

**تمرین ۱**

فرار است برای دو متغیر تصادفی  $X$  و  $Y$  ثابت شود:

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2 \text{Cov}(X, Y)$$

تعريف واریانس در صورت سوال:

$$\text{Var}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$$

حل:

$$\text{Var}(X + Y) = \mathbb{E}[(X + Y)^2] - (\mathbb{E}[X + Y])^2$$

از خطی بودن امید ریاضی میتوان نتیجه گرفت:

$$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$$

پس:

$$\mathbb{E}[(X + Y)^2] = \mathbb{E}[X^2 + 2XY + Y^2] = \mathbb{E}[X^2] + 2\mathbb{E}[XY] + \mathbb{E}[Y^2]$$

$$(\mathbb{E}[X + Y])^2 = (\mathbb{E}[X] + \mathbb{E}[Y])^2 = (\mathbb{E}[X])^2 + 2\mathbb{E}[X]\mathbb{E}[Y] + (\mathbb{E}[Y])^2$$

جایگذاری:

$$\begin{aligned} \text{Var}(X + Y) &= (\mathbb{E}[X^2] + 2\mathbb{E}[XY] + \mathbb{E}[Y^2]) - ((\mathbb{E}[X])^2 + 2\mathbb{E}[X]\mathbb{E}[Y] + (\mathbb{E}[Y])^2) \\ &= \underbrace{(\mathbb{E}[X^2] - (\mathbb{E}[X])^2)}_{\text{Var}(X)} + \underbrace{(\mathbb{E}[Y^2] - (\mathbb{E}[Y])^2)}_{\text{Var}(Y)} + 2 \underbrace{(\mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y])}_{\text{Cov}(X, Y)} \end{aligned}$$

پس:

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2 \text{Cov}(X, Y)$$

**تمرین ۲**

صورت مسئله. تابع هزینه زیر را در نظر بگیرید:

$$J_1(w) = \|y - Xw\|_2^2 + \lambda_2\|w\|_2^2 + \lambda_1\|w\|_1.$$

تعريف:

$$c = (1 + \lambda_2)^{-1/2}, \quad \tilde{X} = c \begin{pmatrix} X \\ \sqrt{\lambda_2} I_d \end{pmatrix}, \quad \tilde{y} = \begin{pmatrix} y \\ 0_{d \times 1} \end{pmatrix},$$

و

$$J_2(w) = \|\tilde{y} - \tilde{X}w\|_2^2 + c\lambda_1\|w\|_1.$$

حل: از آنجا که

$$\tilde{X}w = c \begin{pmatrix} Xw \\ \sqrt{\lambda_2} w \end{pmatrix},$$

داریم:

$$\begin{aligned} \|\tilde{y} - \tilde{X}w\|_2^2 &= \|y - cXw\|_2^2 + \|0 - c\sqrt{\lambda_2} w\|_2^2 \\ &= \|y - cXw\|_2^2 + c^2\lambda_2\|w\|_2^2. \end{aligned}$$

در نتیجه:

$$J_2(w) = \|y - cXw\|_2^2 + c^2\lambda_2\|w\|_2^2 + c\lambda_1\|w\|_1.$$

اکنون تغییر متغیر زیر را اعمال می‌کنیم:

$$v = cw \iff w = \frac{v}{c}.$$

با جایگذاری در تابع هزینه:

$$\begin{aligned} J_2\left(\frac{v}{c}\right) &= \left\|y - cX\frac{v}{c}\right\|_2^2 + c^2\lambda_2\left\|\frac{v}{c}\right\|_2^2 + c\lambda_1\left\|\frac{v}{c}\right\|_1 \\ &= \|y - Xv\|_2^2 + \lambda_2\|v\|_2^2 + \lambda_1\|v\|_1 \\ &= J_1(v). \end{aligned}$$

بنابراین برای هر  $v$  رابطه زیر برقرار است:

$$J_1(v) = J_2(v/c).$$

از این تطابق یک به یک نتیجه می‌گیریم که کمینه‌سازهای دو مسئله با رابطه مقیاسی ساده به هم مرتبط هستند:

$$v^* = \arg \min_v J_1(v) \iff \frac{v^*}{c} = \arg \min_w J_2(w).$$

معادل آن:

$$\boxed{\arg \min_w J_1(w) = c \arg \min_w J_2(w)}.$$

این یعنی مسئله اصلی  $J_1$  را می‌توان با مسئله  $J_2$  حل کرد که نسخه نرمال‌سازی شده آن است و فقط یک ضریب مقیاس بین پاسخهای دو مسئله وجود دارد.

