



دانشگاه صنعتی شریف  
دانشکده مهندسی کامپیوتر

## تمرین تئوری اول

پارسا ملکیان - ۰۱۷۱۰۷۵

یادگیری ماشین  
دکتر شریفی

## ۱. (الف) نادرست

SVM soft-margin دقيقاً برای همین حالت طراحی شده است. وقتی داده‌ها خطی قابل تفکیک نیستند، margin با استفاده از پارامتر  $C$  و متغیرهای slack ( $\xi_i$ ) اجازه می‌دهد برخی نقاط داخل یا سمت اشتباه قرار بگیرند. هدف soft-margin دقیقاً رسیدگی به داده‌های غیرخطی یا نویزدار است.

## (ب) درست

Random Forest به دلیل استفاده از bagging در برابر نویز مقاوم‌تر است. چون از میانگین‌گیری روی چندین درخت مستقل استفاده می‌کند، نویزها خنثی می‌شوند. اما boosted trees روی اشتباهات تمرکز می‌کند و وزن داده‌های اشتباه را افزایش می‌دهد، بنابراین ممکن است روی داده‌های نویزدار overfit کند.

## (پ) درست

در bagging معمولاً از درخت‌های عمیق (fully grown) استفاده می‌کنیم که واریانس بالا و بایاس کم دارند. سپس با ترکیب نتایج  $n$  درخت از طریق averaging یا voting واریانس کل کاهش می‌یابد در حالی که بایاس تقریباً ثابت می‌ماند.

## (ت) درست

استفاده از مدل بسیار پیچیده روی داده محدود منجر به overfitting می‌شود. چنین مدلی توانایی شدن کامل روی داده‌های train را دارد (بایاس کم) اما روی داده‌های جدید عملکرد ضعیفی دارد (واریانس بالا). این همان کلاسیک trade-off bias-variance است.

## (ث) درست

طبق قضیه همگرایی پرسپترون (Perceptron Convergence Theorem)، اگر مجموعه داده به صورت خطی قابل تفکیک باشد، الگوریتم پرسپترون تضمین می‌کند که پس از تعداد محدودی iteration (حداکثر  $\frac{R^2}{\gamma^2}$  که  $R$  شعاع داده‌ها و  $\gamma$  حاشیه است) به یک راه حل همگرا می‌شود.

## ۲. مدل چگالی kNN:

درروش kNN، برای تخمین چگالی نقطه  $x$ ، کره‌ای با مرکز  $x$  و شعاع  $r(x)$  در نظر می‌گیریم که دقیقاً  $k$  نقطه از داده‌ها را شامل شود. تخمین چگالی به صورت زیر است:

$$P(x) = \frac{k}{n \cdot V(x)}$$

که در آن:

-  $k$  تعداد همسایگان نزدیک

- $n$  تعداد کل نمونه‌های آموزشی
- حجم کره با شعاع  $r(x)$  که  $k$  نقطه را دربر می‌گیرد  $V(x)$
- برای فضای  $d$  بعدی:  $C_d \cdot r(x)^d = C_d \cdot V(x)$  ثابت هندسی است.

اثبات نامناسب بودن:

فرض کنید کل داده‌های ما در یک ناحیه محدود  $\mathcal{X}$  قرار دارند. حالا انتگرال کل چگالی را حساب می‌کنیم:

$$\int_{R^d} P(x) dx$$

این انتگرال را به دو قسمت تقسیم می‌کنیم:

$$\int_{\mathcal{X}} P(x) dx + \int_{R^d \setminus \mathcal{X}} P(x) dx$$

قسمت اول (داخل ناحیه محدود):

$$\int_{\mathcal{X}} \frac{k}{n \cdot V(x)} dx < \infty$$

این قسمت محدود است چون روی ناحیه محدود انتگرال می‌گیریم.  
قسمت دوم (بیرون از ناحیه داده‌ها):

برای نقاط  $x$  که بسیار دور از داده‌ها هستند، شعاع  $r(x)$  باید خیلی بزرگ شود تا  $k$  نقطه را پوشش دهد.  
در واقع  $\|x\| \sim r(x)$  وقتی  $\|x\| \rightarrow \infty$ .

