



تمرین اول

پاسخ مسئله‌ی ۱.

الف) تابع فعال سازی **Softmax** اغلب در مسائل دسته‌بندی مورد استفاده قرار می‌گیرد، چرا که می‌توان مقادیر خروجی آن را به عنوان احتمال تعلق نمونه ورودی به هر یک از کلاس‌ها در نظر گرفت. تابع **Softmax** ورودی‌های خود را به مقادیر بین صفر و یک تبدیل می‌کند و مجموع خروجی‌های آن برابر با یک می‌باشد؛ بنابراین خروجی این تابع تمامی خصوصیات یک توزیع احتمالاتی معتبر بر روی کلاس‌های مختلف را دارا می‌باشد. این ویژگی باعث می‌شود که بتوان از آن برای پیش‌بینی احتمال تعلق یک نمونه به هر یک از کلاس‌ها استفاده کرد. همچنین **Softmax** یک تابع همواره مشتق پذیر است که این موضوع سبب می‌شود تا استفاده از آن به هنگام بهینه سازی تابع هزینه (هنگام آموزش مدل) مشکل ساز نباشد.

به‌طور خاص، در مدل‌های یادگیری عمیق مانند شبکه‌های عصبی، از این تابع فعال سازی در لایه آخر استفاده می‌شود تا پیش‌بینی‌های مدل به احتمالات برای کلاس‌های مختلف تبدیل شوند. اگر خروجی یک شبکه عصبی به صورت یک بردار باشد که هر عنصر آن نشان‌دهنده امتیاز (logit) برای یک کلاس است، **Softmax** این امتیازات را به احتمالات نرمال سازی می‌کند.

معادله تابع **Softmax** به شکل زیر است:

$$\text{Softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}$$

که در آن z_i امتیاز یا logit کلاس i و n تعداد کل کلاس‌ها می‌باشد.

ب) واریانس بالا به این معنی است که مدل در هنگام تغییر داده‌های آموزشی، به شدت تغییر می‌کند. به عبارتی دیگر، مدل دچار بیش برازش (overfitting) شده است، یعنی مدل به خوبی به داده‌های آموزشی پاسخ می‌دهد، اما در مواجهه با داده‌های جدید (داده‌های آزمون) عملکرد ضعیفی دارد. در چنین حالتی، مدل بیش از حد به جزئیات و نویزهای موجود در داده‌های آموزشی توجه کرده و در تعمیم به داده‌های جدید (مثلاً داده‌های موجود در مجموعه آزمون) دچار مشکل می‌شود.

یکی از روش‌های موثر برای کاهش واریانس مدل، استفاده از روش Regularization است. Regularization تکنیکی است که از پیچیدگی مدل کاسته و از بیش برازش آن جلوگیری می‌کند. یکی از رایج‌ترین و معروف‌ترین روش‌های Regularization، L_2 regularization (یا Ridge Regression) نام دارد که یک جریمه بر روی مقادیر وزن‌های مدل اعمال می‌کند و آن‌ها را کوچک‌تر نگه می‌دارد تا مدل بیش از حد به داده‌های آموزشی وابسته نشود. یکی دیگر از روش‌های Regularization که به‌ویژه برای شبکه‌های عصبی به شکل گسترده‌ای مورد استفاده قرار می‌گیرد، تکنیک Early Stopping است. در این تکنیک، بهینه سازی مدل به شکل کامل انجام نمی‌شود و بعد از تعدادی مرحله فرایند بهینه سازی را متوقف می‌کنیم تا از حفظ کردن جزئیات داده‌های آموزشی توسط مدل به منظور کاهش هر چه بیشتر تابع هزینه جلوگیری شود.

پ) در ابتدا به تفاوت میان این دو نوع رگرسیون دقت نمایید:

- **رگرسیون Ridge از L2 regularization** استفاده می‌کند که به تابع هزینه یک جریمه به شکل مجموع مربع وزن‌ها اضافه می‌کند:

$$\lambda \sum w_i^2$$

این جریمه وزن‌ها را کوچک می‌کند، اما معمولاً هیچ‌کدام از وزن‌ها را به طور کامل صفر نمی‌کند. در نتیجه همه ویژگی‌ها در مدل باقی می‌مانند ولی اثرگذاری کمتری خواهند داشت.

- **رگرسیون Lasso از L1 regularization** استفاده می‌کند و جریمه‌ای به شکل مجموع قدر مطلق وزن‌ها را به تابع هزینه اضافه می‌نماید:

$$\lambda \sum |w_i|$$

این جریمه باعث می‌شود که برخی از وزن‌ها دقیقاً به صفر برسند. بنابراین Lasso به طور خودکار ویژگی‌های غیرمهم را حذف می‌کند.

با توجه به توضیحات بالا، می‌توان گفت که اگر تمامی ویژگی‌های مدل با خروجی به‌طور قابل توجهی مرتبط باشند و حذف هیچ‌کدام از آن‌ها به کاهش خطای مدل کمک نکند، رگرسیون Ridge گزینه بهتری است. زیرا رگرسیون Lasso تمایل به صفر کردن وزن‌های برخی از ویژگی‌ها دارد که ممکن است منجر به حذف اطلاعات مهمی شود.

همچنین در بعضی موارد رگرسیون Ridge نسبت به رگرسیون Lasso پایدارتر است و در مواجهه با داده‌های جدید عملکرد بهتری دارد، به خصوص در مواقعی که داده‌ها به شدت پرنویز هستند. رگرسیون Ridge تمامی ویژگی‌ها را نگه می‌دارد و تاثیر همه را در نظر می‌گیرد، بنابراین ممکن است مدل نهایی تعمیم‌پذیری بهتری داشته باشد.