### تمرین ۳

فرض کنید برای هر مؤلفه وزن  $w_i$  یک مدل (برنولی-گوسی) قرار دهیم. معادل صریح آن از طریق متغیر  $\{0, 1\} \in z_i$  بیان می‌شود:

$$p(z_i) = \text{Bernoulli}(\pi), \quad p(w_i | z_i) = \begin{cases} \delta(w_i), & z_i = 0, \\ \mathcal{N}(0, \sigma^2), & z_i = 1. \end{cases}$$

بنابراین توزیع  $w_i$  برابر است با

$$p(w_i) = (1 - \pi) \delta(w_i) + \pi \mathcal{N}(0, \sigma^2).$$

برای برازش (MAP) معمولاً کمینه‌سازی منفی لگاریتم احتمال انجام می‌دهیم. اگر  $L(w)$  منفی لگ احتمال داده (یا تابع هزینه داده) باشد، MAP متناظر با کمینه‌سازی

$$\min_{w,z} L(w) - \log p(w, z)$$

است. چون

$$p(w, z) = \prod_i p(z_i) p(w_i | z_i),$$

منفی لگاریتم مبحثی  $(w, z)$  به صورت جمع روی مؤلفه‌ها نوشته می‌شود:

$$-\log p(w, z) = \sum_i \left[ -\log p(z_i) - \log p(w_i | z_i) \right].$$

برای هر  $i$  داریم (تا ضرب در ثابت‌ها):

- اگر  $z_i = 0$  آنگاه  $w_i = 0$  و بخش مربوطه  $\approx -\log(1 - \pi)$ .

- اگر  $z_i = 1$  آنگاه  $w_i$  از گوسی است و بخش مربوطه  $\approx -\log \pi + \frac{w_i^2}{2\sigma^2} + \text{const.}$

پس منفی لگاریتم احتمال مشترک (جزء مربوط به prior) را می‌توان به شکل زیر نوشت:

$$-\log p(w, z) = \sum_i \left[ (1 - z_i)(-\log(1 - \pi)) + z_i \left( -\log \pi + \frac{w_i^2}{2\sigma^2} \right) \right] + \text{const.}$$

اگر حالا به صورت مشترک روی  $w$  و  $z$  مینیمیم بگیریم، برای هر مؤلفه  $i$  به صورت زیر خواهد بود:  $z_i = 1$  هنگامی انتخاب می‌شود که سود کاهش خطأ (یا مقدار  $w_i$  غیرصفر) بیشتر از هزینه فعال‌سازی  $(1 - \pi) - \log(\pi/(1 - \pi))$  باشد. به عبارت دیگر، فعال‌سازی هر مؤلفه هزینه‌ای ثابت  $\lambda_0 := -\log(\pi/(1 - \pi))$  به علاوه یک جریمه گوسی  $w_i^2/(2\sigma^2)$  وارد تابع هزینه می‌کند. در حالت تقریب:

- اگر  $\sigma^2$  را بسیار بزرگ فرض کنیم (slab خیلی پهن باشد) یا جریمه مربعی را کوچک در نظر بگیریم، آنگاه جمله مربعی  $\frac{w_i^2}{2\sigma^2}$  ناچیز می‌شود و تنها هزینه ثابت برای هر وزن  $(1 - z_i) = 1$  باقی می‌ماند.

- در این صورت مجموع هزینه‌های ثابت برای همه مؤلفه‌های فعال برابر است با  $\#\{i : w_i \neq 0\} \cdot \lambda_0$ . یعنی  $\lambda_0 \|w\|_0$ .

