



تمرین سوم

شماره دانشجویی: ۹۸۱۰۱۰۷۴

محمدجواد هزاره

سوال ۱

۱. (c). زیاد بودن ضریب یک ویژگی نشان دهنده‌ی این است که با ثابت نگه داشتن سایر متغیرها و تغییر این متغیر، برچسب پیش‌بینی شده به مقدار زیادی تغییر می‌کند. در نگاه اول ممکن است این موضوع نشان دهنده‌ی این باشد که این ویژگی تاثیر زیادی در مدل دارد اما باید توجه کرد که این بزرگی ضریب می‌تواند ناشی از مقیاس آن ویژگی باشد. اگر مقیاس این ویژگی هماهنگ با سایر ویژگی‌ها نباشد، ضریب آن بسیار بزرگ یا بسیار کوچک خواهد شد. بنابراین بزرگ بودن ضریب لزوماً به معنای مهم یا کم اهمیت بودن آن ویژگی نیست.

۲. (آ) درست. با زیاد شدن داده‌ها مدل از حالت overfitting خارج می‌شود و پیچیدگی کنونی مدل نمی‌تواند تمام اطلاعاتی که داده‌ها در اختیارمان قرار می‌دهد را مدل کند. بنابراین زیاد کردن تعداد داده‌های آموزش، نه تنها باعث کاهش bias نمی‌شود، بلکه آن را افزایش می‌دهد.

(ب) غلط. با کاهش خطا روی داده‌های آموزش، مدل روی این داده‌ها فیت می‌شود و قدرت تعمیم خود را از دست خواهد داد. بنابراین لزوماً کاهش خطا روی داده‌های آموزش باعث کاهش خطا روی داده‌های تست نخواهد شد.

(ج) غلط. افزایش پیچیدگی مدل، باعث کاهش خطای آموزش می‌شود اما همواره باعث کاهش خطای تست نخواهد شد. در ابتدا که پیچیدگی مدل کم است، خطای آموزش و تست هر دو زیاد است یا به عبارتی مدل دچار underfitting شده است. با افزایش پیچیدگی مدل، خطای آموزش به مرور کم می‌شود اما خطای تست تا جایی کم شده و پس از آن شروع به افزایش می‌کند؛ چرا که مدل به سمت فیت شدن روی داده‌های آموزش می‌رود و قدرت تعمیم خود را از دست خواهد داد.

(د) غلط. می‌دانیم با افزایش پیچیدگی مدل، خطای آموزش کم می‌شود. از آن جایی که با مدل چندجمله‌ای درجه ۶ هنوز خطای آموزش به مقدار قابل توجهی زیاد است، کاهش دادن پیچیدگی مدل و استفاده از مدل خطی، نه تنها خطای آموزش را کاهش نداده بلکه باعث افزایش خطای آموزش نیز خواهد شد. برای کاهش خطای آموزش باید پیچیدگی مدل را افزایش داده یا به عبارتی از چندجمله‌ای‌های با درجه‌ی بالاتر استفاده کرد.

سوال ۲

(آ) اگر از SSE برای هزینه استفاده کنیم، آنگاه هزینه بر حسب ω برابر خواهد بود با:

$$\begin{aligned} SSE(\omega) &= \|Y - \hat{Y}\|^2 \\ &= \|Y - X\omega\|^2 \\ &= (Y - X\omega)^T(Y - X\omega) \\ &= Y^T Y - 2Y^T X\omega + \omega^T(X^T X)\omega \end{aligned}$$

برای پیدا کردن ω بهینه کافیست مشتق عبارت بالا را در ω_{opt}^* برابر صفر قرار دهیم:

$$\begin{aligned} \frac{d}{d\omega} SSE(\omega_{opt}^*) &= -2(Y^T X)^T + 2(X^T X)\omega_{opt}^* = 0 \\ \implies (X^T X)\omega_{opt}^* &= X^T Y \\ \implies \boxed{\omega_{opt}^*} &= (X^T X)^{-1} X^T Y \end{aligned}$$