در نهایت در رگرسیون Ridge به دلیل استفاده از جریمه $\lambda \sum w_i^2$ ، تمامی وزن‌ها به شکلی متعادل کوچک می‌شوند. این باعث می‌شود که مدل به طور هموارتر رفتار کند و پیش‌بینی خطی ملایم‌تری داشته باشد. در رگرسیون Lasso به دلیل حذف کامل برخی از ویژگی‌ها، تغییرات مدل می‌تواند ناگهانی‌تر باشد که ممکن است منجر به افزایش واریانس مدل و بیش‌برازش آن شود.

ت) رگولاریزیشن L_2 با اضافه کردن یک جریمه به‌شکل مضربی از مجموع مربعات وزن‌های مدل، پیچیدگی مدل را کاهش می‌دهد. تاثیر این کار بر تعادل بایاس-واریانس به شرح زیر است:

- **کاهش واریانس:** رگولاریزیشن L_2 به کاهش واریانس مدل کمک می‌کند. بدون رگولاریزیشن، یک مدل خطی ممکن است دچار بیش‌برازش شود، به این معنی که مدل بیش از حد به داده‌های آموزشی وابسته می‌شود و نمی‌تواند روی داده‌های جدید به خوبی عمل کند. با کنترل کردن اندازه وزن‌ها و جلوگیری از افزایش بیش از اندازه آنها، پیچیدگی مدل کاهش یافته و مدل کمتر دچار تغییرات شدید با تغییرات جزئی در داده‌های آموزشی می‌شود (اصطلاحاً مدل هموارتر یا smooth تر می‌شود) و در نتیجه واریانس مدل کاهش می‌یابد.

- **افزایش بایاس:** اعمال رگولاریزیشن L_2 می‌تواند بایاس مدل را کمی افزایش دهد. زیرا محدود کردن وزن‌ها به این معنی است که مدل نمی‌تواند به اندازه کافی انعطاف‌پذیر باشد تا تمام جزئیات داده‌ها را فرا بگیرد.

در مجموع می‌توان گفت که رگولاریزیشن L_2 تعادلی بین بایاس و واریانس برقرار می‌کند. اگرچه بایاس کمی افزایش می‌یابد، اما کاهش واریانس معمولاً بهبود عملکرد مدل روی داده‌های جدید را به همراه دارد. این نوع رگولاریزیشن به طور کلی به بهبود تعمیم‌پذیری مدل کمک می‌کند.

پاسخ مسئله‌ی ۲.

الف) دقت کنید که اگر رگرسیون را فقط بر روی ویژگی z ام انجام دهیم، هدف کمینه کردن تابع هزینه زیر خواهد بود:

$$\mathcal{L} = \sum_{i=1}^n (w_j^T x_j^{(i)} - y^{(i)})^2 = \sum_{i=1}^n (x_j^{(i)T} w_j - y^{(i)})^2 = \|X_j^T w_j - y\|^2$$

چرا که در این صورت، مسئله با مسئله رگرسیون خطی زمانی که تنها یک ویژگی داریم (که در اینجا ویژگی z ام در بردار ویژگی‌ها می‌باشد) معادل خواهد شد و می‌دانیم که در آنجا نیز از میانگین مربعات خطا به عنوان تابع هزینه استفاده می‌شود.

حال دقت کنید که می‌توان تابع هزینه را به شکل زیر ساده کرد:

$$\begin{aligned}\mathcal{L} &= \|X_j^T w_j - y\|^2 \\ &= (X_j^T w_j - y)^T (X_j^T w_j - y) \\ &= ((X_j^T w_j)^T - y^T) (X_j^T w_j - y) \\ &= (w_j^T X_j - y^T) (X_j^T w_j - y) \\ &= w_j^T X_j X_j^T w_j - w_j^T X_j y - y^T X_j^T w_j + y^T y\end{aligned}$$

حال دقت کنید که داریم:

$$\begin{aligned}\langle X_j^T w_j, y \rangle &= \langle y, X_j^T w_j \rangle \\ \implies (X_j^T w_j)^T y &= y^T (X_j^T w_j) \\ \implies w_j^T X_j y &= y^T X_j^T w_j\end{aligned}$$

بنابراین خواهیم داشت:

$$\mathcal{L} = w_j^T X_j X_j^T w_j - 2w_j^T X_j y + y^T y$$

برای پیدا کردن w_j بهینه که تابع \mathcal{L} را کمینه کند، گرادینان \mathcal{L} نسبت به w_j را محاسبه کرده و آن را برابر صفر قرار می‌دهیم:

$$\frac{\partial \mathcal{L}}{\partial w_j} = (X_j X_j^T + (X_j X_j^T)^T) w_j - 2X_j y = 2X_j X_j^T w_j - 2X_j y$$

پس می‌توان نوشت: (دقت کنید که X_j یک ماتریس سطری است، در نتیجه $X_j X_j^T$ یک عدد می‌باشد و می‌توان آن را به مخرج کسر منتقل نمود. همچنین فرض می‌کنیم که هیچ یک از ویژگی‌ها تماماً صفر نیستند)

$$\frac{\partial \mathcal{L}}{\partial w_j} = 0 \implies 2X_j X_j^T w_j - 2X_j y = 0 \implies X_j X_j^T w_j = X_j y \implies w_j = \frac{X_j y}{X_j X_j^T}$$

که این همان حکم مسئله است.

