مستقل نیستند، در مثال ۷. ما یک مثال دیگر را در اینجا برای زمانی که دو متغیر X و Y مستقل هستند، اما با توجه به Z مستقل مشروط نیستند، ارائه می‌دهیم. این است که اطلاعات در Z زوج X و Y است، در حالی که در مثال ۷، دانستن Z (بایاس سکه) X و Y را جدا می‌کند (دو تلنگر از سکه بایاس).

مثال ۲۳: تنظیماتی را در نظر بگیرید که در آن سعی می‌کنید قیمت یک خانه را پیش بینی کنید. شما نمونه‌های زیادی از قیمت خانه‌های قبلی دارید، اما بدون هیچ ویژگی مرتبط. بگذارید X و Y قیمت دو خانه متفاوت باشد. بدون هیچ گونه اطلاعات اضافی، این دو متغیر مستقل هستند - در واقع، می‌توانیم آنها را به عنوان $i.i.d.$ در نظر بگیریم. نمونه‌هایی از برخی توزیع‌های اساسی بر روی قیمت مسکن. با این حال، اگر اکنون اطلاعات اضافی به ما داده شود که هر دو خانه مشترک هستند، آن‌ها وابسته می‌شوند. بگذارید Z با متغیری مطابقت داشته باشد که دو خانه در یک همسایگی قرار دارند (یعنی یک متغیر 0 یا 1). اگر $Z = 1$ باشد، آنگاه دانستن قیمت X قطعاً بر میزان توزیع بر قیمت‌ها برای Y تأثیر می‌گذارد. اضافه شدن این ویژگی این دو متغیر تصادفی را به صورت شرطی وابسته می‌کند.

A.4.3. اطلاعات اضافی برای انتظارات و لحظات

در این بخش، چند مثال اضافی از انتظارات توابع یک متغیر تصادفی ارائه می‌کنیم که اغلب در نظر گرفته می‌شوند - به اندازه‌ای که نام‌هایی برای آنها داده شود. به یاد بیاورید که برای یک تابع $f: \mathcal{X} \rightarrow \mathbb{R}$ ، مقدار مورد انتظار $\mathbb{E}[f(X)] = \sum f(x)p(x)$ را داریم. با استفاده از $f(x) = x^k$ در لحظه k -ام، $f(x) = \log 1/p(x)$ تابع آنترپی شناخته شده $H(X)$ یا آنترپی دیفرانسیل برای متغیرهای تصادفی پیوسته و $f(x) = (x - \mathbb{E}[X])^2$ واریانس یک متغیر تصادفی X را که با $V[X]$ نشان داده می‌شود، ارائه می‌دهد. جالب توجه است که احتمال وقوع یک رویداد $A \subseteq \mathcal{X}$ نیز می‌تواند به شکل انتظار بیان شود. یعنی

$$P(A) = \mathbb{E}[1(X \in A)]$$

A . X و Y مستقل هستند، اما با توجه به Z مستقل مشروط نیستند

$P(X = 1)$
a

$$P(Y=y|X=x) = P(Y=y)$$

for example,

$$P(Y=1|X=x) = b$$

$$P(Y=1) = b$$

X	$P(Y=1 X)$
0	b
1	b

$$P(Y=y|X=x, Z=z) \neq P(Y=y|Z=z)$$

for example,

$$P(Y=1|X=1, Z=1) = bc / (1 - c - b(1 - 2c))$$

$$P(Y=1|Z=1) = b(1 - c - a(1 - 2c)) / d$$

$$\text{where } d = P(Z=1)$$

X	Y	$P(Z=1 X, Y)$
0	0	c
0	1	$1 - c$
1	0	$1 - c$
1	1	c

B . X و Z با توجه به Y به صورت مشروط مستقل هستند، اما مستقل نیستند

$P(X = 1)$
a

$$P(Z=z|X=x) \neq P(Z=z)$$

for example,

$$P(Z=1|X=1) = d + ce - cd$$

$$P(Z=1) = d + (e - d)(a(c - b) + b)$$

X	$P(Y=1 X)$
0	b
1	c

$$P(Z=z|X=x, Y=y) = P(Z=z|Y=y)$$

for example,

$$P(Z=1|X=x, Y=1) = e$$

$$P(Z=1|Y=1) = e$$

X	Y	$P(Z=1 X, Y)$
0	0	d
0	1	e
1	0	d
1	1	e

شکل A.1: استقلال در مقابل استقلال شرطی با استفاده از توزیع احتمال شامل سه متغیر تصادفی باینری. توزیع‌های احتمال با استفاده از فاکتورسازی $p(x, y, z) = p(x)p(y|x)p(z|x, y)$ ارائه می‌شوند که در آن همه ثابت‌های $a, b, c, d, e \in [0, 1]$. (الف) متغیرهای X و Y مستقل هستند، اما با توجه به Z به صورت شرطی مستقل نیستند. هنگامی که $c = 0$ ، $Z = X \oplus Y$ که در آن \oplus یک عملگر "انحصاری یا" است. (ب) متغیرهای X و Z با توجه به Y به صورت شرطی مستقل هستند، اما مستقل نیستند.

$f(x)$	Symbol	Name
x	$\mathbb{E}[X]$	Mean
$(x - \mathbb{E}[X])^2$	$V[X]$	Variance
x^k	$\mathbb{E}[X^k]$	k-th moment; $k \in \mathbb{N}$
$(x - \mathbb{E}[X])^k$	$\mathbb{E}[(X - \mathbb{E}[X])^k]$	k-th central moment; $k \in \mathbb{N}$
e^{tx}	$M_X(t)$	Moment generating function
e^{itx}	$\varphi_X(t)$	Characteristic function
$\log \frac{1}{p(x)}$	$H(X)$	(Differential) entropy
$\log \frac{p(x)}{q(x)}$	$D(p q)$	Kullback-Leibler divergence
$\left(\frac{\partial}{\partial \theta} \log p(x \theta)\right)^2$	$\mathcal{I}(\theta)$	Fisher information

جدول A.1: برخی از توابع مهم انتظار $\mathbb{E}[f(X)]$ برای متغیر تصادفی X که با توزیع آن $p(x)$ توصیف شده است. تابع $q(x)$ در تعریف واگرایی Kullback – Leibler غیر منفی است و باید مجموع (ادغام) 1 شود. یعنی خود یک توزیع احتمال است. اطلاعات فیشر برای خانواده‌ای از توزیع‌های احتمالی تعریف شده توسط پارامتر θ تعریف شده است. توجه داشته باشید که تابع مولد لحظه ممکن است برای برخی از توزیع‌ها و همه مقادیر t وجود نداشته باشد. با این حال، تابع مشخصه همیشه وجود دارد، حتی زمانی که تابع چگالی وجود ندارد.

که

$$1(t) = \begin{cases} 1 & \text{t is true} \\ 0 & \text{t is false} \end{cases}$$

یک تابع نشانگر است. با این، می‌توان تابع توزیع تجمعی را به صورت $F_X(t) = \mathbb{E}[1(\mathbf{X} \in (-\infty, t])]$ بیان کرد.

تابع $f(x)$ درون انتظار نیز می‌تواند با ارزش مختلط باشد. برای مثال، $\phi_X(t) = \mathbb{E}[e^{itx}]$ ، جایی که i واحد خیالی است، تابع مشخصه X را تعریف می‌کند. استنتاج آماری چندین تابع انتظار در جدول A.1 خلاصه شده است.