بنابراین MAP با این prior تقریباً معادل حداقل‌سازی

$$\min_w L(w) + \lambda_0 \|w\|_0,$$

که نشان می‌دهد مدل برنولی-گوسی می‌تواند به عنوان منبع منطقی ظهور منظم‌ساز  $\ell_0$  در تابع هزینه تفسیر شود.

## تمرین ۴

چرا منظم‌ساز  $\ell_1$  در لاسو باعث صفر شدن وزن‌ها و ایجاد پراکندگی (sparsity) می‌شود؟

منظم‌ساز  $\ell_1$  به صورت زیر تعریف می‌شود:

$$\lambda \|w\|_1 = \lambda \sum_i |w_i|.$$

تابع قدرمطلق دارای یک گوشه (non-differentiable point) در مقدار صفر است. این خاصیت باعث می‌شود حل بهینه در بسیاری موارد دقیقاً روی محور مختصات بیفت—یعنی برخی ضرایب  $w_i$  دقیقاً صفر شوند. بنابراین لاسو مدل‌هایی با تعداد پارامتر کمتر و ضرایب صفرشونده تولید می‌کند.

تفسیر هندسی و احتمالاتی: چرا  $\ell_1$  نسبت به  $\ell_2$  تمایل بیشتری به صفر کردن ضرایب دارد؟

### دیدگاه هندسی

مجموعه‌ی قیود  $c \leq \|w\|_1$  یک چندضلعی لوزی شکل است:

$$|w_1| + |w_2| \leq c.$$

گوشه‌های این چندضلعی دقیقاً روی محورهای مختصات قرار دارند. وقتی بیضی ناشی از خطای مرربعی با این چندضلعی برخورد می‌کند، احتمال این‌که نقطه‌ی تماس روی گوشه‌ها قرار گیرد زیاد است، و گوشه‌ها همان نقاطی هستند که یک یا چند مؤلفه  $w_i$  برابر صفر می‌شوند.

### دیدگاه احتمالاتی

در تفسیر بیزی، لاسو متناظر با یک پیش‌فرض لاپلاس (Laplace prior) روی وزن‌ها است:

$$p(w_i) \propto e^{-\lambda|w_i|}.$$

این پیش‌فرض شدیداً وزن‌های کوچک را ترجیح می‌دهد و احتمال زیادی برای مقدار دقیقاً صفر فراهم می‌کند. در مقابل، ریج از پیش‌فرض گاووسی استفاده می‌کند:

$$p(w_i) \propto e^{-\frac{\lambda}{2}w_i^2},$$

که وزن‌ها را به صفر \*نزدیک\* می‌کند اما به ندرت مقدار آن‌ها را دقیقاً صفر می‌سازد.

چه کاربردهایی برای منظم‌سازی  $\ell_1$  وجود دارد؟

- انتخاب ویژگی (Feature Selection) به صورت خودکار
- ایجاد مدل‌های ساده‌تر و قابل تفسیرتر
- کاربرد در داده‌های با ابعاد بالا (High-dimensional data) مانند ژنومیکس
- جلوگیری از بیش‌برازش در مدل‌هایی با تعداد ویژگی زیاد

ت) تحلیل رفتاری منظم‌ساز  $\ell_2$  (روش ریج)

منظم‌ساز  $\ell_2$  به صورت زیر است:

$$\lambda\|w\|_2^2 = \lambda \sum_i w_i^2.$$

این منظم‌ساز وزن‌ها را کوچک می‌کند اما هرگز آن‌ها را به صفر نمی‌رساند. هندسی آن یک کره (ball) است و هیچ گوشه‌ای ندارد؛ بنابراین نقطه بینه معمولًاً روی محورها نمی‌افتد.

### مزایای $\ell_2$

- مناسب برای داده‌های چندخطی (multicollinearity)
- توزیع هموارتر وزن‌ها
- بهبود پایداری مدل
- هیچ ویژگی حذف نمی‌شود، بنابراین برای وظایف پیش‌بینی (نه تفسیر) بهتر است