بنابراین:  $V(x) \sim C_d \cdot \|x\|^d$

و چگالی:  $P(x) \sim \frac{k}{n \cdot C_d \cdot \|x\|^d}$

حالا انتگرال قسمت دوم:

$$\int_{R^d \setminus \mathcal{X}} \frac{k}{n \cdot C_d \cdot \|x\|^d} dx$$

با تبدیل به مختصات کروی  $(dx = r^{d-1} dr d\theta, x = r\theta)$ :

$$\int_R^\infty \int_{S^{d-1}} \frac{k}{n \cdot C_d \cdot r^d} \cdot r^{d-1} dr d\theta = \text{const} \cdot \int_R^\infty \frac{1}{r} dr = \infty$$

نتیجه: انتگرال کل واگرا است، پس  $P(x)$  یک توزیع احتمال معتبر نیست.

## اهمیت و نتایج:

### ۱. عدم نرمال سازی:

$$\int_{R^d} P(x)dx \neq 1$$

این یعنی نمی توانیم از آن به عنوان توزیع احتمال واقعی استفاده کنیم.

### ۲. مشکل در استنتاج بیزی:

نمی توانیم مستقیماً از این چگالی در فرمول بیز استفاده کنیم:

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

### ۳. محدودیت در مقایسه مدل ها:

نمی توانیم likelihood یا log-likelihood را به درستی محاسبه کنیم.

### ۴. کاربرد عملی همچنان ممکن:

با وجود این مشکل، kNN برای طبقه بندی همچنان مفید است چون فقط به نسبت چگالی ها نیاز داریم:

$$\frac{P(x|y=1)}{P(x|y=0)}$$

که در آن مخرج های  $V(x)$  ساده می شوند.

## ۳. حل سوال: اثبات نامساوی خطای Ensemble

تعاریف:

$\hat{h}_i(x)$ : مدل  $i$ -ام

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

$h(x)$ : تابع واقعی

اثبات:

گام ۱: بازنویسی خطای ensemble را می نویسیم:

$$\begin{aligned} E_x \left[ \left( h(x) - \hat{H}_m(x) \right)^2 \right] &= E_x \left[ \left( h(x) - \frac{1}{m} \sum_{i=1}^m \hat{h}_i(x) \right)^2 \right] \\ &= E_x \left[ \left( \frac{1}{m} \sum_{i=1}^m \left( h(x) - \hat{h}_i(x) \right) \right)^2 \right] \end{aligned}$$

گام ۲: استفاده از نامساوی جنسن

از آنجایی که تابع  $f(t) = t^2$  تابعی محدب است، از نامساوی جنسن برای میانگین‌ها داریم:

$$\left( \frac{1}{m} \sum_{i=1}^m a_i \right)^2 \leq \frac{1}{m} \sum_{i=1}^m a_i^2$$

با قراردادن  $a_i = h(x) - \hat{h}_i(x)$

$$\left( \frac{1}{m} \sum_{i=1}^m (h(x) - \hat{h}_i(x)) \right)^2 \leq \frac{1}{m} \sum_{i=1}^m (h(x) - \hat{h}_i(x))^2$$

گام ۳: اعمال امید ریاضی  
با گرفتن امید ریاضی از دو طرف:

$$E_x \left[ \left( \frac{1}{m} \sum_{i=1}^m (h(x) - \hat{h}_i(x)) \right)^2 \right] \leq E_x \left[ \frac{1}{m} \sum_{i=1}^m (h(x) - \hat{h}_i(x))^2 \right]$$

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

که همان نامساوی مورد نظر است.

مفهوم این حکم:

۱. کاهش خطای میانگین‌گیری:

خطای مدل ensemble (سمت راست) کمتریا مساوی میانگین خطاهای تک تک مدل‌هاست. این یعنی ترکیب مدل‌ها همیشه بهتریا حداقل به خوبی بهترین مدل فردی عمل می‌کند.

۲. کاهش واریانس:

میانگین‌گیری باعث کاهش واریانس پیش‌بینی‌ها می‌شود. اگر مدل‌ها مستقل و بی‌طرف باشند:

$$\text{error ensemble} = \frac{1}{m} \times \text{error personal of mean}$$

یعنی با افزایش  $m$ ، خطای نسبت  $1/m$  کاهش می‌یابد.

۳. اصل ensemble learning

این اساس روش‌هایی مثل Random Forest، Bagging و مدل‌های ensemble است. با آموزش چند مدل مختلف و میانگین‌گیری از آنها، عملکرد بهتری نسبت به تک تک مدل‌ها داریم.

۴. شرط مهم:

این کاهش خطای زمانی مؤثرتر است که مدل‌ها متنوع و خطاهایشان ناهمبسته باشد. اگر همه مدل‌ها یکسان باشند، هیچ بهبودی حاصل نمی‌شود.