مثال ۲۴: [سه پرتاب یک سکه منصفانه (دوباره)]. دو متغیر تصادفی از مثال‌های ۳ و ۵ را در نظر بگیرید و انتظاری و واریانس را برای X و Y محاسبه کنید. سپس $\mathbb{E}[Y | X = 0]$ را محاسبه کنید. ما با محاسبه $\mathbb{E}[X] = 0 \cdot p_X(0) + 1 \cdot p_X(1) = \frac{1}{2}$ شروع می‌کنیم. به همین ترتیب

$$\begin{aligned} \mathbb{E}[Y] &= \sum_{y=0}^3 y \cdot p_Y(y) \\ &= p_Y(1) + 2p_Y(2) + 3p_Y(3) \\ &= \frac{3}{2} \end{aligned}$$

$f(x, y)$	Symbol	Name
$(x - \mathbb{E}[X])(y - \mathbb{E}[Y])$	$\text{Cov}[X, Y]$	Covariance
$\frac{(x - \mathbb{E}[X])(y - \mathbb{E}[Y])}{\sqrt{V[X]V[Y]}}$	$\text{Corr}[X, Y]$	Correlation
$\log \frac{p(x, y)}{p(x)p(y)}$	$I(X; Y)$	Mutual information
$\log \frac{1}{p(x, y)}$	$H(X, Y)$	Joint entropy
$\log \frac{1}{p(x y)}$	$H(X Y)$	Conditional entropy

جدول A.2: برخی از توابع مهم انتظار $\mathbb{E}[f(X, Y)]$ برای دو متغیر تصادفی X و Y که با توزیع مشترک آنها $p(x, y)$ توضیح داده شده است. اطلاعات متقابل گاهی اوقات به عنوان اطلاعات متقابل متوسط نامیده می‌شود.

انتظار مشروط را می‌توان به صورت پیدا کرد

$$\begin{aligned}\mathbb{E}[Y|X=0] &= \sum_{y=0}^3 y \cdot p_{Y|X}(y|0) \\ &= p_{Y|X}(1|0) + 2p_{Y|X}(2|0) + 3p_{Y|X}(3|0) \\ &= 1\end{aligned}$$

که در آن $p(y|x) = p(x, y)/p(x)$

A.5 مخلوط‌های توزیع

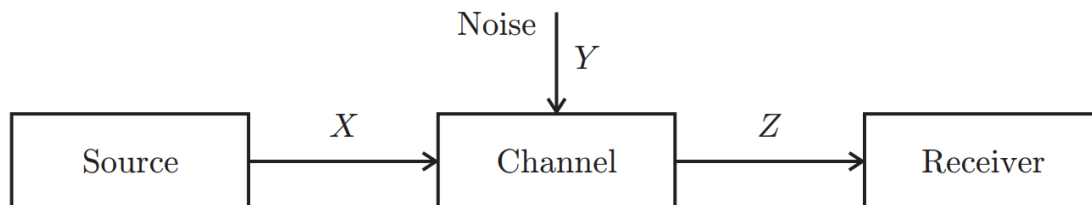
در بخش‌های قبلی دیدیم که متغیرهای تصادفی اغلب با استفاده از خانواده‌های خاصی از توزیع‌های احتمال توصیف می‌شوند. این رویکرد را می‌توان با در نظر گرفتن مخلوطی از توزیع‌ها تعمیم داد. به عنوان مثال، ترکیبات خطی سایر توزیع‌های احتمال. مانند قبل، ما فقط متغیرهای تصادفی را در نظر خواهیم گرفت که توابع جرم یا چگالی خود را دارند

با توجه به مجموعه‌ای از توزیع‌های احتمال m ، $\{p_i(x)\}_{i=1}^m$ ، یک تابع توزیع مخلوط محدود، یا مدل مخلوط، $p(x)$ به صورت تعریف می‌شود.

$$p(x) = \sum_{i=1}^m w_i p_i(x)$$

که در آن $\mathbf{w} = (w_1, w_2, \dots, w_m)$ مجموعه‌ای از اعداد حقیقی غیر منفی است به طوری که $\sum_{i=1}^m w_i = 1$. ما به \mathbf{w} به عنوان ضرایب اختلاط یا گاهی اوقات به عنوان احتمالات اختلاط اشاره می‌کنیم. ترکیب خطی با چنین ضرایبی را ترکیب محدب می‌نامند. راستی آزمایی اینکه تابعی که به این روش تعریف شده است، در واقع یک توزیع احتمال است، ساده است.

در اینجا به طور مختصر به توابع اولیه انتظار از توزیع مخلوط می‌پردازیم. حالت فرضی $\{X_i\}_{i=1}^m$ مجموعه‌ای از m متغیرهای تصادفی است که با توزیع احتمال مربوطه آن‌ها توصیف می‌شوند $\{p_{X_i}(x)\}_{i=1}^m$. همچنین فرض کنید که یک متغیر تصادفی X با توزیع مخلوط با هم کارآمدی w و توزیع احتمال $\{p_{X_i}(x)\}_{i=1}^m$ توصیف شود. سپس با فرض متغیرهای تصادفی پیوسته



$$X: \text{Bernoulli}(\alpha)$$

$$Z = X + Y$$

$$Y: \text{Gaussian}(\mu, \sigma^2)$$

شکل A.2: یک سیستم ارتباطی سیگنال دیجیتال با نویز افزایشی.

در \mathbb{R} تعریف شده است، تابع انتظار به صورت داده شده است

$$\begin{aligned} \mathbb{E}[f(X)] &= \int_{-\infty}^{+\infty} f(x) p_X(x) dx \\ &= \int_{-\infty}^{+\infty} f(x) \sum_{i=1}^m w_i p_{X_i}(x) dx \\ &= \sum_{i=1}^m w_i \int_{-\infty}^{+\infty} f(x) p_{X_i}(x) dx \\ &= \sum_{i=1}^m w_i \mathbb{E}[f(X_i)] \end{aligned}$$

اکنون می‌توانیم این فرمول را برای به دست آوردن میانگین، زمانی که $f(x) = x$ و واریانس، زمانی که $f(x) = (x - E[X])^2$ به دست آوریم، متغیر تصادفی X به عنوان

$$\mathbb{E}[X] = \sum_{i=1}^m w_i \mathbb{E}[X_i]$$

و به این ترتیب

$$V[X] = \sum_{i=1}^m w_i V[X_i] + \sum_{i=1}^m w_i (\mathbb{E}[X_i] - \mathbb{E}[X])^2$$

مثال ۲۵: ارتباطات سیگنالی انتقال یک سیگنال دیجیتال باینری منفرد (بیت) را از طریق یک کانال ارتباطی نویز نشان داده شده در شکل A.2 در نظر بگیرید. بزرگی سیگنال X منتشر شده از منبع به همان اندازه 0 یا 1 ولت است. سیگنال از طریق

یک خط انتقال ارسال می‌شود (به عنوان مثال، ارتباطات رادیویی، فیبر نوری، نوار مغناطیسی) که در آن یک مولفه نویز معمولی توزیع شده صفر به X اضافه می‌شود. توزیع احتمال سیگنال $Z = X + Y$ را که وارد می‌شود استخراج کنید. گیرنده.

ما وضعیت کمی کلی‌تر را در نظر خواهیم گرفت که در آن: $X \sim \text{Bernoulli}(\alpha)$ و Y گاوسی (μ, σ^2) . برای یافتن $p(z)$ از توابع مشخصه متغیرهای تصادفی X, Y و Z استفاده خواهیم کرد که به صورت $\phi_X(t) = \mathbb{E}[e^{itX}]$, $\phi_Y(t) = \mathbb{E}[e^{itY}]$ و $\phi_Z(t) = \mathbb{E}[e^{itZ}]$ نوشته می‌شود. بدون اشتقاق می‌نویسیم

$$\phi_X(t) = 1 - \alpha + \alpha e^{it}$$

$$\phi_Y(t) = e^{it\mu - \frac{\sigma^2 t^2}{2}}$$

و پس از آن

$$\phi_Z(t) = \phi_{X+Y}(t)$$

$$= \phi_X(t) \cdot \phi_Y(t)$$

$$= (1 - \alpha + \alpha e^{it\mu - \frac{\sigma^2 t^2}{2}})$$

$$= \alpha e^{it(\mu+1) - \frac{\sigma^2 t^2}{2}} + (1 - \alpha) e^{it\mu - \frac{\sigma^2 t^2}{2}}$$

با انجام یکپارچه سازی روی $\phi_Z(t)$ می‌توانیم به راحتی آن را تأیید کنیم

$$p(z) = \alpha \cdot \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(z-\mu-1)^2} + (1 - \alpha) \cdot \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(z-\mu)^2}$$

که مخلوطی از دو توزیع نرمال $N(\mu+1, \sigma^2)$ و $N(\mu, \sigma^2)$ با ضرایب $w_1 = \alpha$ و $w_2 = 1 - \alpha$ است. توجه کنید که ترکیب محدب متغیرهای تصادفی $Z = w_1 X + w_2 Y$ به معنای $p_Z(x) = w_1 p_X(x) + w_2 p_Y(x)$ نیست.