البته باید به این نکته توجه نمایید که چون تابع \mathcal{L} یک تابع محدب است، در نقطه‌ای که گرادینانش صفر می‌شود به کمینه سراسری خود خواهد رسید.

ب) در ابتدا دقت کنید که تابع هزینه در این قسمت به شکل زیر می باشد:

$$\mathcal{L} = \sum_{i=1}^n (w^T x^{(i)} - y^{(i)})^2 = \sum_{i=1}^n (x^{(i)T} w - y^{(i)})^2 = \|X^T w - y\|_2^2$$

در ادامه مشابه قسمت الف تابع هزینه را ساده می کنیم:

$$\begin{aligned}\mathcal{L} &= \|X^T w - y\|_2^2 \\ &= (X^T w - y)^T (X^T w - y) \\ &= ((X^T w)^T - y^T) (X^T w - y) \\ &= (w^T X - y^T) (X^T w - y) \\ &= w^T X X^T w - w^T X y - y^T X^T w + y^T y\end{aligned}$$

حال دقت کنید که داریم:

$$\begin{aligned}\langle X^T w, y \rangle &= \langle y, X^T w \rangle \\ \implies (X^T w)^T y &= y^T (X^T w) \\ \implies w^T X y &= y^T X^T w\end{aligned}$$

بنابراین خواهیم داشت:

$$\mathcal{L} = w^T X X^T w - 2w^T X y + y^T y$$

برای پیدا کردن w بهینه که تابع \mathcal{L} را کمینه کند، گرادیان \mathcal{L} نسبت به w را محاسبه کرده و آن را برابر صفر قرار می دهیم:

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial w} &= (X X^T + (X X^T)^T) w - 2X y = 2X X^T w - 2X y \\ \frac{\partial \mathcal{L}}{\partial w} = 0 &\implies 2X X^T w - 2X y = 0 \implies X X^T w = X y\end{aligned}$$

چون در صورت سوال فرض شده است که ویژگی ها دو به دو بر هم عمودند، برای هر $1 \leq i \neq j \leq L$ تساوی $X_i X_j^T = 0$ برقرار می باشد. در نتیجه می توان نوشت:

$$X X^T = \begin{bmatrix} X_1 X_1^T & X_1 X_2^T & \dots & X_1 X_L^T \\ X_2 X_1^T & X_2 X_2^T & \dots & X_2 X_L^T \\ \vdots & \vdots & \ddots & \vdots \\ X_L X_1^T & X_L X_2^T & \dots & X_L X_L^T \end{bmatrix} = \begin{bmatrix} X_1 X_1^T & 0 & \dots & 0 \\ 0 & X_2 X_2^T & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & X_L X_L^T \end{bmatrix}$$

در نتیجه با فرض اینکه هیچ کدام از ویژگی ها تماماً صفر نیستند، $X X^T$ یک ماتریس قطری با درایه های قطری ناصفر بوده و وارون پذیر می باشد. پس خواهیم داشت:

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

از طرفی دقت کنید که داریم:

$$(XX^T)^{-1} = \begin{bmatrix} X_1 X_1^T & \cdot & \dots & \cdot \\ \cdot & X_2 X_2^T & \dots & \cdot \\ \vdots & \vdots & \ddots & \vdots \\ \cdot & \cdot & \dots & X_L X_L^T \end{bmatrix}^{-1} = \begin{bmatrix} \frac{1}{X_1 X_1^T} & \cdot & \dots & \cdot \\ \cdot & \frac{1}{X_2 X_2^T} & \dots & \cdot \\ \vdots & \vdots & \ddots & \vdots \\ \cdot & \cdot & \dots & \frac{1}{X_L X_L^T} \end{bmatrix}$$

همچنین می توان نوشت:

$$Xy = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_L \end{bmatrix} y = \begin{bmatrix} X_1 y \\ X_2 y \\ \vdots \\ X_L y \end{bmatrix}$$

در نهایت با جایگذاری تساوی های بالا در عبارت به دست آمده برای بردار w ، خواهیم داشت:

$$w = (XX^T)^{-1} Xy = \begin{bmatrix} \frac{1}{X_1 X_1^T} & \cdot & \dots & \cdot \\ \cdot & \frac{1}{X_2 X_2^T} & \dots & \cdot \\ \vdots & \vdots & \ddots & \vdots \\ \cdot & \cdot & \dots & \frac{1}{X_L X_L^T} \end{bmatrix} \begin{bmatrix} X_1 y \\ X_2 y \\ \vdots \\ X_L y \end{bmatrix} = \begin{bmatrix} \frac{X_1 y}{X_1 X_1^T} \\ \frac{X_2 y}{X_2 X_2^T} \\ \vdots \\ \frac{X_L y}{X_L X_L^T} \end{bmatrix}$$

$$\implies \forall 1 \leq j \leq L; w_j = \frac{X_j y}{X_j X_j^T}$$

و در نتیجه حکم مسئله درست می باشد.

پ) در این قسمت، تابع هزینه به شکل زیر می باشد:

$$\mathcal{L} = \sum_{i=1}^n (w_j^T x_j^{(i)} + w, - y^{(i)})^2 = \sum_{i=1}^n (x_j^{(i)T} w_j + w, - y^{(i)})^2 = \|X_j^T w_j + w, \mathcal{I} - y\|_2^2$$

به طوری که \mathcal{I} برداری n بعدی است که تمامی مولفه های آن برابر ۱ هستند.

