



دانشگاه صنعتی شریف

دانشکده مهندسی کامپیوتر

یادگیری ماشین

پاییز و زمستان ۱۴۰۴

استاد: علی شریفی زارچی

گردآورندگان: ارشیا یوسف نیا - بنیامین قنبری - زهرا رحمانی - علی باوفا - کیهان هدایی

تمرین اول

یادگیری با نظارت

مهلت ارسال: ۳۰ آبان

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.
- در طول ترم، برای هر تمرین می‌توانید تا ۵ روز تأخیر مجاز داشته باشید و در مجموع حداکثر ۱۵ روز تأخیر مجاز خواهید داشت. توجه داشته باشید که تأخیر در تمرین‌های عملی و تئوری به صورت جداگانه محاسبه می‌شود و مجموع تأخیر هر دو نباید بیشتر از ۱۵ روز شود. پس از اتمام زمان مجاز، دو روز اضافی برای آپلود غیرمجاز در نظر گرفته شده است که در این بازه، به ازای هر ساعت تأخیر، ۲ درصد از نمره تمرین کسر خواهد شد.
- همکاری و هم‌فکری شما در انجام تمرین مانعی ندارد اما پاسخ‌های هر کس حتماً باید توسط خود او نوشته شده باشد.
- در صورت هم‌فکری و یا استفاده از هر منابع خارج درسی، نام هم‌فکران و آدرس منابع مورد استفاده برای حل سوال مورد نظر را ذکر کنید.
- لطفاً تصویری واضح از پاسخ‌های سوالات نظری بارگذاری کنید. در غیر این صورت پاسخ شما تصحیح نخواهد شد.

سوالات نظری (۱۰۰ نمره)

۱. (۱۰ نمره)

- (آ) اگر داده‌ها خطی قابل تفکیک نباشند، هیچ راه‌حلی برای SVM به صورت soft-margin نداریم.
- (ب) برای یک مجموعه داده که نویز زیادی دارد، یعنی تعداد زیادی از داده‌ها برچسب نادرست دارند، استفاده از جنگل تصادفی به طور کلی بهتر از boosted decision trees است.
- (پ) در روش bagging در استفاده از درخت تصمیم‌گیری، معمولاً n درخت با واریانس زیاد و بایاس کمی داریم. انتخاب نهایی بر اساس اجماع آن‌ها با یکدیگر است که حاصل آن، نتیجه‌ای با بایاس کم و واریانس کم می‌شود.
- (ت) در تخمین واریانس-بایاس، استفاده از یک مدل بسیار پیچیده بر روی داده محدود، معمولاً در ابتدا به بایاس کم و واریانس بالا منجر می‌شود.
- (ث) الگوریتم perceptron تضمین می‌کند که برای هر مجموعه داده‌ای که به صورت خطی قابل تفکیک باشد، پس از تعداد محدودی گام، به یک راه‌حل همگرا می‌شود.
- حل.

(آ) نادرست. اگر داده‌ها به صورت خطی جدایی‌ناپذیر نباشند، می‌توان با استفاده از حقه kernel و بردن داده‌ها به ابعاد بالاتر، آن‌ها را به صورت کامل از هم تفکیک کرد.

(ب) درست. جنگل تصادفی درختان زیادی را بر اساس نمونه‌های bootstrap آموزش می‌دهد و میانگین پیش‌بینی‌های آن‌ها را محاسبه می‌کند. این عمل باعث کاهش اثر خطاهای برچسب‌ها در بین درختان می‌شود، بنابراین این مجموعه در برابر نویز مقاوم است. در مقابل، روش‌های boosting وزن بیشتری به نمونه‌هایی که اشتباه پیش‌بینی شده‌اند اختصاص می‌دهند. برچسب‌های نویزی اغلب اشتباه پیش‌بینی می‌شوند و چون این مجموعه داده حاوی تعداد زیادی از این نوع برچسب‌ها است، باعث می‌شود مدل روی آن‌ها تمرکز کند و آن‌ها را بیش‌برازش کند.

- (پ) درست. در bagging معمولاً هر درخت واریانس زیادی دارد، اما چون نتیجه نهایی با اجماع به دست می‌آید، واریانس خروجی کاهش می‌یابد و تخمین پایدارتر می‌شود.
- (ت) درست. استفاده از یک مدل با پیچیدگی بالا برای یک مجموعه داده محدود، منجر به کد شدن نویز به جای الگوی اصلی داده می‌شود که از آن با عنوان واریانس بالا نیز یاد می‌شود.
- (ث) درست. در صورتی که داده‌ها به صورت خطی قابل تفکیک باشند، الگوریتم perceptron قادر است در تعداد محدودی گام، آن خط جداساز را پیدا کند.