A.6 نمایش گرافیکی توزیع‌های احتمال

قبلاً دیدیم که یک توزیع احتمال مشترک را می‌توان با استفاده از قانون زنجیره‌ای از معادله (۱,۲) فاکتور گرفت. چنین فاکتورسازی‌هایی را می‌توان با استفاده از یک نمایش گراف جهت‌دار، که در آن گره‌ها متغیرهای تصادفی را نشان می‌دهند و یال‌ها وابستگی را نشان می‌دهند، تجسم کرد. مثلاً،

$$p(x, y, z) = p(x)p(y|x)p(z|x, y)$$

در شکل A.3A نشان داده شده است. نمایش‌های گرافیکی توزیع‌های احتمال با استفاده از نمودارهای غیر چرخه‌ای جهت‌دار، همراه با توزیع‌های احتمال شرطی، شبکه‌های بیزی یا شبکه‌های باور نامیده می‌شوند. آنها تفسیر و همچنین استنتاج آماری مؤثر را تسهیل می‌کنند.

تجسم روابط بین متغیرها به ویژه زمانی راحت می‌شود که بخواهیم ویژگی‌های استقلال شرطی متغیرها را درک و تحلیل کنیم. شکل A.3B همان فاکتورگیری $p(x, y, z)$ را نشان می‌دهد که در آن متغیر Z مستقل از X با توجه به Y است. با این حال، برای تعیین دقیق ویژگی‌های استقلال شرطی و وابستگی، معمولاً از قوانین جداسازی d برای شبکه‌های اعتقادی استفاده می‌شود. اگرچه روابط اغلب شهودی هستند، گاهی اوقات ویژگی‌های وابستگی به دلیل روابط متعدد بین گره‌ها پیچیده‌تر می‌شوند. به

عنوان مثال، در شکل A.4A، دو گره لبه ندارند، اما به طور مشروط از طریق گره دیگری وابسته هستند. از سوی دیگر، در شکل A.4B، فقدان لبه به معنای استقلال مشروط است. در حال حاضر قوانین جداسازی d را بیشتر بررسی خواهیم کرد. آنها را می‌توان به راحتی در هر کتاب درسی استاندارد در مورد مدل‌های گرافیکی یافت.

شبکه‌های اعتقادی یک تعریف ساده و رسمی دارند. با توجه به مجموعه‌ای از d متغیرهای تصادفی $X = (X_1, \dots, X_d)$ شبکه‌های اعتقادی توزیع احتمال مشترک X را به صورت فاکتور می‌گیرند.

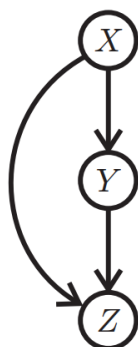
$$p(x) = \prod_{i=1}^d p(x_i | x_{Parents(X_i)})$$

که در آن $Parents(X)$ اجداد مستقیم گره X را در نمودار نشان می‌دهد. در شکل A.3B، گره Y والد Z است، اما گره X والد Z نیست.

ذکر این نکته ضروری است که چندین راه (چند؟) برای فاکتورگیری توزیع وجود دارد. برای مثال، با معکوس کردن ترتیب متغیرهای $p(x, y, z)$ را نیز می‌توان به صورت فاکتور گرفت

$$p(x, y, z) = p(z)p(y|z)p(x|y, z)$$

A. توزیع احتمال گسسته بدون استقلال شرطی



$P(X = 1)$
0.3

X	$P(Y = 1 X)$
0	0.5
1	0.9

X	Y	$P(Z = 1 X, Y)$
0	0	0.3
0	1	0.1
1	0	0.7
1	1	0.4

$$P(X = x, Y = y, Z = z) = P(X = x)P(Y = y|X = x)P(Z = z|X = x, Y = y)$$

B. توزیع احتمال گسسته؛ Z به طور مشروط مستقل از X با Y است



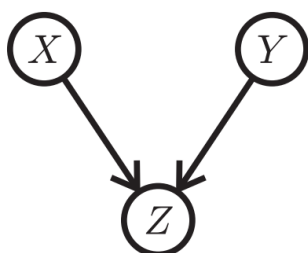
$P(X = 1)$
0.3

Y	$P(Z = 1 Y)$
0	0.2
1	0.7

X	$P(Y = 1 X)$
0	0.5
1	0.9

$$P(X = x, Y = y, Z = z) = P(X = x)P(Y = y|X = x)P(Z = z|Y = y)$$

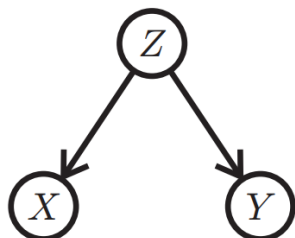
A. X مستقل از Y است، اما Z داده نشده است



$$P(X = x | Y = y) = P(X = x)$$

$$P(X = x | Y = y, Z = z) \neq P(X = x | Z = z)$$

B. X و Y وابسته هستند، اما با توجه به Z به طور مشروط مستقل هستند



$$P(X = x | Y = y, Z = z) = P(X = x | Z = z)$$

$$P(Y = y | X = x, Z = z) = P(Y = y | Z = z)$$

که نمایش گرافیکی متفاوتی دارد و توزیع‌های احتمال شرطی خاص خود را دارد، اما همان توزیع احتمال مشترک با فاکتورگیری قبلی را دارد. انتخاب فاکتورسازی مناسب و تخمین توزیع‌های احتمال شرطی از داده‌ها بعداً به تفصیل مورد بحث قرار خواهد گرفت.

از نمودارهای غیر جهت دار نیز می‌توان برای فاکتورسازی توزیع‌های احتمال استفاده کرد. ایده اصلی در اینجا این است که نمودارها را به دسته‌های حداکثر C (کوچکترین مجموعه‌ای از دسته‌هایی که نمودار را پوشش می‌دهد) تجزیه کنیم و توزیع را به شکل زیر بیان کنیم.

$$p(\mathbf{x}) = \frac{1}{Z} \prod_{C \in \mathcal{C}} \psi_C(\mathbf{x}_C)$$

که در آن هر $\psi_C(\mathbf{x}_C) \geq 0$ تابع پتانسیل دسته و نامیده می‌شود

$$Z = \int_{\mathbf{x}} \prod_{C \in \mathcal{C}} \psi_C(\mathbf{x}_C) d\mathbf{x}$$

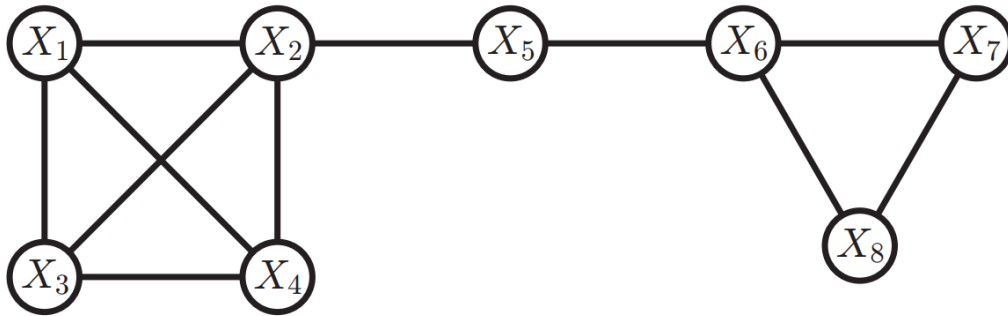
تابع پارتیشن نامیده می‌شود که صرفاً برای اهداف عادی سازی استفاده می‌شود. برخلاف توزیع‌های احتمال شرطی در نمودارهای غیر چرخه‌ای جهت‌دار، پتانسیل‌های دسته معمولاً تفسیر احتمال شرطی ندارند و بنابراین، نرمال‌سازی ضروری است. یک نمونه از تجزیه دسته حداکثر در شکل A.5 نشان داده شده است.