## ۴. سوال ۴: Random Forest و نمونه‌گیری

(الف) دلایل خطای بالای Random Forest روی test:

مشکل: خطای پایین روی train و خطای بالا روی test → Overfitting

دلیل ۱: عمق زیاد درختها

در Random Forest، اگر درخت‌ها خیلی عمیق باشند، هر درخت می‌تواند داده‌های train را کاملاً حفظ (memorize). کند.

راه حل:

- محدود کردن عمق درخت‌ها: تنظیم `max_depth`

- افزایش حداقل نمونه در هر برگ: `min_samples_leaf`

- محدود کردن تعداد نمونه برای `min_samples_split`

دلیل ۲: تعداد کم درخت‌ها یا کمبود تنوع

اگر تعداد درخت‌ها کم باشد یا همبستگی بین آنها زیاد باشد، قدرت ensemble کاهش می‌یابد و مدل به راحتی overfit می‌کند.

راه حل:

- افزایش تعداد درخت‌ها: `n_estimators`

- کاهش `max_features` برای افزایش تنوع بین درخت‌ها

- استفاده از `bootstrap sampling` متنوع‌تر

(ب) محاسبه میانگین و واریانس نمونه‌گیری:

داده شده:

- داده‌های اصلی:  $\text{Var}(x_i) = \sigma^2$  و  $E[x_i] = \mu$  با  $X = \{x_i\}_{i=1}^n$

- نمونه‌گیری:  $z_m = \frac{1}{m} \sum_{i=1}^m a_i$  که به صورت تصادفی و با جایگذاری از  $X$  انتخاب می‌شوند

محاسبه میانگین  $: z_m$

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

نتیجه:  $E[z_m] = \mu$  (نمونه بی‌طرف است)

محاسبه واریانس  $: z_m$

چون  $a_i$  ها مستقل هستند (نمونه‌گیری با جایگذاری):

$$\text{Var}(z_m) = \text{Var} \left( \frac{1}{m} \sum_{i=1}^m a_i \right) = \frac{1}{m^2} \sum_{i=1}^m \text{Var}(a_i)$$

$$= \frac{1}{m^2} \cdot m\sigma^2 = \frac{\sigma^2}{m}$$

$$\text{نتیجه: } \text{Var}(z_m) = \frac{\sigma^2}{m}$$

مقایسه  $\{z_m\}$  با  $X$ :

خوبی‌ها:

۱. کاهش واریانس:

$$\text{Var}(z_m) = \frac{\sigma^2}{m} < \sigma^2 = \text{Var}(x_i)$$

با افزایش  $m$ ، واریانس به ۰ میل می‌کند. این باعث پایداری بیشتر تخمین می‌شود.

۲. بی‌طرفی:

$$E[z_m] = \mu$$

میانگین نمونه، تخمین بی‌طرفی از میانگین جامعه است.

۳. کاهش نویز: با میانگین‌گیری، اثر نقاط پرت (outlier) کاهش می‌یابد.

بدی‌ها:

۱. از دست دادن تنوع: همه  $z_m$  ها نزدیک  $\mu$  هستند، در حالی که  $X$  تنوع بیشتری دارد.

۲. کاهش اطلاعات:  $m$  نمونه اصلی را به یک عدد تبدیل می‌کنیم، اطلاعات زیادی از دست می‌رود.

۳. همبستگی در bootstrap: اگر از نمونه‌گیری با جایگذاری استفاده کنیم، برخی نمونه‌ها تکرار می‌شوند و برخی دیگر حذف می‌شوند.

شرایط استفاده مناسب:

مناسب است زمانی که:

۱. هدف تخمین میانگین است: وقتی می‌خواهیم  $\mu$  را تخمین بزنیم،  $z_m$  تخمین بهتری نسبت به تک نمونه است.

۲. داده دارای نویز زیاد است: میانگین‌گیری باعث کاهش اثر نویز می‌شود.

۳. در ensemble methods: در Random Forest و Bagging، این نمونه‌گیری برای ایجاد تنوع بین مدل‌ها استفاده می‌شود.

۴. برای برآورد عدم قطعیت: با ایجاد چندین  $z_m$ ، می‌توانیم فاصله اطمینان محاسبه کنیم. نامناسب است زمانی که:

۱. نیاز به حفظ تنوع داده‌ها: اگر تنوع و توزیع کامل داده مهم باشد.

۲.  $m$  خیلی کم است: واریانس هنوز زیاد است و تخمین قابل اعتماد نیست.

۳. نمونه اصلی کم است ( $n$  کوچک): نمونه‌گیری مجدد ممکن است به overfitting منجر شود.