(ب) اگر جمله‌ی منظم ساز L2 را به تابع هزینه اضافه کنیم، آنگاه:

$$E(\omega) = SSE(\omega) + \lambda \omega^T \omega$$

با صفر قرار دادن مشتق رابطهی بالا به ازای ω^* خواهیم داشت:

$$\begin{aligned} \frac{d}{d\omega} E(\omega^*) &= -2(Y^T X)^T + 2(X^T X)\omega^* + 2\lambda\omega^* = 0 \\ \implies -X^T Y + (X^T X)\omega^* + \lambda\omega^* &= 0 \\ \implies (X^T X + \lambda I)\omega^* &= X^T Y \\ \implies \boxed{\omega^*} &= (X^T X + \lambda I)^{-1} X^T Y \end{aligned}$$

(ج) اگر: برای این قسمت می‌دانیم رابطۀ $\Sigma X = X F$ را داریم. همچنین از آنجایی که Σ ماتریس کوواریانس است،

ماتریسی متقارن خواهد بود. پس برای ω_{new}^* خواهیم داشت:

$$\begin{aligned}
 \omega_{new}^* &= (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y \\
 &= (X^T X F^{-1})^{-1} X^T \Sigma^{-1} Y & (\Sigma^{-1} X = X F^{-1}) \\
 &= (X^T X F^{-1})^{-1} (F^T)^{-1} X^T Y & (X^T \Sigma^{-1} = (F^T)^{-1} X^T) \\
 &= (F^T (X^T X) F^{-1})^{-1} X^T Y & ((AB)^{-1} = B^{-1} A^{-1}) \\
 &= ((F^T X^T)(X F^{-1}))^{-1} X^T Y \\
 &= (X^T \Sigma \Sigma^{-1} X)^{-1} X^T Y \\
 &= (X^T X)^{-1} X^T Y \\
 &= \omega_{opt}^*
 \end{aligned}$$

فقط اگر: برای این قسمت می دانیم $\omega_{new}^* = \omega_{opt}^*$ ، بنابراین خواهیم داشت:

$$\begin{aligned}
 \omega_{new}^* &= \omega_{opt}^* \\
 \implies (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y &= (X^T X)^{-1} X^T Y \\
 \implies (X^T \Sigma^{-1} X F)^{-1} X^T \Sigma^{-1} &= (X^T X)^{-1} X^T \\
 \implies X^T &= (X^T \Sigma^{-1} X) (X^T X)^{-1} X^T \Sigma \\
 \implies X &= \Sigma X (X^T X)^{-1} (X^T \Sigma^{-1} X) \\
 \implies X (X^T \Sigma^{-1} X)^{-1} (X^T X) &= \Sigma X
 \end{aligned}$$

بنابراین اگر تعریف کنیم $F = (X^T \Sigma^{-1} X)^{-1} (X^T X)$ ، خواسته‌ی مسئله اثبات می‌شود. با توجه به وارون‌پذیر بودن $X^T X$ ، ماتریس F نیز وارون‌پذیر خواهد بود و داریم $\Sigma X = X F$.

(د) اگر داده‌های جدید \tilde{X} با برچسب \tilde{Y} را به داده‌های قبلی اضافه کنیم، آنگاه ماتریس برچسب‌ها و ماتریس داده‌ها به

صورت زیر خواهد بود:

$$\begin{cases} \bar{Y} = \begin{bmatrix} Y \\ \tilde{Y} \end{bmatrix} \\ \bar{X} = \begin{bmatrix} X \\ \tilde{X} \end{bmatrix} \end{cases}$$