۲. (۱۵ نمره) روش تخمین چگالی در knn برآوردی از چگالی یک نقطه از فضای ویژگی‌ها را نسبت به چگالی k همسایه نزدیک‌تر می‌دهد. نشان دهید که مدل چگالی knn را می‌توان یک توزیع نامناسب دید. این موضوع چه اهمیتی دارد و چه نتیجه‌ای می‌دهد؟

توزیع نامناسب (improper distribution) توزیعی است که انتگرال‌گیری روی آن روی کل بازه واگرا خواهد شد.^۱

راهنمایی: برای محاسبه $\mathbb{P}(x_i)$ ، کره‌ای کوچک با مرکزیت خود آن نقطه در نظر بگیرید و شعاع را تا دربرگرفتن k نقطه افزایش دهید. افزون بر این، کره را چنان کوچک فرض کنید که $\mathbb{P}(x_i)$ در آن به تقریب ثابت باشد. حل. هرگاه در kNN بخواهیم چگالی احتمالی را در نقطه‌ای مانند x_i محاسبه کنیم، کره‌ای کوچک در نظر می‌گیریم که نقطه مورد نظر در مرکز آن قرار داشته باشد و شعاع آن را تا جایی افزایش می‌دهیم که این کره k نقطه را در بر گیرد. در این صورت خواهیم داشت:

$$p(x_i) = \frac{k}{NV_i}$$

که N تعداد تمام مشاهدات و V_i حجم کره‌ای به مرکزیت x_i است. فرض می‌کنیم V_i آن قدر کوچک است که $p(x_i)$ در آن ثابت باقی بماند. اکنون خواهیم داشت:

$$\int p(x) dx \approx \sum_{i=1}^N p(x_i) V_i = \sum_{i=1}^N \frac{k}{NV_i} V_i = k \neq 1$$

معادله بالا برقرار خواهد بود اگر و تنها اگر حجم تمامی کره‌ها به اندازه کافی کوچک بوده و N به اندازه کافی بزرگ باشد. هر دو شرط گفته شده نیز در kNN برقرارند. از آنجایی که kNN محدودیتی بر روی توزیع داده اعمال نمی‌کند و تنها بر روی یک بخش محلی متمرکز می‌شود، همان‌گونه که مشاهده کردید می‌بینیم که تخمین چگالی از قوانین توزیع‌های احتمالاتی پیروی نمی‌کند. با توجه به این که انتگرال چگالی تخمین زده شده بر روی کل بازه برابر با یک نمی‌شود، این توزیع نامناسب خواهد بود. به همین علت تفسیر و به کارگیری kNN برای مواردی که در آن‌ها به توزیع‌های احتمالاتی نیاز داریم، دشوار خواهد بود.

۳. (۱۵ نمره) فرض کنید یک تابع رگرسیون $h(x)$ داریم که برای هر بردار ورودی x برچسب $y = h(x)$ را به آن نسبت می‌دهد. از دادگان آموزش استفاده کرده‌ایم و مدل پیش‌بینی‌کننده آموزش داده‌ایم که آن‌ها را با $\hat{h}_1(x), \hat{h}_2(x), \dots, \hat{h}_m(x)$ نشان می‌دهیم. اکنون قرار دهید:

$$\hat{H}_m(x) = \frac{1}{m} \sum_{i=1}^m \hat{h}_i(x)$$

^۱ برای مطالعه بیشتر، مراجعه کنید به: <https://arxiv.org/abs/1711.02064>

ثابت کنید که:

$$\mathbb{E}_{\mathbf{x}} \left[(\hat{H}_m(\mathbf{x}) - h(\mathbf{x}))^2 \right] \leq \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\mathbf{x}} \left[(\hat{h}_i(\mathbf{x}) - h(\mathbf{x}))^2 \right]$$

مفهوم این حکم چیست؟ حل. مقادیر E_{avg} و E_{com} را به صورت زیر تعریف می‌کنیم:

$$E_{com} = \mathbb{E}_x \left[(H_M(x) - h(x))^2 \right], \quad E_{avg} = \frac{1}{M} \sum_{i=1}^M \mathbb{E}_x \left[(y_i(x) - h(x))^2 \right].$$

در این صورت، با تعریف $\epsilon_i(x) = y_i(x) - h(x)$ ، برای E_{com} خواهیم داشت:

$$\begin{aligned} E_{com} &= \mathbb{E}_x \left[(H_M(x) - h(x))^2 \right] \\ &= \mathbb{E}_x \left[\left(\frac{1}{M} \sum_{i=1}^M y_i(x) - h(x) \right)^2 \right] \\ &= \mathbb{E}_x \left[\left(\frac{1}{M} \sum_{i=1}^M y_i(x) - \frac{M}{M} h(x) \right)^2 \right] \\ &= \mathbb{E}_x \left[\left(\frac{1}{M} \sum_{i=1}^M (y_i(x) - h(x)) \right)^2 \right] \\ &= \mathbb{E}_x \left[\left(\frac{1}{M} \sum_{i=1}^M \epsilon_i(x) \right)^2 \right] \\ &= \frac{1}{M^2} \mathbb{E}_x \left[\left(\sum_{i=1}^M \epsilon_i(x) \right)^2 \right]. \end{aligned}$$

همچنین، با همان تعریف $\epsilon_i(x) = y_i(x) - h(x)$ ، برای E_{avg} داریم:

$$E_{avg} = \frac{1}{M} \sum_{i=1}^M \mathbb{E}_x \left[(y_i(x) - h(x))^2 \right] = \frac{1}{M} \sum_{i=1}^M \mathbb{E}_x \left[(\epsilon_i(x))^2 \right].$$

حال با استفاده از نامساوی کوشی-شوارتز به ازای هر x داریم:

$$\left(\sum_{i=1}^M \epsilon_i(x) \right)^2 \leq M \sum_{i=1}^M (\epsilon_i(x))^2.$$