## ۵. (الف) مشتق تابع سیگموئید و مشکل gradient vanishing

تابع سیگموئید:

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

محاسبه مشتق (گام به گام):

گام ۱: از قاعده کسر استفاده می‌کنیم:

$$\frac{d\sigma}{dz} = \frac{d}{dz} \left( \frac{1}{1 + e^{-z}} \right)$$

گام ۲: فرض کنیم  $\sigma = \frac{1}{u}$ ,  $u = 1 + e^{-z}$ , پس

$$\frac{d\sigma}{dz} = -\frac{1}{u^2} \cdot \frac{du}{dz}$$

گام ۳: محاسبه  $\frac{du}{dz}$ :

$$\frac{du}{dz} = \frac{d}{dz}(1 + e^{-z}) = -e^{-z}$$

گام ۴: ترکیب:

$$\frac{d\sigma}{dz} = -\frac{1}{(1 + e^{-z})^2} \cdot (-e^{-z}) = \frac{e^{-z}}{(1 + e^{-z})^2}$$

گام ۵: ساده‌سازی با استفاده از  $\sigma(z) = \frac{1}{1+e^{-z}}$

$$\frac{e^{-z}}{(1 + e^{-z})^2} = \frac{1}{1 + e^{-z}} \cdot \frac{e^{-z}}{1 + e^{-z}}$$

$$= \frac{1}{1 + e^{-z}} \cdot \frac{1 + e^{-z} - 1}{1 + e^{-z}} = \sigma(z) \cdot (1 - \sigma(z))$$

نتیجه نهایی:

$$\frac{d\sigma(z)}{dz} = \sigma(z)(1 - \sigma(z))$$

رسم نمودارها و تحلیل:

نمودار تابع سیگموئید:

- برای  $\sigma(z) \rightarrow 0$ :  $z \rightarrow -\infty$

- برای  $\sigma(0) = 0.5$ :  $z = 0$

- برای  $\sigma(z) \rightarrow 1$ :  $z \rightarrow +\infty$

نمودار مشتق:

$\sigma'(0) = 0.25$ : بیشترین مقدار در  $z = 0$

$$\sigma'(z) \approx 0 : |z| > 3$$

مشکل gradient: vanishing

وقتی  $|z|$  بزرگ باشد (مثلاً  $z > 5$  یا  $z < -5$ ):

$$\sigma'(z) \approx 0$$

در بروزرسانی وزن‌ها با descent: gradient

$$\theta_{new} = \theta_{old} - \alpha \frac{\partial L}{\partial \theta}$$

که:

$$\frac{\partial L}{\partial \theta} = \frac{\partial L}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial z} \cdot \frac{\partial z}{\partial \theta}$$

مشکل: وقتی  $0 \approx \sigma'(z)$ , گرادیان خیلی کوچک می‌شود  $\Rightarrow$  بروزرسانی وزن‌ها خیلی کند  $\Rightarrow$  یادگیری متوقف می‌شود.

این مشکل در لایه‌های اولیه شبکه‌های عمیق شدیدتر است (gradient) به صورت ضربی کاهش می‌یابد).

## (ب) همگرایی کند رگرسیون لجستیک

تابع هزینه رگرسیون لجستیک:

$$L(\theta) = -\frac{1}{n} \sum_{i=1}^n [y_i \log(\sigma(z_i)) + (1 - y_i) \log(1 - \sigma(z_i))]$$

$$. z_i = \theta^T x_i$$

دلیل همگرایی کند برای داده‌های خطی قابل تفکیک:

وقتی داده‌ها کاملاً قابل تفکیک هستند، مدل تلاش می‌کند احتمال  $P(y = 1|x)$  را برای کلاس مثبت به 1 و برای کلاس منفی به 0 برساند.

برای این کار،  $|\theta^T x|$  باید به  $\infty$  برود:

$$\theta^T x \rightarrow +\infty \Rightarrow \sigma(\theta^T x) \rightarrow 1$$

$$\theta^T x \rightarrow -\infty \Rightarrow \sigma(\theta^T x) \rightarrow 0$$

مشکل:

هیچ مقدار محدودی از  $\theta$  نمی‌تواند  $\sigma$  را دقیقاً به 0 یا 1 برساند، پس الگوریتم همیشه تلاش می‌کند  $||\theta||$  را افزایش دهد.