بنابراین برای خطای L_1 داده‌های جدید داریم:

$$L(\omega, \lambda) = \|\bar{Y} - \bar{X}\omega\|^2 + \lambda\|\omega\|_1$$

اگر بخواهیم این تابع با تابع خطای داده شده برابر باشد، ضریب λ تابع بالا را همان ضریب خطای L_1 در تابع داده شده در نظر گرفته و داریم:

$$\begin{aligned}\|\bar{Y} - \bar{X}\omega\|^2 + \lambda_2\|\omega\|_1 &= \|Y - X\omega\|^2 + \lambda_1\|\omega\|_2^2 + \lambda_2\|\omega\|_1 \\ \|Y - X\omega\|^2 + \|\tilde{Y} - \tilde{X}\omega\|^2 &= \|Y - X\omega\|^2 + \lambda_1\|\omega\|_2^2 \\ \|\tilde{Y} - \tilde{X}\omega\|^2 &= \lambda_1\|\omega\|_2^2 \\ \tilde{Y}^T\tilde{Y} - 2\tilde{Y}^T\tilde{X}\omega + \omega^T(\tilde{X}^T\tilde{X})\omega &= \lambda_1\omega^T\omega\end{aligned}$$

بنابراین بین داده‌های جدید و برچسب‌هایشان باید رابطه‌ی زیر برقرار باشد تا تابع خطاهای مذکور با یکدیگر برابر شوند:

$$\tilde{Y}^T(\tilde{Y} - 2\tilde{X}\omega) = \omega^T(\lambda_1\mathbf{I} - \tilde{X}^T\tilde{X})\omega$$

اگر برچسب داده‌های جدید یعنی \tilde{Y} را ماتریس صفر در نظر بگیریم، سمت چپ تساوی بالا صفر شده و اگر ماتریس داده‌های جدید یعنی \tilde{X} را برابر $\sqrt{\lambda_1}\mathbf{I}$ در نظر بگیریم، سمت راست نیز صفر خواهد شد. بنابراین در این صورت خطای L_1 داده‌های \bar{X} و \bar{Y} برابر با خطای داده شده روی X و Y خواهد بود.

سوال ۳

با توجه به واگرایی Kullback-Liebler برای توزیع‌های $p(x)$ و $q(x)$ که یک بردار n بعدی است داریم:

$$KL(p||q) = \int p(x) \log\left(\frac{p(x)}{q(x)}\right) dx = \int p(x) [\log(p(x)) - \log(q(x))] dx$$

با توجه به اینکه توزیع $q(x)$ یک توزیع نمایی با میانگین μ و ماتریس کوواریانس Σ است داریم:

$$\begin{aligned}\log(q(x)) &= \log\left((2\pi)^{-\frac{n}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right)\right) \\ &= -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma|) - \frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\end{aligned}$$

بنابراین برای واگرایی KL که تابعی از μ و Σ خواهد بود داریم:

$$KL(p||q) = f(\mu, \Sigma) = \int p(x) \left(\log(p(x)) + \frac{n}{2} \log(2\pi) + \frac{1}{2} \log(|\Sigma|) + \frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right) dx$$

برای پیدا کردن کمینه‌ی این تابع باید گرادیان f را محاسبه کرده و برابر صفر قرار دهیم. برای محاسبه‌ی گرادیان، با توجه به این که انتگرال روی متغیر x گرفته می‌شود، می‌توان نخست از تابع زیر انتگرال مشتق گرفت و سپس انتگرال را محاسبه کرد. هم‌چنین مبنای لگاریتم را عدد نپر در نظر می‌گیریم و در نتیجه بجای \log داریم \ln . بنابراین برای گرادیان f داریم:

$$\begin{cases} \frac{\partial}{\partial \mu} f(\mu, \Sigma) = \int p(x) (\Sigma^{-1} (x - \mu)) dx \\ \frac{\partial}{\partial \Sigma} f(\mu, \Sigma) = \int p(x) \left(\frac{1}{2|\Sigma|} (|\Sigma| \Sigma^{-1}) - \frac{1}{2} (x - \mu)(x - \mu)^T \Sigma^{-1} \Sigma^{-1} \right) dx \end{cases}$$