طرفین را در $\frac{1}{M^2}$ ضرب می‌کنیم:

$$\frac{1}{M^2} \left(\sum_{i=1}^M \epsilon_i(x) \right)^2 \leq \frac{1}{M} \sum_{i=1}^M (\epsilon_i(x))^2.$$

اکنون از امیدریاضی نسبت به x می‌گیریم:

$$\mathbb{E}_x \left[\frac{1}{M^2} \left(\sum_{i=1}^M \epsilon_i(x) \right)^2 \right] \leq \mathbb{E}_x \left[\frac{1}{M} \sum_{i=1}^M (\epsilon_i(x))^2 \right].$$

بنابراین،

$$\frac{1}{M^2} \mathbb{E}_x \left[\left(\sum_{i=1}^M \epsilon_i(x) \right)^2 \right] \leq \frac{1}{M} \sum_{i=1}^M \mathbb{E}_x \left[(\epsilon_i(x))^2 \right],$$

و در نتیجه

$$E_{com} \leq E_{avg}.$$

۴. (۱۰ نمره) به سؤالاتی که در ادامه آمده‌اند با بیان توضیحات کامل پاسخ دهید.

(آ) یک مدل Random Forest روی تعداد زیادی داده برای یک مسئله دسته‌بندی آموزش داده‌ایم. میزان خطای مدل روی دادگان آموزشی پایین است ولی روی دادگان آزمون خطای زیادی گزارش شده است. دو دلیل برای این مشکل بیان کنید و برای هر کدام راه‌حلی بگویید.

(ب) فرض کنید دادگان $X = \{x_i\}_{i=1}^n$ با میانگین و واریانس μ, σ^2 داده شده‌اند. بگیرید:

$$z_m = \frac{1}{m} \left(\sum_{i=1}^m a_i \right)$$

که در آن a_i ها تصادفی و متمایز از X گزینش شده‌اند. میانگین و واریانس z_m را بیابید. اکنون بیان کنید مجموعه دادگان $\{z_m\}$ چه خوبی و بدیهایی نسبت به X دارد و در چه شرایطی استفاده از آن برای آموزش مناسب است.

حل.

(آ) مورد اول: پیچیدگی بالای مدل جنگل تصادفی. این یعنی مدل ما به داده‌ها overfit می‌شود و خاصیت تعمیم‌پذیری ندارد.

راه‌حل: محدود کردن عمق و تعداد درخت‌ها تا واریانس کاهش یابد و افزایش حداقل تعداد نمونه لازم برای جدا کردن گره‌ها. این کارها باعث می‌شود مدل محدودیت بیشتری پیدا کند و از بیش‌برازش جلوگیری شود.

مورد دوم: وجود تعداد زیادی ویژگی نامناسب.

راه‌حل: استفاده از روش‌هایی مانند PCA و سایر روش‌های کاهش بُعد برای استخراج تعدادی ویژگی مفید و مهم و ادامه کار تنها بر روی آن‌ها.

مورد سوم: تعداد داده‌ها کم است یا کل دامنه رفتار را پوشش نمی‌دهند.

راه‌حل: اضافه کردن داده‌ها از سراسر دامنه ممکن (گردآوری داده بیشتر و متنوع‌تر).

(ب) فرض کنیم داده‌ها از مجموعه‌ی $X = \{x_i\}_{i=1}^n$ با میانگین μ و واریانس σ^2 انتخاب شده‌اند. آماره‌ی z_m به صورت زیر تعریف می‌شود:

$$z_m = \frac{1}{m} \sum_{i=1}^m a_i,$$

که در آن a_i نمونه‌هایی تصادفی و متمایز (نمونه‌گیری بدون جایگذاری) از X هستند. میانگین آماره‌ی z_m برابر با میانگین جامعه μ است، زیرا z_m یک برآوردگر نااریب برای μ است:

$$\mathbb{E}[z_m] = \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m a_i \right] = \frac{1}{m} \sum_{i=1}^m \mathbb{E}[a_i] = \frac{1}{m} \sum_{i=1}^m \mu = \mu.$$

واریانس آماری z_m در نمونه‌گیری تصادفی ساده بدون جایگذاری از یک جامعه‌ی متناهی با اندازه‌ی n از رابطه‌ی زیر به دست می‌آید:

$$\text{Var}(z_m) = \frac{\sigma^2}{m} \left(\frac{n-m}{n-1} \right),$$

که در آن عامل $\frac{n-m}{n-1}$ به عنوان «عامل تصحیح جامعه‌ی متناهی» (FPC) شناخته می‌شود.

مزایا و معایب مجموعه‌ی داده‌های $\{z_m\}$ نسبت به X :

- **واریانس کاهش‌یافته (دقت بالاتر):** واریانس $\text{Var}(z_m)$ به دلیل تقسیم بر اندازه‌ی نمونه m و اعمال FPC بسیار کمتر از واریانس σ^2 داده‌های اصلی X است؛ بنابراین z_m برآوردگر دقیق‌تری برای μ محسوب می‌شود.
- **نارایی:** چون $\mathbb{E}[z_m] = \mu$ ، به طور متوسط این آماره دقیقاً برابر با مقدار واقعی میانگین جامعه است.
- **از دست دادن جزئیات:** z_m تنها میانگین نمونه را نشان می‌دهد و اطلاعات مربوط به پراکندگی و ویژگی‌های منحصر به فرد (مانند داده‌های پرت) تک‌تک نقاط در X را نادیده می‌گیرد.