مشابه قسمت های قبل، در ابتدا تابع هزینه را ساده می کنیم:

$$\begin{aligned} \mathcal{L} &= \|X_j^T w_j + w, \mathcal{I} - y\|_2^2 \\ &= (X_j^T w_j + w, \mathcal{I} - y)^T (X_j^T w_j + w, \mathcal{I} - y) \\ &= (w_j^T X_j + \mathcal{I}^T w, - y^T) (X_j^T w_j + w, \mathcal{I} - y) \\ &= w_j^T X_j X_j^T w_j + w_j^T X_j w, \mathcal{I} - w_j^T X_j y + \mathcal{I}^T w, \mathcal{I} X_j^T w_j + \mathcal{I}^T w, \mathcal{I} \\ &\quad - \mathcal{I}^T w, y - y^T X_j^T w_j - y^T w, \mathcal{I} + y^T y \end{aligned}$$

حال دقت کنید که مشابه آنچه که پیش تر نشان دادیم، در اینجا نیز داریم:

$$\langle X_j^T w_j, w, \mathcal{I} \rangle = \langle w, \mathcal{I}, X_j^T w_j \rangle \implies w_j^T X_j w, \mathcal{I} = \mathcal{I}^T w, X_j^T w_j$$

$$\langle X_j^T w_j, y \rangle = \langle y, X_j^T w_j \rangle \implies w_j^T X_j y = y^T X_j^T w_j$$

$$\langle w, \mathcal{I}, y \rangle = \langle y, w, \mathcal{I} \rangle \implies \mathcal{I}^T w, y = y^T w, \mathcal{I}$$

بنابراین خواهیم داشت:

$$\mathcal{L} = w_j^T X_j X_j^T w_j + \mathcal{I}^T w, \mathcal{I} + y^T y + 2w_j^T X_j w, \mathcal{I} - 2w_j^T X_j y - 2\mathcal{I}^T w, y$$

برای پیدا کردن w_j, w بهینه که تابع \mathcal{L} را کمینه کنند، گرادینان \mathcal{L} نسبت به هر کدام از این دو متغیر را محاسبه کرده و آنها را برابر صفر قرار می دهیم: (دقت کنید که چون w و w_j عدد هستند، می توان آن ها را در ضرب با ماتریس ها جابجا کرد و همچنین تساوی های $w_j^T = w_j$ و $w, = w$ برقرار می باشند)

$$\frac{\partial \mathcal{L}}{\partial w_j} = 2w_j X_j X_j^T + 2w, X_j \mathcal{I} - 2X_j y$$

$$\frac{\partial \mathcal{L}}{\partial w,} = 2w, \mathcal{I}^T \mathcal{I} + 2w_j X_j \mathcal{I} - 2\mathcal{I}^T y$$

$$\frac{\partial \mathcal{L}}{\partial w_j} = 0 \implies w_j X_j X_j^T + w, X_j \mathcal{I} = X_j y$$

$$\frac{\partial \mathcal{L}}{\partial w,} = 0 \implies w, \mathcal{I}^T \mathcal{I} + w_j X_j \mathcal{I} = \mathcal{I}^T y$$

در نهایت با حل دستگاه دو معادله دو مجهول بالا، خواهیم داشت:

$$w, = \frac{d_{\text{r}} - \frac{B d_{\text{l}}}{A}}{C - \frac{B^{\text{r}}}{A}} = \frac{d_{\text{r}} A - d_{\text{l}} B}{CA - B^{\text{r}}}$$

$$w_j = \frac{d_{\text{l}} - B \left(\frac{d_{\text{r}} - \frac{B d_{\text{l}}}{A}}{C - \frac{B^{\text{r}}}{A}} \right)}{A} = \frac{d_{\text{l}} C - d_{\text{r}} B}{CA - B^{\text{r}}}$$

به طوری که در عبارات بالا داریم:

$$A = X_j X_j^T, \quad B = X_j \mathcal{I}, \quad C = \mathcal{I}^T \mathcal{I}, \quad d_{\text{l}} = X_j y, \quad d_{\text{r}} = \mathcal{I}^T y$$

پاسخ مسئله‌ی ۳.

دقت کنید که برای قرارگیری نقطه به مختصات (x_1, x_2) در داخل و یا روی مرز ناحیه مشخص شده، تمامی نابرابری های زیر باید به صورت همزمان برقرار باشند:

$$x_2 \leq 1 \quad x_2 \geq -1$$

$$x_2 \leq 2 - x_1 \quad x_2 \leq 2 + x_1$$

$$x_2 \geq -2 + x_1 \quad x_2 \geq -2 - x_1$$

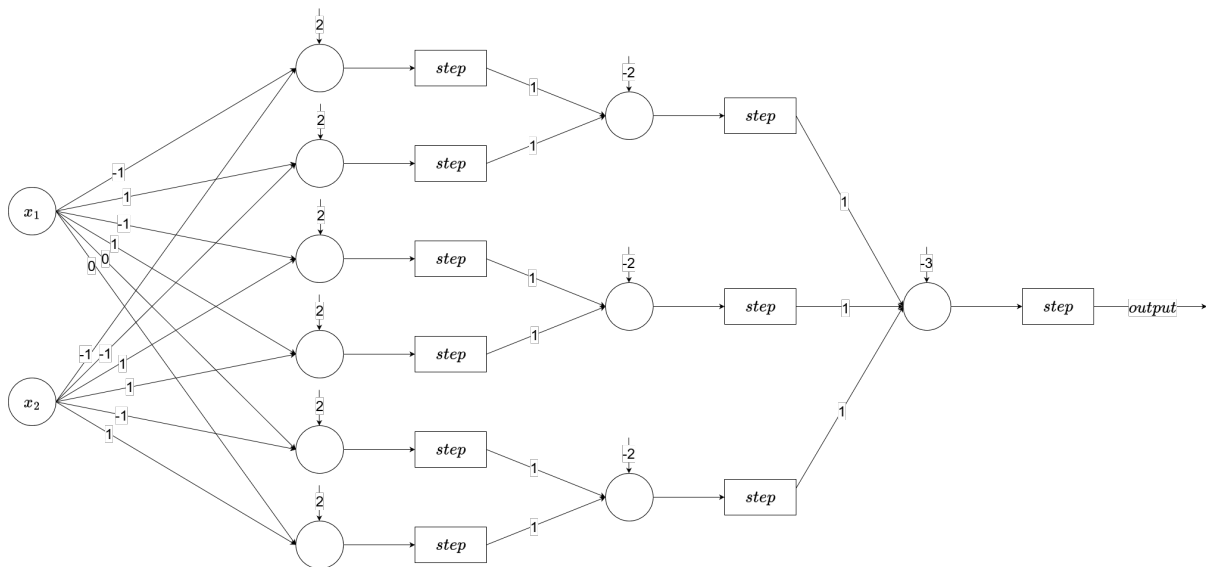
یا به طور معادل باید داشته باشیم:

$$1 - x_2 \geq 0 \quad x_2 + 1 \geq 0$$

$$2 - x_1 - x_2 \geq 0 \quad 2 + x_1 - x_2 \geq 0$$

$$x_2 + 2 - x_1 \geq 0 \quad x_2 + 2 + x_1 \geq 0$$

خروجی شبکه عصبی زیر که دارای دو لایه مخفی می باشد، تنها زمانی یک می شود که مختصات نقطه ورودی در تمامی نابرابری های بالا صدق کند و در غیر این صورت خروجی آن صفر خواهد شد.



شکل ۱: شبکه عصبی با دو لایه مخفی

دقت کنید که در طراحی شبکه بالا، تعریف تابع پله به شکل زیر در نظر گرفته شده است:

$$step(x) = \begin{cases} 1 & x \geq 0 \\ 0 & x < 0 \end{cases}$$

پاسخ مسئله‌ی ۴.

الف) در دو حالت گفته شده، مشتق تابع softmax را نسبت به z_k به دست می‌آوریم.

• حالت $k = i$:

$$\begin{aligned}
 \frac{\partial \hat{y}_k}{\partial z_k} &= \frac{\partial \frac{e^{z_k}}{\sum_{j=1}^n e^{z_j}}}{\partial z_k} \\
 &= \left[\frac{\partial e^{z_k}}{\partial z_k} \cdot \sum_{j=1}^n e^{z_j} - \frac{\partial \left(\sum_{j=1}^n e^{z_j} \right)}{\partial z_k} \cdot e^{z_k} \right] \cdot \frac{1}{\left(\sum_{j=1}^n e^{z_j} \right)^2} \\
 &= \left[e^{z_k} \cdot \sum_{j=1}^n e^{z_j} - \sum_{j=1}^n e^{z_j} \cdot e^{z_k} \right] \cdot \frac{1}{\left(\sum_{j=1}^n e^{z_j} \right)^2} \\
 &= \frac{e^{z_k}}{\sum_{j=1}^n e^{z_j}} - \frac{e^{z_k}}{\left(\sum_{j=1}^n e^{z_j} \right)} \\
 &= \frac{e^{z_k}}{\sum_{j=1}^n e^{z_j}} - \left(\frac{e^{z_k}}{\sum_{j=1}^n e^{z_j}} \right) \\
 &= \hat{y}_k - \hat{y}_k \\
 &= \hat{y}_k (1 - \hat{y}_k)
 \end{aligned}$$

• حالت $k \neq i$:

$$\begin{aligned}
 \frac{\partial \hat{y}_i}{\partial z_k} &= \frac{\partial \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}}{\partial z_k} \\
 &= \left[\frac{\partial e^{z_i}}{\partial z_k} \cdot \sum_{j=1}^n e^{z_j} - \frac{\partial \left(\sum_{j=1}^n e^{z_j} \right)}{\partial z_k} \cdot e^{z_i} \right] \cdot \frac{1}{\left(\sum_{j=1}^n e^{z_j} \right)^2} \\
 &= [-e^{z_k} \cdot e^{z_i}] \cdot \frac{1}{\left(\sum_{j=1}^n e^{z_j} \right)^2} \\
 &= -\frac{e^{z_i}}{\left(\sum_{j=1}^n e^{z_j} \right)} \cdot \frac{e^{z_k}}{\left(\sum_{j=1}^n e^{z_j} \right)} \\
 &= -\hat{y}_i \hat{y}_k
 \end{aligned}$$