توابع بالقوه معمولاً کاملاً مثبت در نظر گرفته می‌شوند، و به این صورت بیان می‌شوند.

$$\psi_C(\mathbf{x}_C) = \exp(-E(\mathbf{x}_C))$$

که در آن $E(\mathbf{x}_C)$ یک تابع انرژی مشخص شده توسط کاربر بر روی دسته متغیرهای تصادفی \mathbf{X}_C است. این منجر به توزیع احتمال شکل زیر می‌شود

$$p(\mathbf{x}) = \frac{1}{Z} \exp\left(\sum_{C \in \mathcal{C}} \log \psi_C(\mathbf{x}_C)\right)$$



$$\mathbf{X}_{C_1} = \{X_1, X_2, X_3, X_4\}$$

$$\mathbf{X}_{C_3} = \{X_5, X_6\}$$

$$\mathbf{X}_{C_2} = \{X_2, X_5\}$$

$$\mathbf{X}_{C_4} = \{X_6, X_7, X_8\}$$

همانطور که فرمول بندی شد، این توزیع احتمال، توزیع بولتزمن یا توزیع گیبس نامیده می‌شود.

تابع انرژی $E(\mathbf{x})$ باید برای مقادیر \mathbf{x} که احتمال بیشتری دارند کمتر باشد. همچنین ممکن است شامل پارامترهایی باشد که سپس از داده‌های آموزشی موجود تخمین زده می‌شود. البته، در یک مسئله پیش‌بینی، یک گراف بدون جهت باید ایجاد شود تا متغیرهای هدف را که در اینجا زیرمجموعه‌ای از \mathbf{X} در نظر گرفته می‌شوند، نیز در بر گیرد.

اکنون هر توزیع احتمال را بر روی تمام پیکربندی‌های ممکن بردار تصادفی \mathbf{X} با نمایش گرافیکی زیربنایی آن در نظر بگیرید. اگر اموال زیر

$$p(\mathbf{x}_i | \mathbf{x}_{-\mathbf{x}_i}) = p(\mathbf{x}_i | \mathbf{x}_{N(\mathbf{x}_i)}) \quad (A.4)$$

راضی است، توزیع احتمال به عنوان شبکه مارکوف یا میدان تصادفی مارکوف نامیده می‌شود. در معادله بالا

$$\mathbf{X}_{-\mathbf{x}_i} = (X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_d)$$

و $N(X)$ مجموعه‌ای از متغیرهای تصادفی مجاور X در نمودار است. یعنی یک لبه بین X و هر گره در $N(X)$ وجود دارد. به مجموعه متغیرهای تصادفی در $N(X)$ پتوی مارکوف X نیز گفته می‌شود.

می‌توان نشان داد که هر توزیع گیبس ویژگی معادله (A.4) را برآورده می‌کند و برعکس، برای هر توزیع احتمالی که معادله (A.4) برای آن وجود دارد، می‌توان به عنوان یک توزیع گیبس با برخی از پارامترها نشان داد. این هم ارزی توزیع‌های گیبس و شبکه‌های مارکوف توسط قضیه هامرسلی-کلیفورد ایجاد شد.

پیوست B

پس زمینه بهینه سازی

B.1 قوانین اساسی برای گرادیان

برای مشتقات، قوانین مفیدی وجود دارد که با آنها آشنا هستید، مانند $\frac{d}{dw} e^w = e^w$ ، $\frac{d}{dw} w^2 = 2w$ ، $\frac{d}{dw} aw = a$ ما می توانیم به طور مشابه چنین قوانینی را برای تنظیم چند متغیره بنویسیم تا محاسبه گرادیان ها را بدون نیاز به محاسبه هر مشتق جزئی ساده کنیم. هر یک از قوانین زیر را می توان با محاسبه مشتقات جزئی، با قوانینی که برای حالت تک متغیره به آن ها عادت دارید، تأیید کرد. ما قوانین کلیدی این سند را در اینجا خلاصه می کنیم. برای مرجع کامل تر، به کتاب آشپزی ماتریسی [۱۶] مراجعه کنید.

برخی از این قوانین در جدول B.1 خلاصه شده است. این فهرست جامع نیست، اما برخی از قوانین اضافی را امکان پذیر می سازد. به عنوان مثال، برای به دست آوردن مشتق تابع $f(x) = x^T A$ ، ابتدا می توان مشتق $f(x)^T = A^T x$ را بدست آورد و سپس جابجایی آن را گرفت زیرا

$$(\nabla f(x))^T = \nabla(f(x)^T)$$

بنابراین، چون برای $\nabla f(x)^T = A^T$ ما $\nabla f(x) = A$ را دریافت می کنیم

$f(x)$	$\frac{\partial f}{\partial x}$
$x^T x$	$2x$
Ax	A
$x^T Ax$	$Ax + A^T x$

جدول B.1: فرمول های مشتق مفید بردار ها با توجه به بردار ها. مشتق تابع با مقدار برداری $f: \mathbb{R}^{d \times 1} \rightarrow \mathbb{R}^{m \times 1}$ با توجه به بردار $x \in \mathbb{R}^{d \times 1}$ یک ماتریس $m \times d$ به اندازه M با مولفه های $M_{ij} = \partial y_j / \partial x_i$ ، $i \in \{1, 2, \dots, d\}$ و $j \in \{1, 2, \dots, m\}$ است. یک مشتق از اسکالر با توجه به یک بردار، که در آن $m = 1$ ، یک مورد خاص از این وضعیت است که منجر به بردار ستونی $d \times 1$ می شود. توجه داشته باشید که در جدول، m برای هر ردیف یکسان نیست. به عنوان مثال، $f(x) = x^T x$ یک اسکالر است، در حالی که برای یک ماتریس عمومی $A \in \mathbb{R}^{m \times d}$ ، $f(x) = Ax$ یک بردار m بعدی است.

پیوست C

پس زمینه جبر خطی

1. دیدگاه جبری

ابزار قدرتمند دیگری برای تجزیه و تحلیل و درک رگرسیون خطی از جبر خطی خطی و کاربردی می آید. در این بخش ما مسیری انحرافی می کنیم تا به مبانی جبر خطی بپردازیم و سپس این مفاهیم را برای عمیق تر کردن درک خود از رگرسیون به کار ببریم. در جبر خطی، ما اغلب علاقه مند به حل مجموعه معادلات زیر هستیم که در زیر به صورت ماتریسی آورده شده است.

$$Ax = b \quad (C.1)$$

در اینجا، A یک ماتریس $m \times n$ ، b یک بردار $m \times 1$ و x یک بردار $n \times 1$ است که باید پیدا شود. تمام عناصر A ، x و b به عنوان اعداد واقعی در نظر گرفته می شوند. ما با یک سناریوی ساده شروع می کنیم و فرض می کنیم A یک مربع، ماتریس 2×2 است. این مجموعه از معادلات را می توان به صورت بیان کرد

$$a_{11}x_1 + a_{12}x_2 = b_1$$

$$a_{21}x_1 + a_{22}x_2 = b_2$$

به عنوان مثال، ممکن است ما علاقه مند به حل باشیم

$$x_1 + 2x_2 = 3$$

$$x_1 + 3x_2 = 5$$

این یک فرمول مناسب زمانی است که ما می خواهیم سیستم را حل کنیم، به عنوان مثال. با حذف گاوسی با این حال، فرمول مناسبی برای درک مسئله وجود راه حل ها نیست. برای اینکه بتوانیم این کار را انجام دهیم، مفاهیم اساسی جبر خطی را به اختصار مرور می کنیم.