در مجموع، آماره‌های $\{z_m\}$ که میانگین نمونه‌های m -تایی هستند، به دلیل ناریب بودن و واریانس پایین (که با افزایش m کاهش می‌یابد)، برآوردگرهای بسیار خوبی برای میانگین جامعه μ به شمار می‌آیند. این ویژگی آن‌ها را برای استفاده در فرایندهای آموزشی مانند Mini-Batch SGD مناسب می‌سازد، زیرا استفاده از آن‌ها نویز محاسباتی (واریانس گرادیان‌ها) را کاهش داده و به همگرایی سریع‌تر و پایدارتر بهینه‌ساز کمک می‌کند؛ مشروط بر این‌که اندازه‌ی m به قدر کافی بزرگ انتخاب شود تا نمونه نماینده‌ی خوبی از جامعه باشد.

۵. (۱۰ نمره) به سؤالاتی که در ادامه آمده‌اند با بیان توضیحات کامل پاسخ دهید.

(آ) با توجه به تابع سیگموید به سؤالات زیر پاسخ دهید.

- مشتق تابع سیگموید را بدست آورید (تمامی مراحل به طور کامل نوشته شود).
- تابع سیگموید و مشتق آن را رسم کنید. از نمودار بدست آمده کمک بگیرید و به طور شهودی توضیح دهید که مشکل تابع سیگموید در بروز رسانی وزن‌ها با کاهش گرادیان چیست.
- (ب) چرا رگرسیون لجستیک که با الگوریتم کاهش گرادیان آموزش داده می‌شود، بعضی وقت‌ها که دادگان تقریباً به طور خطی قابل تفکیک هستند، کند همگرا می‌شود؟
- (ج) در یک مسئله طبقه‌بندی چندکلاسه با K کلاس، به جای استفاده از چندین مدل دودویی، می‌توان از یک مدل چندکلاسه استفاده کرد. نشان دهید که اگر به همه بردارهای پارامتر $\theta_1, \theta_2, \dots, \theta_K$ یک بردار ثابت c اضافه کنیم، مقادیر احتمالات تغییری نمی‌کنند.

حل.

(آ) (بخش اول)

$$y = \sigma(x) = \frac{1}{1 + e^{-x}}$$

$$u = 1 + e^{-x}, \quad y = \frac{1}{u}$$

$$\frac{du}{dx} = -e^{-x}, \quad \frac{dy}{du} = -\frac{1}{u^2}$$

با استفاده از قاعده زنجیره‌ای داریم:

$$\frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx} = \left(-\frac{1}{u^2}\right) (-e^{-x}) = \frac{e^{-x}}{u^2} = \frac{e^{-x}}{(1 + e^{-x})^2}.$$

از طرفی:

$$\sigma(x) = \frac{1}{1 + e^{-x}}, \quad 1 - \sigma(x) = \frac{e^{-x}}{1 + e^{-x}},$$

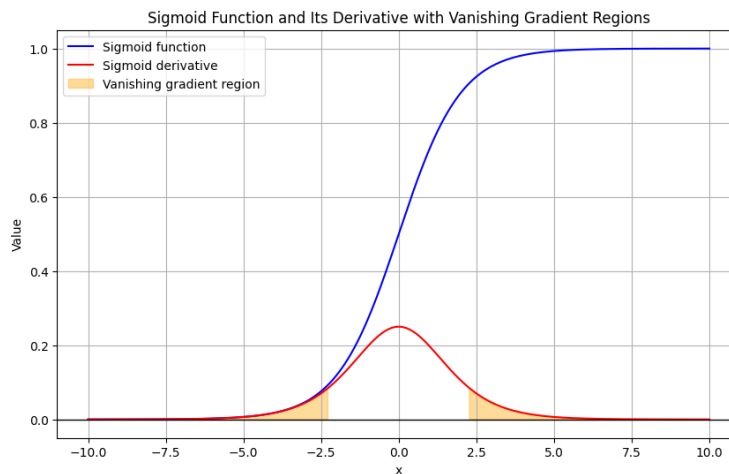
پس می‌توان نوشت:

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \sigma(x)(1 - \sigma(x)).$$

در نتیجه:

$$\sigma'(x) = \sigma(x)(1 - \sigma(x)).$$

(بخش دوم)



ناحیه نارنجی‌رنگ نشان‌دهنده بخش‌هایی است که مشتق $\sigma'(x)$ بسیار کوچک (نزدیک به صفر) است. در این نواحی، گرادیان‌هایی که برای به‌روزرسانی وزن‌ها و بایاس‌ها استفاده می‌شوند بسیار کوچک می‌شوند، در نتیجه مدل بسیار آهسته یاد می‌گیرد یا حتی یادگیری تقریباً متوقف می‌شود.