ب) با توجه به مشتق های به دست آمده در قسمت قبل، خواهیم داشت:

$$\begin{aligned}
 \frac{\partial L}{\partial z_k} &= -\frac{\partial (\sum_{i=1}^n y_i \log \hat{y}_i)}{\partial z_k} \\
 &= -\frac{\partial (y_k \log \hat{y}_k)}{\partial z_k} - \frac{\partial (\sum_{i \neq k} y_i \log \hat{y}_i)}{\partial z_k} \\
 &= -\frac{\partial (y_k \log \hat{y}_k)}{\partial z_k} - \sum_{i \neq k} \frac{\partial (y_i \log \hat{y}_i)}{\partial z_k} \\
 &= -y_k \cdot \frac{\partial \log \hat{y}_k}{\partial \hat{y}_k} \cdot \frac{\partial \hat{y}_k}{\partial z_k} - \sum_{i \neq k} y_i \cdot \frac{\partial \log \hat{y}_i}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial z_k} \\
 &= -y_k \cdot \frac{1}{\hat{y}_k} \cdot \hat{y}_k (1 - \hat{y}_k) + \sum_{i \neq k} y_i \cdot \frac{1}{\hat{y}_i} \cdot \hat{y}_k \hat{y}_i \\
 &= y_k (\hat{y}_k - 1) + \sum_{i \neq k} y_i \hat{y}_k \\
 &= \sum_{i=1}^n y_i \hat{y}_k - \hat{y}_k
 \end{aligned}$$

پاسخ مسئله‌ی ۵.

الف) در ابتدا تابع هزینه را برای رگرسیون Ridge و رگرسیون خطی می نویسیم و با محاسبه گرادیان این توابع و برابر قرار دادن گرادیان آنها با صفر، پارامتر بهینه را برای هر کدام به دست می آوریم. دقت کنید که فرض می کنیم ماتریس X و بردارهای y و β به شکل زیر تعریف شده اند:

$$y = X\beta + e, \quad e \sim \mathcal{N}(\cdot, \sigma^2 I), \quad x_i \in \mathbb{R}^L, \quad y \in \mathbb{R}^n$$

$$X = [x_1, x_2, \dots, x_n]^T, \quad y = [y_1, y_2, \dots, y_n]^T, \quad \beta = [\beta_1, \beta_2, \dots, \beta_L]^T$$

همچنین فرض می کنیم که ویژگی ها مستقل از یکدیگرند (چرا که در غیر این صورت، رگرسیون خطی پاسخ نخواهد داشت).

• رگرسیون خطی:

$$\begin{aligned} \mathcal{L}_{LS} &= \|X\hat{\beta}_{LS} - y\|_2^2 \\ &= (X\hat{\beta}_{LS} - y)^T (X\hat{\beta}_{LS} - y) \\ &= \hat{\beta}_{LS}^T X^T X \hat{\beta}_{LS} - \hat{\beta}_{LS}^T X^T y - y^T X \hat{\beta}_{LS} + y^T y \\ &= \hat{\beta}_{LS}^T X^T X \hat{\beta}_{LS} - 2\hat{\beta}_{LS}^T X^T y + y^T y \end{aligned}$$

$$\begin{aligned} \frac{\partial \mathcal{L}_{LS}}{\partial \hat{\beta}_{LS}} = 0 &\implies 2X^T X \hat{\beta}_{LS} - 2X^T y = 0 \\ &\implies X^T X \hat{\beta}_{LS} = X^T y \\ &\implies \hat{\beta}_{LS} = (X^T X)^{-1} X^T y \end{aligned}$$

• رگرسیون Ridge:

$$\begin{aligned} \mathcal{L}_{Ridge} &= \|X\hat{\beta}_{Ridge} - y\|_2^2 + \lambda \|\hat{\beta}_{Ridge}\|_2^2 \\ &= (X\hat{\beta}_{Ridge} - y)^T (X\hat{\beta}_{Ridge} - y) + \lambda \hat{\beta}_{Ridge}^T \hat{\beta}_{Ridge} \\ &= \hat{\beta}_{Ridge}^T X^T X \hat{\beta}_{Ridge} - \hat{\beta}_{Ridge}^T X^T y - y^T X \hat{\beta}_{Ridge} + y^T y + \lambda \hat{\beta}_{Ridge}^T \hat{\beta}_{Ridge} \\ &= \hat{\beta}_{Ridge}^T X^T X \hat{\beta}_{Ridge} - 2\hat{\beta}_{Ridge}^T X^T y + y^T y + \lambda \hat{\beta}_{Ridge}^T \hat{\beta}_{Ridge} \end{aligned}$$

$$\begin{aligned} \frac{\partial \mathcal{L}_{Ridge}}{\partial \hat{\beta}_{Ridge}} = 0 &\implies 2X^T X \hat{\beta}_{Ridge} - 2X^T y + 2\lambda \hat{\beta}_{Ridge} = 0 \\ &\implies X^T X \hat{\beta}_{Ridge} + \lambda \hat{\beta}_{Ridge} = X^T y \\ &\implies (X^T X + \lambda I) \hat{\beta}_{Ridge} = X^T y \\ &\implies \hat{\beta}_{Ridge} = (X^T X + \lambda I)^{-1} X^T y \end{aligned}$$

در ادامه، واریانس هر کدام از این ضرایب را محاسبه خواهیم کرد.