اما وقتی  $|\theta^T x|$  خیلی بزرگ شود:

$$\sigma'(\theta^T x) \approx 0$$

گرادیان خیلی کوچک می‌شود:

$$\nabla_{\theta} L = \frac{1}{n} \sum_{i=1}^n (\sigma(z_i) - y_i) x_i \approx 0$$

نتیجه: بروزرسانی‌ها خیلی کوچک  $\Rightarrow$  همگرایی خیلی کند  $\Rightarrow \|\theta\| \rightarrow \infty$  به آرامی.  
راحل: استفاده از regularization (مثل  $L_2$ ) که  $\|\theta\|$  را محدود می‌کند.

### (ج) در طبقه‌بندی چندکلاسه invariance

مدل softmax:

$$P(y = k|x) = \frac{e^{\theta_k^T x}}{\sum_{j=1}^K e^{\theta_j^T x}}$$

هدف: نشان دهیم اگر به همه  $\theta_k$  یک بردار ثابت  $c$  اضافه کنیم، احتمالات تغییر نمی‌کنند.  
اثبات:

فرض کنیم  $\theta'_k = \theta_k + c$  برای همه  $k = 1, \dots, K$

احتمال جدید:

$$\begin{aligned} P'(y = k|x) &= \frac{e^{(\theta_k+c)^T x}}{\sum_{j=1}^K e^{(\theta_j+c)^T x}} \\ &= \frac{e^{\theta_k^T x + c^T x}}{\sum_{j=1}^K e^{\theta_j^T x + c^T x}} \\ &= \frac{e^{\theta_k^T x} \cdot e^{c^T x}}{\sum_{j=1}^K e^{\theta_j^T x} \cdot e^{c^T x}} \\ &= \frac{e^{\theta_k^T x} \cdot e^{c^T x}}{e^{c^T x} \cdot \sum_{j=1}^K e^{\theta_j^T x}} \\ &= \frac{e^{\theta_k^T x}}{\sum_{j=1}^K e^{\theta_j^T x}} = P(y = k|x) \end{aligned}$$

نتیجه: در نتیجه softmax فقط به تفاوت نسبی بین  $\theta_k$  ها وابسته است، نه مقادیر مطلق آنها.

این یعنی پارامترها unique نیستند و بی‌نهایت راه حل وجود دارد که همه یک پیش‌بینی می‌دهند.  
راحل عملی: معمولاً  $\theta_K = 0$  قرار می‌دهیم (یک کلاس را reference می‌گیریم).

## ۶. (الف) SVM در Geometric Margin

بخش ۱: چرا geometric margin برابر است با  $\gamma = \min_i \frac{y_i(w^T x_i + b)}{\|w\|}$  فاصله یک نقطه  $x_i$  از ابرصفحه  $w^T x + b = 0$  به صورت زیر محاسبه می‌شود:

$$d_i = \frac{|w^T x_i + b|}{\|w\|}$$

برای طبقه‌بندی صحیح، می‌خواهیم  $y_i(w^T x_i + b) > 0$  (نقطه در سمت درست باشد). فاصله (با علامت) signed:

$$d_i^{\text{signed}} = \frac{y_i(w^T x_i + b)}{\|w\|}$$

کمترین فاصله از ابرصفحه تا نزدیک‌ترین نقطه Geometric margin:

$$\gamma = \min_i d_i^{\text{signed}} = \min_i \frac{y_i(w^T x_i + b)}{\|w\|}$$

توضیح:

- صورت کسر  $y_i(w^T x_i + b)$  فاصله functional است
- مخرج  $\|w\|$  برای نرمال‌سازی و تبدیل به فاصله واقعی هندسی
- حداقل روی همه نقاط  $\Rightarrow$  نزدیک‌ترین نقطه به مرز

بخش ۲: چرا کاهش  $\frac{1}{2} \|w\|^2$  باعث افزایش margin می‌شود؟

از شرط مسئله:  $y_i(w^T x_i + b) \geq 1$  برای همه  $i$ . پس:

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

رابطه معکوس:

$$\gamma \geq \frac{1}{\|w\|} \Rightarrow \gamma \text{ کاهش افزایش} \Rightarrow \|w\| \text{ کاهش}$$

استدلال هندسی:

مسئله بهینه‌سازی:

$$\min_{w,b} \frac{1}{2} \|w\|^2 \quad \text{s.t.} \quad y_i(w^T x_i + b) \geq 1$$

شرط ۱:  $y_i(w^T x_i + b) \geq 1$  functional margin یعنی  $y_i(w^T x_i + b) \geq 1$  است.