1.1. چهار زیرفضای اساسی

هدف این بخش بررسی مختصر چهار زیرفضای اساسی در جبر خطی (فضای ستون، فضای ردیف، خالی، فضای خالی سمت چپ) و رابطه متقابل آنهاست. ما با مثال خود از بالا شروع می کنیم و سیستم معادلات خطی را به صورت زیر می نویسیم

$$\begin{bmatrix} 1 \\ 1 \end{bmatrix} x_1 + \begin{bmatrix} 2 \\ 3 \end{bmatrix} x_2 = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$

اکنون می بینیم که با حل $Ax = b$ به دنبال مقادیر مناسب بردارهای $(1, 1)$ و $(2, 3)$ هستیم تا ترکیب خطی آنها $(3, 5)$ را ایجاد کند. این مقادیر $x_1 = -1$ و $x_2 = 2$ هستند. اجازه دهید $a_1 = (1, 1)$ و $a_2 = (2, 3)$ را به عنوان بردارهای ستون A تعریف کنیم. یعنی $A = [a_1 a_2]$. بنابراین، هر زمان که b را بتوان به صورت ترکیب خطی از بردارهای ستون a_1 و a_2 بیان کرد، $Ax = b$ قابل حل خواهد بود.

تمام ترکیبات خطی ستون های ماتریس A فضای ستون A ، $C(A)$ را با بردارهای a_1, \dots, a_n تشکیل می دهند اساس این فضا است. هر دو b و $C(A)$ در فضای m بعدی \mathbb{R}^m قرار دارند. بنابراین، آنچه $Ax = b$ می گوید این است که b باید در فضای ستون A قرار داشته باشد تا معادله جواب داشته باشد. در مثال بالا، اگر ستون های A به طور خطی مستقل باشند، راه حل منحصر به فرد است، یعنی تنها یک ترکیب خطی از بردارهای ستون وجود دارد که b را نشان می دهد. در غیر این صورت، چون A یک ماتریس مربع است، سیستم هیچ راه حلی ندارد. نمونه ای از چنین وضعیتی است

$$\begin{bmatrix} 1 & 2 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$

که در آن $a_1 = (1, 1)$ و $a_2 = (2, 2)$. در اینجا، a_1 و a_2 (خطی) وابسته هستند زیرا $2a_1 - a_2 = 0$. بین فضاهای ایجاد شده توسط مجموعه ای از بردارها و ویژگی های ماتریس A ارتباط عمیقی وجود دارد. در حال حاضر، با استفاده از مثال بالا، کافی است بگوییم که اگر a_1 و a_2 مستقل باشند، ماتریس A غیر مفرد است (تفرد را فقط برای ماتریس های مربع می توان مورد بحث قرار داد)، که دارای رتبه کامل است.

به شیوه ای معادل با فضای ستون، تمام ترکیبات خطی ردیف های A فضای ردیف را تشکیل می دهند که با $C(A^T)$ نشان داده می شود، که در آن هر دو x و $C(A^T)$ در \mathbb{R}^n هستند. تمام راه حل های $Ax = 0$ ، فضای خالی ماتریس، $N(A)$ را تشکیل می دهند، در حالی که تمام راه حل های $A^T y = 0$ فضای خالی سمت چپ A ، $N(A^T)$ را تشکیل می دهند. واضح است که $C(A)$ و $N(A^T)$ در \mathbb{R}^m تعبیه شده اند، در حالی که $C(A^T)$ و $N(A)$ در \mathbb{R}^n هستند. با این حال، جفت فضاهای فرعی متعامد هستند (بردارهای u و v متعامد هستند اگر $u^T v = 0$). یعنی هر بردار در $C(A)$ با تمام بردارهای $N(A^T)$ متعامد است و هر بردار در $C(A^T)$ نسبت به همه بردارهای $N(A)$ متعامد است. این به راحتی قابل مشاهده است: اگر $x \in N(A)$ ، پس طبق تعریف $Ax = 0$ ، و بنابراین هر ردیف از A متعامد به x است. اگر هر سطر بر x متعامد باشد، تمام ترکیبات خطی سطرها نیز متعامد هستند.

متعامد بودن ویژگی کلیدی چهار زیرفضا است، زیرا تجزیه مفید بردارهای x و b را از معادله فراهم می کند. (C.1) با توجه به A (ما در بخش بعدی از آن استفاده خواهیم کرد). برای مثال، هر $x \in \mathbb{R}^n$ را می توان به صورت تجزیه کرد

که در آن $x_r \in C(A^T)$ و $x_n \in N(A)$ ، به طوری که $\|x\|_2^2 = \|x_r\|_2^2 + \|x_n\|_2^2$. به طور مشابه، هر $b \in \mathbb{R}^m$ را می توان به صورت تجزیه کرد

$$b = b_c + b_l$$

$$\|b\|_2^2 = \|b_c\|_2^2 + \|b_l\|_2^2 \text{ و } b_l \in N(A^T), b_c \in C(A)$$

در بالا اشاره کردیم که ویژگی های فضاهای بنیادی با ویژگی های ماتریس A ارتباط تنگاتنگی دارند. برای نتیجه گیری این بخش، اجازه دهید به طور خلاصه در مورد رتبه یک ماتریس و رابطه آن با ابعاد زیرفضاهای بنیادی بحث کنیم. اساس فضا کوچکترین مجموعه بردارهایی است که فضا را در بر می گیرد (این مجموعه بردارها منحصر به فرد نیستند). اندازه پایه را بعد فضا نیز می گویند. در مثال ابتدای این بخش، یک فضای ستون دو بعدی با بردارهای پایه $a_1 = (1, 1)$ و $a_2 = (2, 3)$ داشتیم. از سوی دیگر،

برای $a_1 = (1,1)$ و $a_2 = (2,3)$ یک فضای ستون یک بعدی، یعنی یک خط، به طور کامل توسط هر یک از بردارهای پایه تعیین می شود. جای تعجب نیست که بعد فضای پوشیده شده توسط بردارهای ستون برابر با رتبه ماتریس A است. یکی از نتایج اساسی در جبر خطی این است که رتبه A با بعد $C(A)$ یکسان است که به نوبه خود با بعد $C(A^T)$ بعد.

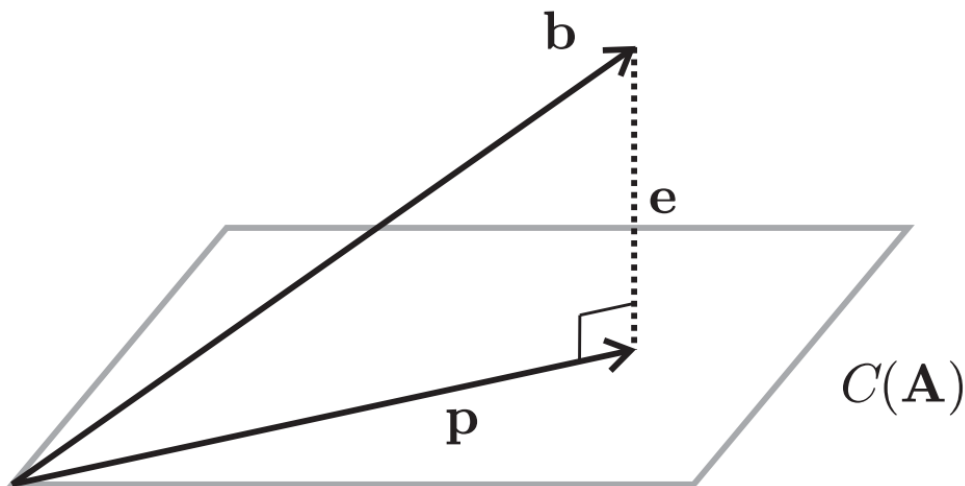
1.2.C به حداقل رساندن $\|Ax - b\|_2^2$

اجازه دهید اکنون دوباره به راه حل های $Ax = b$ نگاه کنیم. به طور کلی، سه نتیجه متفاوت وجود دارد:

۱. هیچ راه حلی برای سیستم وجود ندارد

۲. یک راه حل منحصر به فرد برای سیستم وجود دارد، و

۳. بی نهایت راه حل وجود دارد.