• رگرسیون خطی:

$$\begin{aligned}\hat{\beta}_{LS} &= (X^T X)^{-1} X^T y \\ &= (X^T X)^{-1} X^T (X\beta + e) \\ &= (X^T X)^{-1} (X^T X)\beta + (X^T X)^{-1} X^T e \\ &= \beta + (X^T X)^{-1} X^T e\end{aligned}$$

$$\begin{aligned}\Rightarrow \text{Var}(\hat{\beta}_{LS}) &= \text{Var}(\beta + (X^T X)^{-1} X^T e) \\ &= \text{Var}((X^T X)^{-1} X^T e) \\ &= (X^T X)^{-1} X^T \text{Var}(e) [(X^T X)^{-1} X^T]^T \\ &= (X^T X)^{-1} X^T \text{Var}(e) X (X^T X)^{-1} \\ &= \sigma^2 (X^T X)^{-1} (X^T X) (X^T X)^{-1} \\ &= \sigma^2 (X^T X)^{-1}\end{aligned}$$

• رگرسیون Ridge: در ابتدا دقت کنید که می توان نوشت:

$$\begin{aligned}\hat{\beta}_{Ridge} &= (X^T X + \lambda I)^{-1} X^T y \\ &= (X^T X + \lambda I)^{-1} (X^T X) (X^T X)^{-1} X^T y \\ &= (X^T X + \lambda I)^{-1} X^T X \hat{\beta}_{LS}\end{aligned}$$

در نتیجه خواهیم داشت:

$$\begin{aligned}\text{Var}(\hat{\beta}_{Ridge}) &= \text{Var}((X^T X + \lambda I)^{-1} X^T X \hat{\beta}_{LS}) \\ &= (X^T X + \lambda I)^{-1} X^T X \text{Var}(\hat{\beta}_{LS}) [(X^T X + \lambda I)^{-1} X^T X]^T \\ &= (X^T X + \lambda I)^{-1} X^T X \text{Var}(\hat{\beta}_{LS}) X^T X (X^T X + \lambda I)^{-1} \\ &= (X^T X + \lambda I)^{-1} X^T X \sigma^2 (X^T X)^{-1} X^T X (X^T X + \lambda I)^{-1} \\ &= \sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1}\end{aligned}$$

حال نشان می دهیم که $\text{Var}(\hat{\beta}_{LS}) - \text{Var}(\hat{\beta}_{Ridge})$ یک ماتریس مثبت معین می باشد.

در ابتدا تعریف می کنیم

$$W = X^T X (X^T X + \lambda I)^{-1}$$

بنابراین به دست می آید:

$$\begin{aligned}\text{Var}(\hat{\beta}_{Ridge}) &= \sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1} \\ &= \sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X)^{-1} X^T X (X^T X + \lambda I)^{-1} \\ &= \sigma^2 W^T (X X^T)^{-1} W\end{aligned}$$

حال می توان نوشت:

$$\begin{aligned}
& \text{Var}(\hat{\beta}_{LS}) - \text{Var}(\hat{\beta}_{Ridge}) \\
&= \sigma^2 (X^T X)^{-1} - \sigma^2 W^T (X^T X)^{-1} W \\
&= \sigma^2 \{ W^T (W^T)^{-1} (X^T X)^{-1} W - W^T (X^T X)^{-1} W \} \\
&= \sigma^2 W^T \{ (W^T)^{-1} (X^T X)^{-1} W - (X^T X)^{-1} \} W \\
&= \sigma^2 W^T \{ (X^T X)^{-1} (X^T X + \lambda I) (X^T X)^{-1} (X^T X + \lambda I) (X^T X)^{-1} - (X^T X)^{-1} \} W \\
&= \sigma^2 W^T \{ (I + \lambda (X^T X)^{-1}) (X^T X)^{-1} (I + \lambda (X^T X)^{-1}) - (X^T X)^{-1} \} W \\
&= \sigma^2 W^T \{ ((X^T X)^{-1} + \lambda (X^T X)^{-2}) (I + \lambda (X^T X)^{-1}) - (X^T X)^{-1} \} W \\
&= \sigma^2 W^T \{ (X^T X)^{-1} + \lambda (X^T X)^{-2} + \lambda (X^T X)^{-2} + \lambda^2 (X^T X)^{-3} - (X^T X)^{-1} \} W \\
&= \sigma^2 W^T \{ 2\lambda (X^T X)^{-2} + \lambda^2 (X^T X)^{-3} \} W \\
&= \sigma^2 (X^T X + \lambda I)^{-1} X^T X \{ 2\lambda (X^T X)^{-2} + \lambda^2 (X^T X)^{-3} \} X^T X (X^T X + \lambda I)^{-1} \\
&= \sigma^2 (X^T X + \lambda I)^{-1} \{ 2\lambda I + \lambda^2 (X^T X)^{-1} \} (X^T X + \lambda I)^{-1}
\end{aligned}$$

دقت کنید که می دانیم:

$$\forall v \neq 0 : u = (X^T X + \lambda I)^{-1} v \neq 0$$

و در نتیجه خواهیم داشت:

$$\begin{aligned}
v^T [\text{Var}(\hat{\beta}_{LS}) - \text{Var}(\hat{\beta}_{Ridge})] v &= \sigma^2 u^T \{ 2\lambda I + \lambda^2 (X^T X)^{-1} \} u \\
&= \sigma^2 \lambda u^T u + \sigma^2 \lambda^2 u^T (X^T X)^{-1} u > 0
\end{aligned}$$

بنابراین $\text{Var}(\hat{\beta}_{LS}) - \text{Var}(\hat{\beta}_{Ridge})$ یک ماتریس مثبت معین می باشد و حکم درست است.