(ب) وقتی داده‌ها تقریباً به صورت خطی قابل تفکیک باشند، لجستیک رگرسیون سعی می‌کند پارامترهایی پیدا کند که دو کلاس را به طور کامل از هم جدا کنند. برای رسیدن به این هدف، مقدار وزن‌ها به‌طور مداوم بزرگ‌تر می‌شود تا احتمال‌های پیش‌بینی‌شده برای نمونه‌های هر کلاس تا حد ممکن به ۰ یا ۱ نزدیک شوند.

اما با بزرگ شدن وزن‌ها، سیگموئید تابع در نواحی اشباع خود قرار می‌گیرد؛ یعنی خروجی آن حتی با تغییرات بزرگ در وزن‌ها، تغییر اندکی می‌کند. در این وضعیت، گرادیان‌ها بسیار کوچک می‌شوند و هر گام الگوریتم (Gradient Descent) گرادیان کاهش فقط یک به‌روزرسانی بسیار کوچک در پارامترها ایجاد می‌کند.

به همین دلیل، فرآیند آموزش ظاهراً «گیر می‌کند» و بسیار کند پیش می‌رود، در حالی که در واقع مدل به آرامی در جهت وزن‌های بسیار بزرگ حرکت می‌کند. افزودن رگولاریزیشن (Regularization) رگولاریزیشن با محدود کردن رشد وزن‌ها، از این رفتار جلوگیری کرده و به مدل کمک می‌کند سریع‌تر و پایدارتر همگرا شود.

(ج) اگر برای همه کلاس‌ها $\theta'_j = \theta_j + c$ (برای یک بردار ثابت c) داشته باشیم، آنگاه:

$$e^{\theta'_j T x} = e^{(\theta_j + c)^T x} = e^{c^T x} e^{\theta_j^T x}.$$

در نتیجه در مدل Softmax داریم:

$$P'(y = k) = \frac{e^{\theta'_k T x}}{\sum_j e^{\theta'_j T x}} = \frac{e^{c^T x} e^{\theta_k^T x}}{e^{c^T x} \sum_j e^{\theta_j^T x}} = \frac{e^{\theta_k^T x}}{\sum_j e^{\theta_j^T x}} = P(y = k).$$

بنابراین اضافه کردن یک بردار ثابت یکسان به همه θ_j ها، احتمال‌های خروجی را تغییر نمی‌دهد و تنها یک افزونگی در پارامترها ایجاد می‌کند؛ یعنی مدل از نظر آماری معادل باقی می‌ماند.

۶. (۱۵ نمره) به سؤالاتی که در ادامه آمده‌اند با بیان توضیحات کامل پاسخ دهید.

(آ) برای یک مجموعه binary classification به صورت (\mathbf{x}_i, y_i) که $y_i \in \{+1, -1\}$ مسئله hard-margin^۲ را در نظر بگیرید که به صورت مسئله بهینه‌سازی

$$\min \left\{ \frac{1}{2} \|\mathbf{w}\|^2 \right\} \quad s.t. \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

تعریف می‌شود.

• توضیح دهید چرا Geometric Margin مربوط به این classifier برابر با

$$\gamma = \min_i \frac{y_i(\mathbf{w}^T \mathbf{x}_i + b)}{\|\mathbf{w}\|}$$

است.

• چرا با کاهش دادن $\frac{1}{2} \|\mathbf{w}\|^2$ با توجه به شروط داده شده، این margin افزایش می‌یابد؟

(ب) فرض کنید مدل Linear Regression را داریم:

$$\mathbf{y} = \mathbf{X}\beta + \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

• فرم بسته‌ای برای تخمین‌گر least squares برای پارامتر β به دست آورید.

• نشان دهید که این تخمین‌گر unbiased است و ماتریس covariance آن را نیز به دست آورید.

حل.

(آ) بخش اول

می‌دانیم Geometric Margin برای نقطه‌ی i برابر با فاصله‌ی آن نقطه از ابرصفحه‌ی classifier است. با توجه به این که مقدار $w^T x + b$ می‌تواند مثبت یا منفی باشد و $y_i \in \{-1, +1\}$ ، به صورت زیر عمل می‌کنیم:

• نقطه‌ی x را روی ابرصفحه در نظر می‌گیریم:

$$w^T x + b = 0.$$

• بردار از x به x_i برابر است با:

$$x_i - x_{..}$$

^۲https://en.wikipedia.org/wiki/Support_vector_machine

- برای به دست آوردن فاصله از ابرصفحه، تصویر این بردار را روی بردار واحد $\frac{w}{\|w\|}$ محاسبه می‌کنیم:

$$\gamma_i = \left| \left\langle \frac{w}{\|w\|}, x_i - x_* \right\rangle \right| = \left| \frac{w^T (x_i - x_*)}{\|w\|} \right|.$$

- با توجه به این که $w^T x_* + b = 0$ داریم:

$$\gamma_i = \left| \frac{w^T x_i + b}{\|w\|} \right|.$$

از آنجا که $w^T x_i + b$ هم علامت با y_i است، می‌توان نوشت:

$$\gamma_i = \frac{y_i (w^T x_i + b)}{\|w\|}.$$

در نهایت، Geometric Margin برای یک مجموعه‌ی داده به صورت کمینه‌ی margin‌های نقاط تعریف می‌شود:

$$\gamma = \min_i \left\{ \frac{y_i (w^T x_i + b)}{\|w\|} \right\}.$$

شرط مسئله‌ی بهینه‌سازی داده شده برابر است با:

$$y_i (w^T x_i + b) \geq 1.$$

بنابراین:

$$\forall i : y_i (w^T x_i + b) \geq 1 \implies \min_i \{y_i (w^T x_i + b)\} \geq 1,$$