که با صفر قرار دادن روابط بالا به ازای μ^* و Σ^* داریم:

$$\frac{\partial}{\partial \mu} f(\mu^*, \Sigma^*) = \Sigma^{*-1} \int p(x) (x - \mu^*) dx = 0$$

$$\implies \mu^* \int p(x) dx = \int x p(x) dx$$

$$\implies \boxed{\mu^* = \mathbb{E}_p[x]}$$

$$\frac{\partial}{\partial \Sigma} f(\mu^*, \Sigma^*) = \frac{1}{2} \left(\int p(x) \left(I - (x - \mu^*)(x - \mu^*)^T \Sigma^{*-1} \right) dx \right) \Sigma^{*-1} = 0$$

$$\implies I \int p(x) dx = \left(\int (x - \mu^*)(x - \mu^*)^T p(x) dx \right) \Sigma^{*-1}$$

$$\implies I = \text{Var}_p(x) \Sigma^{*-1}$$

$$\implies \boxed{\Sigma^* = \text{Var}_p(x)}$$

بنابراین بهترین تخمینی که از توزیع $p(x)$ با استفاده از توزیع‌های نرمال می‌توان داشت، توزیع نرمالی با میانگین μ^* و ماتریس کوواریانس Σ^* است.

سوال ۴

آ) اگر برای مدل داشته باشیم $\hat{y}^{(i)} = w_j x_j^{(i)} + \epsilon$ ، و برای تابع هزینه از MSE استفاده کنیم، آنگاه داریم:

$$\begin{aligned} MSE(w_j) &= \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2 \\ &= \frac{1}{n} \sum_{i=1}^n (y^{(i)} - w_j x_j^{(i)})^2 \end{aligned}$$

که برای کمینه کردن خطا داریم:

$$\begin{aligned} \frac{d}{dw_j} MSE(w_j^*) &= -\frac{1}{n} \sum_{i=1}^n (y^{(i)} - w_j^* x_j^{(i)}) x_j^{(i)} = 0 \\ \implies w_j^* \sum_{i=1}^n x_j^{(i)} x_j^{(i)} &= \sum_{i=1}^n y^{(i)} x_j^{(i)} \\ \implies w_j^* (\mathbf{x}_j^T \mathbf{x}_j) &= \mathbf{x}_j^T \mathbf{y} \\ \implies w_j^* &= \frac{\mathbf{x}_j^T \mathbf{y}}{\mathbf{x}_j^T \mathbf{x}_j} \end{aligned}$$

ب) می دانیم جواب بهینه برای بردار w اگر همه ی ویژگی ها را لحاظ کنیم، از رابطه ی زیر بدست می آید:

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

حال اگر ستون های X را با بردارهای ستونی \mathbf{x}_i نشان دهیم، آنگاه برای ماتریس $X^T X$ داریم:

$$(X^T X)_{ij} = \mathbf{x}_i^T \mathbf{x}_j$$

که اگر بدانیم ستون های X متعامد هستند، آنگاه ماتریس $X^T X$ برابر خواهد بود با:

$$(X^T X)_{ij} = \begin{cases} 0 & i \neq j \\ \mathbf{x}_i^T \mathbf{x}_i & i = j \end{cases}$$

حال برای حاصل ضرب وارون این ماتریس در ماتریس X^T خواهیم داشت:

$$X^\dagger = \begin{pmatrix} \frac{1}{\mathbf{x}_1^T \mathbf{x}_1} & 0 & \cdots & 0 \\ 0 & \frac{1}{\mathbf{x}_2^T \mathbf{x}_2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{\mathbf{x}_m^T \mathbf{x}_m} \end{pmatrix}_{m \times m} \times \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_m^T \end{pmatrix}_{m \times n} = \begin{pmatrix} \frac{\mathbf{x}_1^T}{\mathbf{x}_1^T \mathbf{x}_1} \\ \frac{\mathbf{x}_2^T}{\mathbf{x}_2^T \mathbf{x}_2} \\ \vdots \\ \frac{\mathbf{x}_m^T}{\mathbf{x}_m^T \mathbf{x}_m} \end{pmatrix}_{m \times n}$$