دقت کنید که $\text{trace}\{\text{Var}(\hat{\beta}_{LS}) - \text{Var}(\hat{\beta}_{Ridge})\} = \sum_i \lambda_i (\text{Var}(\hat{\beta}_{LS}) - \text{Var}(\hat{\beta}_{Ridge}))$ که برابر مجموع اختلاف واریانس ها می باشد، نیز مثبت است؛ چرا که بخاطر مثبت معین بودن این ماتریس تمامی مقادیر ویژه آن مثبت اند.

ب) در ابتدا دقت کنید که داریم:

$$\hat{Y}(\lambda) = X \hat{\beta}_{Ridge}$$

بنابراین با توجه به واریانس به دست آمده برای $\hat{\beta}_{Ridge}$ در قسمت قبل، می توان نوشت:

$$\begin{aligned}
\text{Var}[\hat{Y}(\lambda)] &= \text{Var}(X \hat{\beta}_{Ridge}) \\
&= X \text{Var}(\hat{\beta}_{Ridge}) X^T \\
&= \sigma^2 X (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1} X^T
\end{aligned}$$

حال تجزیه SVD ماتریس X را به شکل زیر در نظر بگیرید:

$$X = U \Sigma V^T$$

به طوری که U و V ماتریس هایی متعامدند (به عبارت دیگر تساوی های $U^T = U^{-1}$ و $V^T = V^{-1}$ برقرار هستند) و Σ یک ماتریس قطری است که درایه های واقع بر روی قطر اصلی در آن، مقادیر تکین X می باشند.

بنابراین داریم:

$$X^T X = (U \Sigma V^T)^T (U \Sigma V^T) = V \Sigma (U^T U) \Sigma V^T = V \Sigma^T V^T$$

در ادامه دقت کنید که با توجه به خاصیت تغییرناپذیری trace با شیفیت چرخشی در ضرب ماتریس ها، خواهیم داشت:

$$\begin{aligned} \text{trace} \left\{ \text{Var} \left[\hat{Y}(\lambda) \right] \right\} &= \text{trace} \left\{ \sigma^T X (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1} X^T \right\} \\ &= \sigma^T \text{trace} \left\{ X (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1} X^T \right\} \\ &= \sigma^T \text{trace} \left\{ X^T X (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1} \right\} \\ &= \sigma^T \text{trace} \left\{ (X^T X (X^T X + \lambda I)^{-1})^T \right\} \\ &= \sigma^T \text{trace} \left\{ (V \Sigma^T V^T (V \Sigma^T V^T + \lambda I)^{-1})^T \right\} \\ &= \sigma^T \text{trace} \left\{ (V \Sigma^T V^T (V \Sigma^T V^T + \lambda V V^T)^{-1})^T \right\} \\ &= \sigma^T \text{trace} \left\{ (V \Sigma^T V^T (V (\Sigma^T + \lambda I) V^T)^{-1})^T \right\} \\ &= \sigma^T \text{trace} \left\{ (V \Sigma^T (V^T V) (\Sigma^T + \lambda I)^{-1} V^T)^T \right\} \\ &= \sigma^T \text{trace} \left\{ (V \Sigma^T (\Sigma^T + \lambda I)^{-1} V^T)^T \right\} \\ &= \sigma^T \text{trace} \left\{ V \Sigma^T (\Sigma^T + \lambda I)^{-1} (V^T V) \Sigma^T (\Sigma^T + \lambda I)^{-1} V^T \right\} \\ &= \sigma^T \text{trace} \left\{ V \Sigma^T (\Sigma^T + \lambda I)^{-1} \Sigma^T (\Sigma^T + \lambda I)^{-1} V^T \right\} \\ &= \sigma^T \text{trace} \left\{ V^T V \Sigma^T (\Sigma^T + \lambda I)^{-1} \Sigma^T (\Sigma^T + \lambda I)^{-1} \right\} \\ &= \sigma^T \text{trace} \left\{ \Sigma^T (\Sigma^T + \lambda I)^{-1} \Sigma^T (\Sigma^T + \lambda I)^{-1} \right\} \\ &= \sigma^T \text{trace} \left\{ \Sigma^T (\Sigma^T + \lambda I)^{-2} \right\} \quad (*) \\ &= \sigma^T \text{trace} \left\{ \begin{bmatrix} \sigma_1(X)^T & \dots & \cdot \\ \vdots & \ddots & \vdots \\ \cdot & \dots & \sigma_p(X)^T \end{bmatrix} \begin{bmatrix} \sigma_1(X)^T + \lambda & \dots & \cdot \\ \vdots & \ddots & \vdots \\ \cdot & \dots & \sigma_p(X)^T + \lambda \end{bmatrix}^{-2} \right\} \\ &= \sigma^T \text{trace} \left\{ \begin{bmatrix} \sigma_1(X)^T & \dots & \cdot \\ \vdots & \ddots & \vdots \\ \cdot & \dots & \sigma_p(X)^T \end{bmatrix} \begin{bmatrix} \frac{1}{(\sigma_1(X)^T + \lambda)^2} & \dots & \cdot \\ \vdots & \ddots & \vdots \\ \cdot & \dots & \frac{1}{(\sigma_p(X)^T + \lambda)^2} \end{bmatrix} \right\} \\ &= \sigma^T \sum_{i=1}^p \sigma_i(X)^T [\sigma_i(X)^T + \lambda]^{-2} \end{aligned}$$

که این همان حکم مسئله می باشد (البته دقت کنید که در صورت سوال، باید به جای مقادیر ویژه، مقادیر تکین قرار بگیرد).