و در نتیجه:

$$\gamma = \frac{1}{\|w\|} \min_i \{y_i (w^T x_i + b)\} \geq \frac{1}{\|w\|}.$$

بخش دوم

هدف مسئله‌ی بهینه‌سازی مینیمم کردن $\frac{1}{4} \|w\|^2$ است. این کار معادل با کوچک کردن $\|w\|$ و در نتیجه بزرگ کردن $\frac{1}{\|w\|}$ است. با توجه به این که $\gamma \geq \frac{1}{\|w\|}$ ، کمینه کردن $\|w\|$ در حضور قیدهای

$$y_i (w^T x_i + b) \geq 1$$

منجر به حداکثر شدن Geometric Margin روی مجموعه داده می‌شود. پس با حل این مسئله‌ی بهینه‌سازی، ابرصفحه‌ای را می‌یابیم که Geometric Margin آن حداکثر است.

(ب) بخش اول

ابتدا تابع loss (تابع هزینه) را به صورت زیر در نظر می‌گیریم:

$$J(\beta) = (y - X\beta)^T (y - X\beta).$$

آن را بسط می‌دهیم:

$$J(\beta) = y^T y - 2\beta^T X^T y + \beta^T X^T X \beta.$$

اکنون نسبت به β مشتق می‌گیریم:

$$\frac{\partial J}{\partial \beta} = -2X^T y + 2X^T X \beta.$$

شرط بهینگی $\frac{\partial J}{\partial \beta} = 0$ را اعمال می‌کنیم:

$$-2X^T y + 2X^T X \hat{\beta} = 0 \implies X^T X \hat{\beta} = X^T y \implies \hat{\beta} = (X^T X)^{-1} X^T y.$$

بخش دوم

برای نشان دادن این که $\hat{\beta}$ یک برآوردگر unbiased است، باید نشان دهیم:

$$\mathbb{E}[\hat{\beta}] = \beta.$$

مدل را به صورت $y = X\beta + \epsilon$ می‌نویسیم، که در آن $\mathbb{E}[\epsilon] = 0$ و $\text{Var}(\epsilon) = \sigma^2 I$ است. آنگاه:

$$\hat{\beta} = (X^T X)^{-1} X^T y = (X^T X)^{-1} X^T (X\beta + \epsilon) = (X^T X)^{-1} X^T X \beta + (X^T X)^{-1} X^T \epsilon.$$

از امیدریاضی می‌گیریم:

$$\mathbb{E}[\hat{\beta}] = \mathbb{E}[(X^T X)^{-1} X^T X \beta] + \mathbb{E}[(X^T X)^{-1} X^T \epsilon] = \beta + (X^T X)^{-1} X^T \mathbb{E}[\epsilon] = \beta.$$

پس:

$$\mathbb{E}[\hat{\beta}] = \beta \implies \hat{\beta} \text{ unbiased است.}$$

واریانس $\hat{\beta}$ نیز به صورت زیر به دست می‌آید:

$$\text{Var}(\hat{\beta}) = \text{Var}((X^T X)^{-1} X^T \epsilon) = (X^T X)^{-1} X^T \text{Var}(\epsilon) X (X^T X)^{-1}.$$

با توجه به $\text{Var}(\epsilon) = \sigma^2 I$ داریم:

$$\text{Var}(\hat{\beta}) = (X^T X)^{-1} X^T (\sigma^2 I) X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1}.$$

۷. (۱۰ نمره) همان‌طور که می‌دانید اعتبارسنجی متقاطع^۳ با حذف یک نمونه یا LOOCV را می‌توان بدین گونه نشان داد:

$$e(S_n) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}_{-i}(\mathbf{x}_i))^2$$

که در آن $S_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ (مجموعه داده‌های آموزش)، $S_n^{-i} = S_n \setminus \{(\mathbf{x}_i, y_i)\}$ (مجموعه داده‌ها پس از حذف نمونه i ام)، و \hat{f}_{-i} مدل آموزش دیده روی S_n^{-i} می‌باشند.

(آ) با فرض اینکه فقط نمونه اول حذف شده است، خطا به شکل

$$e_1(S_n) = (y_1 - \hat{f}_{-1}(\mathbf{x}_1))^2$$

است. نشان دهید:

$$\mathbb{E}[e_1(S_n)] = \mathbb{E}[(y - \hat{f}_{S_{n-1}}(\mathbf{x}))^2]$$

که در آن (\mathbf{x}, y) یک نمونه تصادفی از توزیع داده‌هاست و $\hat{f}_{S_{n-1}}$ مدل آموزش دیده روی S_{n-1} است.

(ب) با استفاده از نتیجه بالا، نشان دهید که:

$$\mathbb{E}[e(S_n)] = \mathbb{E}[(y - \hat{f}_{S_{n-1}}(\mathbf{x}))^2]$$

^۳[https://en.wikipedia.org/wiki/Cross-validation_\(statistics\)](https://en.wikipedia.org/wiki/Cross-validation_(statistics))

حل.