بنابراین برای w داریم:

$$\mathbf{w}^* = X^\dagger \mathbf{y} = \begin{pmatrix} \frac{\mathbf{x}_1^T}{\mathbf{x}_1^T \mathbf{x}_1} \\ \frac{\mathbf{x}_2^T}{\mathbf{x}_2^T \mathbf{x}_2} \\ \vdots \\ \frac{\mathbf{x}_m^T}{\mathbf{x}_m^T \mathbf{x}_m} \end{pmatrix}_{m \times n} \times \mathbf{y}_{n \times 1} = \begin{pmatrix} \frac{\mathbf{x}_1^T \mathbf{y}}{\mathbf{x}_1^T \mathbf{x}_1} \\ \frac{\mathbf{x}_2^T \mathbf{y}}{\mathbf{x}_2^T \mathbf{x}_2} \\ \vdots \\ \frac{\mathbf{x}_m^T \mathbf{y}}{\mathbf{x}_m^T \mathbf{x}_m} \end{pmatrix}_{m \times 1}$$

بنابراین همانطور که دیده می‌شود، w_j^* برابر با حالتی است که فقط از ویژگی j ام برای پیدا کردن مقدار بهینه w_j استفاده کرده بودیم.

ج) اگر از یکی از ویژگی‌ها و بایاس استفاده کرده باشیم، مدل ما به صورت $\hat{y}^{(i)} = w_0 + w_j x_j^{(i)}$ خواهد بود. بنابراین برای MSE داریم:

$$MSE(w_0, w_j) = \frac{1}{n} \sum_{i=1}^n \left(y^{(i)} - w_0 - w_j x_j^{(i)} \right)^2$$

برای پیدا کردن مقدار بهینه w_0 و w_j باید گرادیان MSE در نقاط بهینه صفر شود، بنابراین:

$$\begin{cases} \frac{\partial}{\partial w_0} MSE(w_0, w_j) = -\frac{2}{n} \sum_{i=1}^n \left(y^{(i)} - w_0 - w_j x_j^{(i)} \right) \\ \frac{\partial}{\partial w_j} MSE(w_0, w_j) = -\frac{2}{n} \sum_{i=1}^n \left(y^{(i)} - w_0 - w_j x_j^{(i)} \right) x_j^{(i)} \end{cases}$$

با صفر قراردادن روابط بالا برای w_0^* و w_j^* داریم:

$$\begin{aligned}\frac{\partial}{\partial w_0} MSE(w_0^*, w_j^*) &= 0 \\ \implies \frac{\sum_{i=1}^n y^{(i)}}{n} - w_0^* \frac{\sum_{i=1}^n 1}{n} - w_j^* \frac{\sum_{i=1}^n x_j^{(i)}}{n} &= 0 \\ \implies w_0^* &= \mathbb{E}[y] - w_j^* \mathbb{E}[x_j] \quad (\star)\end{aligned}$$

و با استفاده از مقدار بدست آمده برای w_0^* داریم:

$$\begin{aligned}\frac{\partial}{\partial w_j} MSE(w_0^*, w_j^*) &= 0 \\ \implies \frac{\sum_{i=1}^n y^{(i)} x_j^{(i)}}{n} - w_0^* \frac{\sum_{i=1}^n x_j^{(i)}}{n} - w_j^* \frac{\sum_{i=1}^n x_j^{(i)} x_j^{(i)}}{n} &= 0 \\ \implies \mathbb{E}[y x_j] - w_0^* \mathbb{E}[x_j] - w_j^* \mathbb{E}[x_j^2] &= 0 \\ \stackrel{(\star)}{\implies} \mathbb{E}[y x_j] - \mathbb{E}[y] \mathbb{E}[x_j] + w_j^* \mathbb{E}[x_j] \mathbb{E}[x_j] - w_j^* \mathbb{E}[x_j^2] &= 0 \\ \implies w_j^* &= \frac{\mathbb{E}[y x_j] - \mathbb{E}[y] \mathbb{E}[x_j]}{\mathbb{E}[x_j^2] - \mathbb{E}[x_j]^2} = \frac{\text{Cov}(y, x_j)}{\text{Var}(x_j)}\end{aligned}$$

بنابراین $\boxed{w_0^* = \mathbb{E}[y] - w_j^* \mathbb{E}[x_j]}$ و $\boxed{w_j^* = \frac{\text{Cov}(y, x_j)}{\text{Var}(x_j)}}$