چون  $\gamma = \frac{\text{margin functional}}{\|w\|}$  با ثابت نگه داشتن صورت کسر ( $\geq 1$ ) functional margin و کاهش دادن مخرج ( $\|w\|$ ) افزایش margin می‌یابد:

$$\gamma \text{ کمینه‌سازی} \Leftrightarrow \frac{1}{\|w\|} \text{ بیشینه‌سازی} \Leftrightarrow \gamma \text{ بیشینه‌سازی}$$

نتیجه: با کمینه کردن  $\frac{1}{2} \|w\|^2$  (که معادل کمینه کردن  $\|w\|$  است)، geometric margin را بیشینه می‌کنیم.

## (ب) تخمین‌گر Least Squares

مدل:

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

که  $\beta \in R^p, X \in R^{n \times p}, y \in R^n$

بخش ۱: فرم بسته تخمین‌گر Least Squares

تابع هزینه:

$$L(\beta) = \|y - X\beta\|^2 = (y - X\beta)^T(y - X\beta)$$

بازکردن:

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

گرفتن مشتق نسبت به  $\beta$ :

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

صرف قراردادن برای یافتن نقطه بهینه:

$$-2X^T y + 2X^T X \beta = 0$$

$$X^T X \beta = X^T y$$

فرم بسته (اگر  $X^T X$  معکوس پذیر باشد):

$$\boxed{\hat{\beta} = (X^T X)^{-1} X^T y}$$

بخش ۲: اثبات unbiased بودن

امید ریاضی از  $\hat{\beta}$ :

$$E[\hat{\beta}] = E[(X^T X)^{-1} X^T y]$$

جایگذاری  $y = X\beta + \epsilon$

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

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

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

$E[\epsilon] = 0$  است و deterministic  $X$  چون

$$E[\hat{\beta}] = \beta + (X^T X)^{-1} X^T E[\epsilon] = \beta + 0 = \beta$$

$$\boxed{E[\hat{\beta}] = \beta \Rightarrow \text{unbiased}} \text{ است تخمین‌گر}$$

### بخش ۳: ماتریس Covariance

$$\text{Cov}(\hat{\beta}) = E[(\hat{\beta} - E[\hat{\beta}])(\hat{\beta} - E[\hat{\beta}])^T]$$

چون  $E[\hat{\beta}] = \beta$

$$\text{Cov}(\hat{\beta}) = E[(\hat{\beta} - \beta)(\hat{\beta} - \beta)^T]$$

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

محاسبه:

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

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

پس:

$$\text{Cov}(\hat{\beta}) = E[(X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1}]$$

چون  $X$  ثابت است:

$$= (X^T X)^{-1} X^T E[\epsilon \epsilon^T] X (X^T X)^{-1}$$

$E[\epsilon \epsilon^T] = \text{Cov}(\epsilon) = \sigma^2 I$  و

$$= (X^T X)^{-1} X^T (\sigma^2 I) X (X^T X)^{-1}$$

$$= \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1}$$

$$= \sigma^2 (X^T X)^{-1}$$

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

نتیجه: واریانس تخمین‌گر به  $\sigma^2$  (واریانس نویز) و (هندسه داده‌ها) بستگی دارد.

۷. (الف) اثبات برای نمونه اول

هدف: نشان دهیم:

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

که در آن  $(x, y)$  یک نمونه تصادفی از توزیع داده است.

اثبات:

خطای نمونه اول:

$$e_1(S_n) = (y_1 - \hat{f}_{S_{n-1}}(x_1))^2$$

امید ریاضی از  $e_1$ :

$$E[e_1(S_n)] = E \left[ (y_1 - \hat{f}_{S_{n-1}^{-1}}(x_1))^2 \right]$$

نکته کلیدی: امید ریاضی روی همه نمونه‌های تصادفی ممکن  $S_n$  است.  
برای محاسبه این امید، می‌توانیم آن را به صورت زیر بنویسیم:

$$E[e_1(S_n)] = E_{S_n} \left[ (y_1 - \hat{f}_{S_{n-1}^{-1}}(x_1))^2 \right]$$

این را می‌توان به صورت امید مشروط نوشت:

$$= E_{(x_1, y_1)} \left[ E_{S_{n-1}^{-1}|(x_1, y_1)} \left[ (y_1 - \hat{f}_{S_{n-1}^{-1}}(x_1))^2 \right] \right]$$

اما چون  $\hat{f}_{S_{n-1}^{-1}}$  تابعی از  $S_{n-1}^{-1}$  است (که شامل  $(x_1, y_1)$  نیست):

$$= E_{(x_1, y_1)} \left[ E_{S_{n-1}^{-1}} \left[ (y_1 - \hat{f}_{S_{n-1}^{-1}}(x_1))^2 | (x_1, y_1) \right] \right]$$