(آ) اگر A و B متغیرهای تصادفی (RVs) با توزیع احتمال یکسان باشند، آنگاه $\mathbb{E}[f(A)] = \mathbb{E}[f(B)]$ خواهد بود. این امر را می‌توان با نوشتن انتگرال متناظر، روشن‌تر بیان کرد:

$$\mathbb{E}[f(A)] = \int f(x) p_A(x) dx = \int f(x) p_B(x) dx = \mathbb{E}[f(B)].$$

این حکم زمانی که A و B خود مجموعه‌ای از متغیرهای تصادفی باشند نیز برقرار است. به‌ویژه، اگر

$$A = \{S_n^{-1}, (x_1, y_1)\}$$

باشد، به این معنا که روی S_n^{-1} آموزش می‌دهیم و روی (x_1, y_1) آزمون می‌گیریم، و این مجموعه همان توزیع

$$B = \{S_{n-1}, (x, y)\}$$

را داشته باشد، که در آن روی مجموعه‌ی دیگری از $n-1$ نمونه به نام S_{n-1} آموزش می‌دهیم و روی (x, y) آزمون می‌گیریم (و همه‌ی این‌ها از یک توزیع زیربنایی مشترک نمونه‌برداری شده‌اند)، آنگاه همچنان خواهیم داشت:

$$\mathbb{E}[f(A)] = \mathbb{E}[f(B)].$$

(ب) این نتیجه اساساً از قسمت (آ) و با استفاده از خاصیت خطی بودن امید ریاضی به‌دست می‌آید. ابتدا توجه می‌کنیم که استدلال قسمت (آ) به حالت کلی زیر تعمیم می‌یابد:

$$\text{error}_1(S_n) = \text{error}_2(S_n) = \dots = \text{error}_i(S_n) = \dots = \text{error}_n(S_n),$$

که در آن $\text{error}_i(S_n)$ خطای مدلی است که روی $S_n \setminus \{(x_i, y_i)\}$ یاد گرفته شده و روی (x_i, y_i) ارزیابی شده است.

بنابراین برای امید ریاضی خطای LOOCV خواهیم داشت:

$$\mathbb{E}[\text{error}_{\text{LOOCV}}(S_n)] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(y_i - \hat{f}_{-i}(x_i))^2] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\text{error}_i(S_n)] = \mathbb{E}[\text{error}_1(S_n)].$$

در برابری اول، از قانون خطی بودن امید ریاضی استفاده کرده‌ایم تا عملگر \mathbb{E} را به داخل علامت جمع منتقل کنیم.

۸. (۱۵ نمره) جدول دادگان زیر را برای یک مسئله دسته‌بندی غذا در نظر بگیرید:

دما	اندازه	مزه	خوشایندی
گرم	کوچک	شور	خیر
سرد	بزرگ	شیرین	خیر
سرد	بزرگ	شیرین	خیر
سرد	کوچک	ترش	بله
گرم	کوچک	ترش	بله
گرم	بزرگ	شور	خیر
گرم	بزرگ	ترش	بله
سرد	کوچک	شیرین	بله
سرد	کوچک	شیرین	بله
گرم	بزرگ	شور	خیر

- (آ) آنتروپی^۴ اولیه ستون مربوط به خوشایند بودن غذا چند است؟
- (ب) فرض کنید ویژگی مزه را ریشه درخت گرفته‌ایم، اکنون information gain^۵ مربوط به این ویژگی را مشخص کنید.
- (پ) درخت تصمیم کامل را برای اجرای پیش‌بینی رسم کنید و همه محاسبات و مراحل را بیاورید.
- (ت) اگر یک غذا شیرین، کوچک، و سرد باشد، پیش‌بینی ما با درخت قسمت قبل چه خواهد بود؟
- (ث) یک نمونه bootstrap^۶ از دادگان جدول گرفته و درخت تصمیم‌گیری را دوباره بسازید.

حل.

(آ)

$$p(\text{خوشایند}) = \frac{5}{10} = \frac{1}{2}, \quad p(\text{ناخوشایند}) = \frac{5}{10} = \frac{1}{2}$$

با فرض استفاده از لگاریتم مبنای ۲ داریم:

$$H(y) = -\left(\frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2}\right) = -2 \left(\frac{1}{2} \log \frac{1}{2}\right) = \log 2 = 1.$$

(ب)

$$IG(y, m) = H(y) - H(y | m).$$

برای ویژگی «مزه» داریم:

$$p(\text{شور}) = \frac{3}{10} \Rightarrow \mathbb{P}(\text{ناخوشایند} | \text{شور}) = 1,$$

$$p(\text{شیرین}) = \frac{4}{10} \Rightarrow \mathbb{P}(\text{خوشایند} | \text{شیرین}) = \frac{1}{2},$$

$$p(\text{ترش}) = \frac{3}{10} \Rightarrow \mathbb{P}(\text{خوشایند} | \text{ترش}) = 1.$$

(توجه: بر طبق قرارداد، $0 \log 0 = 0$ در نظر گرفته می‌شود.)