نکته 1.C: تصویر طرح ریزی بردار b به فضای ستون ماتریس A . بردارهای $p(b_c)$ و $e(b_l)$ به ترتیب نقطه طرح و خطا را نشان می دهند.

این نتایج به رابطه بین رتبه (r) و مقدار A و ابعاد m و n بستگی دارد. ما قبلاً می دانیم که وقتی $r = m = n$ (ماتریس مربع، معکوس، رتبه کامل A) یک راه حل منحصر به فرد برای سیستم وجود دارد، اما اجازه دهید موقعیت های دیگر را بررسی کنیم. به طور کلی، هنگامی که $r = n < m$ (رتبه ستون کامل)، سیستم یا یک راه حل دارد یا هیچ راه حلی ندارد، همانطور که به صورت لحظه ای خواهیم دید. وقتی $r = m < n$ (رتبه ردیف کامل)، سیستم بی نهایت راه حل دارد. در نهایت، در مواردی که $r < n$ و $r < m$ ، یا هیچ راه حلی وجود ندارد یا بی نهایت راه حل وجود دارد. از آنجا که $Ax = b$ ممکن است قابل حل نباشد، ما حل $Ax = b$ را به مینیمم کردن $\|Ax - b\|_2$ تعمیم می دهیم. به این ترتیب می توان همه موقعیت ها را در یک چارچوب یکپارچه در نظر گرفت.

اجازه دهید مثال زیر را در نظر بگیریم

$$\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix}, \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

که نمونه ای را نشان می دهد که بعید است راه حلی برای $\mathbf{Ax} = \mathbf{b}$ داشته باشیم، مگر اینکه محدودیتی در b_1, b_2 و b_3 وجود داشته باشد. در اینجا، محدودیت $b_3 = 2b_2 - b_1$ است. در این وضعیت، $C(\mathbf{A})$ یک صفحه دو بعدی در \mathbb{R}^3 است که توسط بردارهای ستون $\mathbf{a}_1 = (1,1,1)$ و $\mathbf{a}_2 = (2,3,4)$ پوشیده شده است. اگر محدودیت عناصر \mathbf{b} برآورده نشود، هدف ما این است که سعی کنیم نقطه ای را در $C(\mathbf{A})$ پیدا کنیم که نزدیکترین نقطه به \mathbf{b} باشد. همانطور که در شکل C.1 نشان داده شده است، این نقطه ای است که \mathbf{b} به $C(\mathbf{A})$ پیش بینی می شود. ما به طرح ریزی \mathbf{b} به $C(\mathbf{A})$ به صورت \mathbf{p} اشاره خواهیم کرد. حال با استفاده از نماد جبری استاندارد معادلات زیر را داریم

$$\mathbf{b} = \mathbf{p} + \mathbf{e}$$

$$\mathbf{p} = \mathbf{Ax}$$

از آنجایی که \mathbf{p} و \mathbf{e} متعامد هستند، می دانیم که $\mathbf{p}^T \mathbf{e} = 0$. اجازه دهید \mathbf{x} را حل کنیم

$$(\mathbf{Ax})^T (\mathbf{b} - \mathbf{Ax}) = 0$$

$$\mathbf{x}^T \mathbf{A}^T \mathbf{b} - \mathbf{x}^T \mathbf{A}^T \mathbf{Ax} = 0$$

$$\mathbf{x}^T (\mathbf{A}^T \mathbf{b} - \mathbf{A}^T \mathbf{Ax}) = 0$$

و بنابراین

$$\mathbf{x}^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$$

این دقیقاً همان راه حلی است که مجموع مجذور خطاها را به حداقل رساند و احتمال را به حداکثر رساند. ماتریکس

$$\mathbf{A}^\dagger = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$$

شبه معکوس مور-پنروز یا به سادگی یک شبه معکوس نامیده می شود. این یک ماتریس مهم است زیرا همیشه وجود دارد و منحصر به فرد است، حتی در شرایطی که معکوس $\mathbf{A}^T \mathbf{A}$ وجود ندارد.

این زمانی اتفاق می افتد که \mathbf{A} دارای ستون های وابسته باشد (از لحاظ فنی، \mathbf{A} و $\mathbf{A}^T \mathbf{A}$ فضای خالی یکسانی دارند که حاوی بیش از مبدأ سیستم مختصات است؛ بنابراین رتبه $\mathbf{A}^T \mathbf{A}$ کمتر از n است). اجازه دهید برای لحظه ای به بردار طرح ریزی \mathbf{p} نگاه کنیم. ما داریم

$$\mathbf{p} = \mathbf{Ax}$$

$$= \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$$

که در آن $\mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ ماتریسی است که \mathbf{b} را به فضای ستون \mathbf{A} می فرستد

در حالی که ما به همان نتیجه ای که در بخش های قبلی بود رسیدیم، ابزار جبر خطی به ما اجازه می دهد تا رگرسیون OLS را در سطح عمیق تری مورد بحث قرار دهیم. اجازه دهید لحظه ای وجود و تعدد راه حل ها را بررسی کنیم

$$\arg \min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2 \quad (C.2)$$

واضح است که راه حل این مشکل همیشه وجود دارد. با این حال، اکنون خواهیم دید که راه حل این مشکل به طور کلی منحصر به فرد نیست و به رتبه A بستگی دارد. x را یک راه حل برای معادله در نظر بگیرید. (C.2). به یاد بیاورید که $x = x_r + x_n$ و در A ضرب می شود. بنابراین، هر بردار $x = x_r + \alpha x_n$ که $\alpha \in \mathbb{R}$ ، نیز یک راه حل است. توجه داشته باشید که x_r در همه این راه حل ها مشترک است. اگر نمی توانید آن را ببینید، فرض کنید بردار دیگری از فضای ردیف وجود دارد و نشان دهید که امکان پذیر نیست. اگر ستون های A مستقل باشند، راه حل منحصر به فرد است زیرا فضای خالی فقط مبدا را در بر می گیرد. در غیر این صورت، بی نهایت راه حل وجود دارد. در چنین مواردی، با فرافکنی b به $C(A)$ دقیقاً چه راه حلی پیدا می شود؟ اجازه دهید به آن نگاه کنیم:

$$\begin{aligned} x^* &= A^\dagger b \\ &= (A^T A)^{-1} A^T (p + e) \\ &= (A^T A)^{-1} A^T p \\ &= x_r \end{aligned}$$

به عنوان $p = Ax_r$ با توجه به اینکه x_r منحصر به فرد است، راه حلی که با بهینه سازی مینیمم مربعات یافت می شود، راه حلی است که به طور همزمان $\|Ax - b\|_2$ و $\|x\|_2$ را به حداقل می رساند (مشاهده کنید که $\|x\|_2$ به حداقل می رسد زیرا راه حل هر جزء از فضای خالی را نادیده می گیرد). بنابراین، مشکل رگرسیون OLS گاهی اوقات به عنوان مسئله مینیمم نرم مینیمم مربع نامیده می شود

حال بیا بید موقعیت هایی را در نظر بگیریم که در آن $Ax = b$ راه حل های بی نهایت زیادی دارد. یعنی وقتی $b \in C(A)$. این معمولاً زمانی رخ می دهد که $n > m \geq r$ باشد. در اینجا، به دلیل اینکه b قبلاً در فضای ستون A قرار دارد، تنها سؤال این است که با روش کمینه سازی چه راه حل خاصی پیدا می شود. همانطور که در بالا دیدیم، نتیجه فرآیند کمینه سازی راه حلی با مینیمم نرم L_2 $\|x\|_2$ است.