چون  $S_{n-1}^{-1}$  و  $(x_1, y_1)$  مستقل هستند:

$$= E_{(x_1, y_1), S_{n-1}^{-1}} \left[ (y_1 - \hat{f}_{S_{n-1}^{-1}}(x_1))^2 \right]$$

حالا  $(x_1, y_1)$  یک نمونه تصادفی از توزیع است و  $S_{n-1}^{-1}$  مجموعه‌ای از  $n-1$ -نمونه مستقل دیگر است.  
اگر  $(x, y)$  را یک نمونه تصادفی جدید از همان توزیع  $S_{n-1}$  را مجموعه‌ای از  $n-1$ -نمونه بنامیم، داریم:

$$E[e_1(S_n)] = E_{(x, y), S_{n-1}} \left[ (y - \hat{f}_{S_{n-1}}(x))^2 \right]$$

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

تفسیر: خطای مورد انتظار برای نمونه اول برابر است با خطای مورد انتظار برای یک نمونه تصادفی جدید وقتی مدل روی  $1-n$ -نمونه آموزش دیده است.

(ب) تعمیم به LOOCV کامل

هدف: نشان دهیم:

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

اثبات:

خطای LOOCV:

$$e(S_n) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}_{S_{n-1}^{-i}}(x_i))^2 = \frac{1}{n} \sum_{i=1}^n e_i(S_n)$$

امید ریاضی:

$$E[e(S_n)] = E\left[\frac{1}{n} \sum_{i=1}^n e_i(S_n)\right]$$

با استفاده از خطی بودن امید ریاضی:

$$= \frac{1}{n} \sum_{i=1}^n E[e_i(S_n)]$$

نکته کلیدی: به دلیل تقارن، همه نمونه‌ها از همان توزیع هستند و نقش یکسانی دارند:

$$E[e_1(S_n)] = E[e_2(S_n)] = \dots = E[e_n(S_n)]$$

بنابراین:

$$E[e(S_n)] = \frac{1}{n} \sum_{i=1}^n E[e_i(S_n)] = \frac{1}{n} \cdot n \cdot E[e_1(S_n)]$$

$$= E[e_1(S_n)]$$

از نتیجه قسمت (الف):

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

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

تفسیر و اهمیت:

۱. تخمین بی‌طرف خطای تعمیم:

LOOCV یک تخمین‌گر (تقریباً) بی‌طرف از خطای تعمیم است. خطای مورد انتظار آن برابر است با خطای واقعی مدل آموزش دیده روی  $1 - n$  نمونه.

۲. مقایسه با error: test

اگر مدل روی  $n$  نمونه آموزش ببیند، خطای test برابر است با:

$$\text{Error Test} = E[(y - \hat{f}_{S_n}(x))^2]$$

خطای LOOCV کمی بدینانه تراست (چون فقط  $1 - n$  نمونه استفاده می‌کند).

۳. واریانس بالا:

با اینکه LOOCV تخمین‌گر بی‌طرف است، اما واریانس بالایی دارد چون مدل‌های  $\hat{f}_{S_{n-1}^{(i)}}$  همبستگی زیادی با هم دارند (تقریباً روی همان داده‌ها آموزش دیده‌اند).

۴. cross-validation: k-fold

به همین دلیل معمولاً از CV k-fold (مثلاً  $k = 5$  یا  $k = 10$ ) استفاده می‌شود که trade-off بین variance و bias دارد.

## ۸. سوال ۸: درخت تصمیم و Bootstrap

داده‌ها:

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

تعداد کل:  $n = 10$ , بله: ۵, خیر: ۵

(الف) آنتروپی اولیه

$$H(S) = - \sum_i p_i \log_2(p_i)$$

$$H(S) = -\frac{5}{10} \log_2 \left( \frac{5}{10} \right) - \frac{5}{10} \log_2 \left( \frac{5}{10} \right)$$

$$= -0.5 \times (-1) - 0.5 \times (-1) = 0.5 + 0.5$$

$H(S) = 1$

(ب) Gain Information برای ویژگی مزه

تقسیم براساس مزه:

مزه = شور: نمونه‌های  $\{1, 6, 10\}$  ⇒ همه خیر (۳ نمونه)

$$H(\text{salty}) = 0$$

مزه = شیرین: نمونه‌های  $\{2, 3, 8, 9\}$  ⇒  $\{1, 6, 10\}$  خیر (۴ نمونه)

$$H(\text{sweet}) = -\frac{3}{4} \log_2 \left( \frac{3}{4} \right) - \frac{1}{4} \log_2 \left( \frac{1}{4} \right)$$

$$= -0.75 \times (-0.415) - 0.25 \times (-2) = 0.311 + 0.5 = 0.811$$

مزه = ترش: نمونه‌های  $\{4, 5, 7\}$  ⇒ همه بله (۳ نمونه)

$$H(\text{sour}) = 0$$

آنتروپی شرطی:

$$\begin{aligned} H(S|\text{taste}) &= \frac{3}{10} \times 0 + \frac{4}{10} \times 0.811 + \frac{3}{10} \times 0 \\ &= 0.324 \end{aligned}$$

Gain: Information

$$IG(\text{taste}) = H(S) - H(S|\text{taste}) = 1 - 0.324 = 0.676$$

(پ) ساخت درخت تصمیم کامل

گام ۱: انتخاب ریشه

باید Gain Information برای همه ویژگی‌ها محاسبه کنیم:

دما: - گرم:  $\Rightarrow \{2, 3, 4, 8, 9\}$  - سرد:  $\Rightarrow \{1, 5, 6, 7, 10\}$   $H = 0.971$

$$H = 0.971$$

$$H(S|\text{temperature}) = \frac{5}{10} \times 0.971 + \frac{5}{10} \times 0.971 = 0.971$$

$$IG(\text{temperature}) = 1 - 0.971 = 0.029$$

اندازه: - کوچک:  $\Rightarrow \{2, 3, 6, 7, 10\}$  - بزرگ:  $\Rightarrow \{1, 4, 5, 8, 9\}$   $H = 0.722$

$$H = 0.722$$

$$H(S|\text{size}) = \frac{5}{10} \times 0.722 + \frac{5}{10} \times 0.722 = 0.722$$

$$IG(\text{size}) = 1 - 0.722 = 0.278$$

مزه: (از قسمت قبل)  $IG = 0.676$

نتیجه: بیشترین IG مربوط به مزه است ⇒ ریشه = مزه

گام ۲: شاخه شور (pure) همه خیر ⇒ بزرگ = خیر

گام ۳: شاخه ترش (pure) همه بله ⇒ بزرگ = بله

گام ۴: شاخه شیرین نمونه‌های  $\{2, 3, 8, 9\}$ :  $\Rightarrow \{3\}$  بله،  $\{2\}$  خیر

محاسبه IG برای ویژگی‌های باقیمانده:

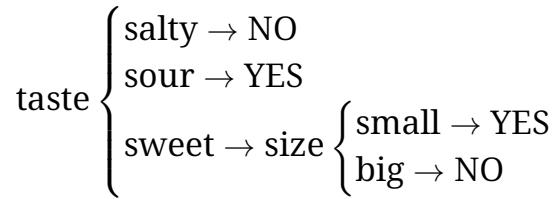
دما در شیرین: - سرد:  $\Rightarrow \{2, 3, 8, 9\}$  همه سرد  $\Rightarrow \{3\}$  بله،  $\{2\}$  خیر

اندازه در شیرین: - کوچک:  $\Rightarrow \{8, 9\}$  همه خیر  $\Rightarrow \{3\}$  بزرگ

$$H(S|\text{size}) = \frac{2}{4} \times 0 + \frac{2}{4} \times 0 = 0$$

$$IG(\text{size}) = 0.811 - 0 = 0.811$$

اندازه بهترین split است.  
درخت نهایی:



(ت) پیش‌بینی برای غذای شیرین، کوچک، سرد

مسیر در درخت:  
 $\text{sweet} = \text{taste} \rightarrow \text{small} = \text{size} \rightarrow \boxed{\text{YES}}$

(ث) نمونه Bootstrap

نمونه‌گیری Bootstrap: با جایگذاری ۱۰ نمونه از داده اصلی انتخاب می‌کنیم.  
 فرض کنید نمونه bootstrap ما:  $\{1, 1, 3, 4, 5, 6, 7, 8, 9, 10\}$  است.  
 (نمونه ۱ دوبار، نمونه ۲ حذف شده)  
 داده جدید:  
 - بله: ۵، خیر: ۵  
 $H(S) = 1 -$

محاسبات مشابه قبل انجام می‌شود و درخت ساخته می‌شود.  
 نکته: به دلیل تکرار برخی نمونه‌ها و حذف برخی دیگر، ممکن است درخت متفاوتی به دست آید که تنوع را در Random Forest ایجاد می‌کند.  
 در این مثال خاص، چون نمونه ۲ و ۳ مشابه بودند، احتمالاً درخت مشابهی ساخته می‌شود، اما در حالت کلی تفاوت‌هایی خواهد داشت.