بنابراین آنتروپی شرطی برحسب مزه برابر است با:

$$H(y | m) = -\frac{3}{10} [0] - \frac{4}{10} \left(\frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2}\right) - \frac{3}{10} [0] = -\frac{4}{10} \log \frac{1}{2} = \frac{4}{10} = \frac{2}{5}.$$

پس:

$$IG(y(\text{خوشایندی}), m(\text{مزه})) = H(y) - H(y | m) = 1 - \frac{2}{5} = \frac{3}{5} = 0.6.$$

- (پ) برای انتخاب بهترین ویژگی به‌عنوان ریشه، IG همه ویژگی‌ها را محاسبه می‌کنیم و ویژگی با بیشترین IG را به‌عنوان ریشه درخت تصمیم برمی‌گزینیم.

اطلاعات متقابل ویژگی اندازه:

$$H(y | \text{اندازه}) = -\frac{1}{2} \left(\frac{4}{5} \log \frac{4}{5} + \frac{1}{5} \log \frac{1}{5}\right) - \frac{1}{2} \left(\frac{4}{5} \log \frac{4}{5} + \frac{1}{5} \log \frac{1}{5}\right).$$

^۴ <https://en.wikipedia.org/wiki/Entropy>

^۵ [https://en.wikipedia.org/wiki/Information_gain_\(decision_tree\)](https://en.wikipedia.org/wiki/Information_gain_(decision_tree))

^۶ [https://en.wikipedia.org/wiki/Bootstrapping_\(statistics\)](https://en.wikipedia.org/wiki/Bootstrapping_(statistics))

بنابراین:

$$H(y | \text{اندازه}) = \frac{4}{5} \log \frac{5}{4} + \frac{1}{5} \log 5 \approx 0.7219,$$

$$IG(y, \text{اندازه}) = H(y) - H(y | \text{اندازه}) \approx 1 - 0.7219 = 0.2781.$$

اطلاعات متقابل ویژگی دما:

$$H(y | \text{دما}) = -\frac{1}{2} \left(\frac{3}{5} \log \frac{3}{5} + \frac{2}{5} \log \frac{2}{5} \right) \times 2,$$

پس:

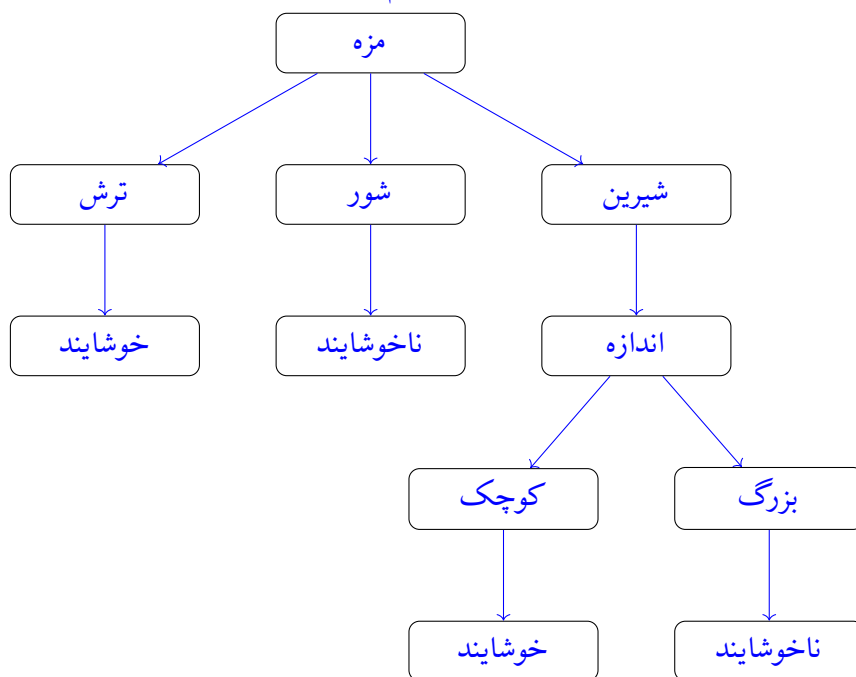
$$H(y | \text{دما}) = \frac{3}{5} \log \frac{5}{3} + \frac{2}{5} \log \frac{5}{2} \approx 0.971,$$

$$IG(y, \text{دما}) = H(y) - H(y | \text{دما}) \approx 1 - 0.971 = 0.029.$$

از بخش (ب) داشتیم:

$$IG(y, \text{مزه}) = 0.6.$$

بنابراین بهترین ویژگی برای ریشه، ویژگی مزه است و درخت تصمیم کامل به صورت زیر است:



(ت) بر اساس درخت تصمیم کامل رسم شده در قسمت (پ)، پیش‌بینی برای نمونه داده شده برابر است با:

خوشایند.

(ث) مجموعه داده bootstrap (نمونه‌گیری با جایگذاری) به صورت زیر است:
درخت تصمیم به دست آمده از این مجموعه bootstrap:

اندازه	مزه	دما	خوشایند بودن غذا
کوچک	شور	گرم	خیر
بزرگ	شیرین	سرد	خیر
بزرگ	شیرین	سرد	خیر
بزرگ	شیرین	سرد	خیر
کوچک	ترش	سرد	بله
بزرگ	شور	گرم	خیر
بزرگ	ترش	گرم	بله
بزرگ	ترش	گرم	بله
کوچک	شیرین	سرد	بله
کوچک	شیرین	سرد	بله