به طور خلاصه، اجازه دهید ابتدا به نماد اصلی خود بازگردیم که در آن X ماتریس و w وزن هایی هستند که باید پیدا شوند. هدف از مسئله رگرسیون OLS حل $Xw = y$ است، اگر قابل حل باشد. وقتی $d < n$ این یک سناریوی واقع بینانه در عمل نیست. بنابراین، ما نیاز را کاهش دادیم و سعی کردیم نقطه ای را در فضای ستون $C(X)$ که نزدیکترین به y است را پیدا کنیم. معلوم شد که این معادل به حداقل رساندن مجموع خطاهای مربع (یا فاصله اقلیدسی) بین بردارهای n بعدی Xw و y است. همچنین مشخص شد که معادل راه حل حداکثر احتمال ارائه شده در بخش ۵.۱ است. وقتی $d < n$ ، یک وضعیت معمول در عمل این است که بی نهایت راه حل وجود دارد. در این شرایط، الگوریتم بهینه سازی ما الگوریتمی را با مینیمم نرم L_2 پیدا می کند.

پیوست D

جزئیات در مورد رویکردهای نمایندگی بدون نظارت با استفاده از فاکتورسازی

انواع مختلفی از الگوریتم های یادگیری بدون نظارت وجود دارد که در واقع با فاکتورسازی ماتریس داده ها مطابقت دارند که در جدول D.1 خلاصه می کنیم. در بسیاری از موارد، آنها به سادگی با تعریف یک هسته جالب در \mathbf{X} ، و سپس فاکتورگیری آن هسته (یعنی هسته PCA) مطابقت دارند. اگر ورودی خالی باشد، هیچ قاعده سازی و هیچ محدودیتی را مشخص نمی کند. برای فهرست کامل تر، [۱۸] و [۲۰] را ببینید. همانند تنظیمات رگرسیون، می توانیم این اتلاف اقلیدسی را به هر زیان محدب تعمیم دهیم

$$L_x(\mathbf{H}, \mathbf{D}, \mathbf{X}) = \sum_{i=1}^n L_x(\mathbf{H}_i; \mathbf{D}, \mathbf{X}_i;)$$

جایی که در بالا استفاده کردیم

$$L_x(\mathbf{H}, \mathbf{D}, \mathbf{X}) = \sum_{i=1}^n ||\mathbf{H}_i \mathbf{D} - \mathbf{X}_i|| = ||\mathbf{H} \mathbf{D} - \mathbf{X}||_F^2$$

الگوریتم هایی برای یادگیری لغت نامه ها

تمرکز ما بر پیش بینی است، و بنابراین مایلیم از این نمایش ها برای یادگیری تحت نظارت (یا نیمه نظارت) استفاده کنیم. ما از یک رویکرد دو مرحله ای استفاده می کنیم، که در آن ابتدا نمایش جدید به روشی بدون نظارت یاد می شود و سپس با الگوریتم های یادگیری نظارت شده استفاده می شود. این دو مرحله را می توان با یادگیری فرهنگ لغت تحت نظارت در یک مرحله ترکیب کرد. برای بحث در مورد این رویکرد پیشرفته تر به [۱۳] مراجعه کنید.

رایج ترین استراتژی برای یادگیری این مدل های فرهنگ لغت، انجام یک کمینه سازی متناوب بر روی متغیرها است. بهینه سازی روی \mathbf{D} و \mathbf{H} به طور مشترک محدب نیست. با این حال، در هر متغیر به طور جداگانه محدب است. استراتژی این است که یک متغیر را ثابت کنیم، مثلاً \mathbf{H} ، و در دیگری نزول کنیم، مثلاً \mathbf{D} ، و سپس سوئیچ کنیم، \mathbf{D} را ثابت کنیم و در \mathbf{H} نزولی کنیم. این کمینه سازی متناوب تا زمان همگرایی ادامه می یابد. اگرچه این یک بهینه سازی غیر محدب است، شواهد اخیری وجود دارد که نشان می دهد این روش در واقع مینیمم مطلق را برمی گرداند (به عنوان مثال [۱۲] مراجعه کنید). ما این روش را در الگوریتم ۵ خلاصه می کنیم.

هنگامی که فرهنگ لغت \mathbf{D} و نمایش جدید \mathbf{H} یاد گرفتیم، می توانیم وزن های نظارت شده $\mathbf{W} \in \mathbb{R}^{k \times m}$ را برای به دست آوردن $\mathbf{Y} \approx \mathbf{H}\mathbf{W}$ یاد بگیریم. این را می توان با هر یک از روش های رگرسیون خطی یا طبقه بندی که تاکنون آموخته ایم انجام داد.

در نهایت، ما باید بدانیم که چگونه از این مدل های آموخته شده برای پیش بینی خارج از نمونه (یعنی برای نمونه های جدید) استفاده کنیم. ماتریس های \mathbf{D} و \mathbf{W} حاوی تمام اطلاعات لازم برای انجام پیش بینی خارج از نمونه هستند و \mathbf{H} نیازی به ذخیره سازی ندارد، زیرا این نمایش مختص داده های آموزشی بود. برای نمونه جدید \mathbf{x}_{new} ، نمایش را می توان با استفاده از آن به دست آورد

$$\mathbf{h}_{\text{new}} = \underset{\mathbf{h} \in \mathbb{R}^k}{\operatorname{argmin}} L_{\mathbf{x}}(\mathbf{h}\mathbf{D}, \mathbf{x})$$

با نمایش این نمونه، می توانیم $f(h_{\text{new}}\mathbf{W})$ را پیش بینی کنیم.

Algorithm	Loss and constraints
CCA \equiv orthonormal PLS	$\left\ \begin{bmatrix} \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \\ \mathbf{Y}(\mathbf{Y}^\top \mathbf{Y})^{-1} \end{bmatrix} - \mathbf{H}\mathbf{D} \right\ _F^2$
Isomap	$\left\ \mathbf{K} - \mathbf{H}\mathbf{D} \right\ _F^2$ $\mathbf{K} = -\frac{1}{2}(\mathbf{I} - \mathbf{e}\mathbf{e}')\mathbf{S}(\mathbf{I} - \mathbf{e}\mathbf{e}')$ with $\mathbf{S}_{i,j} = \ \mathbf{X}_i - \mathbf{X}_j\ $
K-means clustering	$\left\ \mathbf{X} - \mathbf{H}\mathbf{D} \right\ _F^2$ with $\mathbf{H} \in \{0, 1\}^{n \times k}$, $\mathbf{H}\mathbf{1} = \mathbf{1}$
K-medians clustering	$\left\ \mathbf{X} - \mathbf{H}\mathbf{D} \right\ _{1,1}$ with $\mathbf{H} \in \{0, 1\}^{n \times k}$, $\mathbf{H}\mathbf{1} = \mathbf{1}$
Laplacian eigenmaps \equiv Kernel LPP	$\left\ \mathbf{K} - \mathbf{H}\mathbf{D} \right\ _F^2$ for $\mathbf{K} = \mathbf{L}^\dagger$
Metric multi-dimensional scaling	$\left\ \mathbf{K} - \mathbf{H}\mathbf{D} \right\ _F^2$ for isotropic kernel \mathbf{K}
Normalized-cut	$\left\ (\mathbf{\Lambda}^{-1}\mathbf{X} - \mathbf{H}\mathbf{D})\mathbf{\Lambda}^{1/2} \right\ _F^2$ with $\mathbf{H} \in \{0, 1\}^{n \times k}$, $\mathbf{H}\mathbf{1} = \mathbf{1}$
Partial least squares	$\left\ \mathbf{X}\mathbf{Y}' - \mathbf{D}\mathbf{H} \right\ _F^2$
PCA	$\left\ \mathbf{X} - \mathbf{H}\mathbf{D} \right\ _F^2$
Kernel PCA	$\left\ \mathbf{K} - \mathbf{H}\mathbf{D} \right\ _F^2$
Ratio cut	$\left\ \mathbf{K} - \mathbf{H}\mathbf{D} \right\ _F^2$ for $\mathbf{K} = \mathbf{L}^\dagger$

شکل D.1. الگوریتم های یادگیری بدون نظارت که با فاکتورسازی ماتریسی مطابقت دارند.

Algorithm 5: Alternating minimization for dictionary learning

Input:

inner dimension k
loss L , where $L(\mathbf{H}\mathbf{D}) = L_x(\mathbf{H}, \mathbf{D}, \mathbf{X})$
 R_D , the regularizer on \mathbf{D}
 R_H , the regularizer on \mathbf{H}
the regularization weight λ
convergence tolerance
fixed positive step-sizes η_D, η_H
dataset x_1, \dots, x_n

Initialization:

$\mathbf{D}, \mathbf{H} \leftarrow$ full-rank random matrices with inner dimension k
prevobj $\leftarrow \infty$

Loop until convergence within tolerance or reach maximum number of iterations:

Update \mathbf{D} using one step of gradient descent
Update \mathbf{H} using one step of gradient descent
currentobj $\leftarrow L(\mathbf{H}\mathbf{D}) + \lambda R_D(\mathbf{D}) + \frac{\lambda}{n} R_H(\mathbf{H})$
If $|\text{currentobj} - \text{prevobj}| < \text{tolerance}$, Then break
prevobj $\leftarrow \text{currentobj}$

Output:

\mathbf{D}, \mathbf{H}

پیوست E

تخمین بیزی

رویکردهای حداکثر پسین و حداکثر درستنمایی راه حلی را گزارش می کنند که به ترتیب با حالت توزیع پسین و تابع درستنمایی مطابقت دارد. با این حال، این رویکرد امکان توزیع های اریب، توزیع های چندوجهی یا صرفاً مناطق بزرگ با مقادیر مشابه $p(f|D)$ را در نظر نمی گیرد. تخمین بیزی به این نگرانی ها می پردازد.

ایده اصلی در آمار بیزی، به حداقل رساندن ریسک پسین است

$$R = \int_{\mathcal{F}} \ell(f, \hat{f}) \cdot p(f|D) df$$

که در آن \hat{f} تخمین ما و $\ell(f, \hat{f})$ مقدار تابع هزینه بین دو مدل است. هنگامی که $\ell(f, \hat{f}) = (f - \hat{f})^2$ (سوء استفاده از علامت گذاری را نادیده بگیرید)، می توانیم خطر بعدی را به صورت زیر به حداقل برسانیم.

$$\frac{\partial}{\partial \hat{f}} R = 2\hat{f} - 2 \int_{\mathcal{F}} f \cdot p(f|D) df = 0$$

که از آن می توان به دست آورد که مینیمم کردن ریسک پسین، تابع میانگین پسین است. یعنی

$$f_B = \int_{\mathcal{F}} f \cdot p(f|D) df = \mathbb{E}[F|D]$$

که در آن F یک متغیر تصادفی است که مدل را نشان می دهد. ما باید به f_B به عنوان تخمینگر بیز اشاره کنیم. ذکر این نکته ضروری است که محاسبه میانگین پسین معمولاً شامل حل انتگرال های پیچیده است. در برخی شرایط، این انتگرال ها را می توان به صورت تحلیلی حل کرد. در برخی دیگر، ادغام عددی ضروری است.

مثال ۲۶: اجازه دهید $\mathcal{D} = \{2, 5, 9, 5, 4, 8\}$ دوباره یک $i.i.d$ باشد. نمونه از پواسون (λ_0) فرض کنید دانش قبلی در مورد پارامتر توزیع را می توان با استفاده از توزیع گاما با پارامترهای $k = 3$ و $\theta = 1$ بیان کرد. تخمین بیزی λ_0 را بیابید.

می خواهیم $\mathbb{E}[\lambda|D]$ را پیدا کنیم. اجازه دهید ابتدا توزیع پسین را به صورت بنویسیم

$$\begin{aligned} p(\lambda|D) &= \frac{p(D|\lambda)p(\lambda)}{p(D)} \\ &= \propto \{ p(D|\lambda)p(\lambda) \} \int_0^\infty p(D|\lambda)p(\lambda) d\lambda \end{aligned}$$

که در آن، همانطور که در مثال های قبلی نشان داده شد، آن را داریم

$$p(\mathcal{D}|\lambda) = \frac{\lambda^{\sum_{i=1}^n x_i} \cdot e^{-n\lambda}}{\prod_{i=1}^n x_i!}$$

و

$$p(\lambda) = \frac{\lambda^{k-1} e^{-\frac{\lambda}{\theta}}}{\theta^k \Gamma(k)}$$

قبل از محاسبه $p(\mathcal{D})$ ، ابتدا به این نکته توجه کنیم

$$\int_0^{\infty} x^{\alpha-1} e^{-\beta x} dx = \frac{\Gamma(\alpha)}{\beta^{\alpha}}$$

اکنون، ما می توانیم آن را استخراج کنیم

$$\begin{aligned} p(\mathcal{D}) &= \int_0^{\infty} p(\mathcal{D}|\lambda) p(\lambda) d\lambda \\ &= \int_0^{\infty} \frac{\lambda^{\sum_{i=1}^n x_i} \cdot e^{-n\lambda}}{\prod_{i=1}^n x_i!} \cdot \frac{\lambda^{k-1} e^{-\frac{\lambda}{\theta}}}{\theta^k \Gamma(k)} d\lambda \\ &= \frac{\Gamma(k + \sum_{i=1}^n x_i)}{\theta^k \Gamma(k) \prod_{i=1}^n x_i! \left(n + \frac{1}{\theta}\right)^{\sum_{i=1}^n x_i + k}} \end{aligned}$$

و متعاقباً آن

$$\begin{aligned} p(\lambda|\mathcal{D}) &= \frac{p(\mathcal{D}|\lambda) p(\lambda)}{p(\mathcal{D})} \\ &= \frac{\lambda^{\sum_{i=1}^n x_i} \cdot e^{-n\lambda}}{\prod_{i=1}^n x_i!} \cdot \frac{\lambda^{k-1} e^{-\frac{\lambda}{\theta}}}{\theta^k \Gamma(k)} \cdot \frac{\theta^k \Gamma(k) \prod_{i=1}^n x_i! (n + \frac{1}{\theta})^{\sum_{i=1}^n x_i + k}}{\Gamma(k + \sum_{i=1}^n x_i)} \\ &= \frac{\lambda^{k-1 + \sum_{i=1}^n x_i} \cdot e^{-\lambda(n + \frac{1}{\theta})}}{\Gamma(k + \sum_{i=1}^n x_i)} \cdot \left(n + \frac{1}{\theta}\right)^{\sum_{i=1}^n x_i + k} \end{aligned}$$

سرانجام،

$$\begin{aligned} \mathbb{E}[\lambda|\mathcal{D}] &= \int_0^{\infty} \lambda p(\lambda|\mathcal{D}) d\lambda \\ &= \frac{k + \sum_{i=1}^n x_i}{n + \frac{1}{\theta}} = 5.14 \end{aligned}$$

که تقریباً همان راه حلی است که برآورد MAP در مثال ۹ یافت شد.

از مثال قبلی مشهود است که انتخاب توزیع قبلی پیامدهای مهمی در محاسبه میانگین پسین دارد. ما توزیع گاما را تصادفی انتخاب نکرده ایم. یعنی وقتی احتمال در قبلی ضرب شد، توزیع حاصل در همان کلاس توابع قبلی باقی می ماند. ما به چنین توزیع های قبلی به عنوان پیشین های مزدوج اشاره خواهیم کرد. پیشین های مزدوج نیز ریاضیات را ساده می کنند. در واقع، این دلیل اصلی توجه آنهاست. جالب اینجاست که علاوه بر توزیع پواسون، توزیع گاما مزدوج قبل از توزیع نمایی و همچنین خود توزیع گاما است.