



## یادگیری ماشین

پاییز ۱۴۰۳

استاد: علی شریفی زارچی

مسئول تمرین:

مهلت ارسال نهایی: ۲۵ مهر

## فصل اول

تمرین اول

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روزهای 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.
- اگر بخش عملی یا تئوری تمرین را قبل از مهلت ارسال امتیازی آپلود کنید، ۲۰ درصد نمره اضافی به آن بخش تعلق خواهد گرفت و پس از آن، ویدئویی تحت عنوان راهنمایی برای حل تمرین منتشر خواهد شد.
- حتماً تمرین‌ها را بر اساس موارد ذکر شده در صورت سوالات حل کنید. در صورت وجود هرگونه ابهام، آن را در صفحه تمرین در سایت کوئرا مطرح کنید و به پاسخ‌هایی که از سوی دستیار آموزشی مربوطه ارائه می‌شود، توجه کنید.
- 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.
- فایل پاسخ‌های سوالات نظری را در قالب یک فایل pdf به فرمت `HW1_T_[STD_ID].pdf` آماده کنید و برای سوالات عملی، هریک را در یک فایل zip جداگانه قرار دهید و فایل zip اول را به فرمت `HW1_P1_[STD_ID].zip` و فایل zip دوم را به فرمت `HW1_P2_[STD_ID].zip` نامگذاری کرده و هرکدام را به صورت جداگانه آپلود کنید.
- گردآوردندگان تمرین: مبینا سلیمی‌پناه، محمد مولوی، امیرعلی لقمانی، فاطمه السادات موسوی، عرشیا قارونی

## سوالات نظری (۱۰۰ نمره)

۱. (۲۰ نمره) به سوالات زیر پاسخ کوتاه دهید.

- الف) چرا از تابع softmax اغلب برای مسائل دسته‌بندی استفاده می‌شود؟
- ب) بالا بودن واریانس در مدل چه معنایی دارد؟ یک روش ممکن برای کاهش واریانس در مدل خود بیان کنید.
- پ) چرا در حالتی که تمام ویژگی‌ها تا حد خوبی با خروجی مرتبط هستند، رگرسیون Ridge به رگرسیون Lasso ترجیح داده می‌شود؟
- ت) چگونه رگولاریزیشن  $L_2$  در  $classifier$ های خطی بر روی تعادل بایاس-واریانس تأثیر می‌گذارد؟
- حل.

- الف) خروجی softmax محدودیت‌های یک توزیع احتمالی را برآورده می‌کند (یا پاسخ‌های دیگری که اشاره می‌کنند خروجی softmax مقادیر بین ۰ تا ۱ تولید می‌کند که جمع آن‌ها برابر ۱ است).
- ب) این بدان معناست که مدل به داده‌های آموزش بیش از حد تطبیق یافته است (overfitting) و قابلیت تعمیم ندارد. برای بخش دوم هر پاسخی که به قابلیت تعمیم کمک کند، قابل قبول است مانند اضافه کردن داده‌های بیشتر، dropout، ساختن مدل کوچکتر و غیره.

پ) زیرا در این حالت رگولاریزیشن Lasso از بین متغیرهایی که ارتباط زیادی باهم دارند یکی را نگه می‌دارد و بقیه را حذف می‌کند در حالیکه ممکن است تمام این متغیرها در خروجی اثر بزرگی داشته باشند. در رگولاریزیشن Ridge حتی اگر متغیرها همبستگی زیادی داشته باشند تمام متغیرها حفظ می‌شوند. بنابراین امکان حذف متغیر پراهمیت وجود ندارد.

ت) رگولاریزیشن  $L_2$  با اعمال جریمه‌ای به مقادیر بزرگ وزن‌ها، پیچیدگی مدل را کنترل می‌کند. این امر از بیش‌برازش (*overfitting*) جلوگیری کرده و واریانس مدل را کاهش می‌دهد بدون اینکه بایاس به‌طور قابل‌توجهی افزایش یابد. در نتیجه، *Generalization* مدل بهبود می‌یابد.

۲. (۲۰ نمره) در یک مسئله رگرسیون خطی داریم:

$$y = \underline{w}^T \underline{x}, \quad \underline{x} \in \mathbb{R}^L, \quad y \in \mathbb{R}$$

$$X = [\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N], \quad \underline{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}, \quad \underline{w} = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_L \end{pmatrix}$$

الف) اگر رگرسیون را فقط بر روی ویژگی  $j$  انجام دهیم، نشان دهید که:

$$w_j = \frac{X_j y}{X_j X_j^T}$$

که  $X_j$  سطر  $j$  ماتریس داده‌ها است.

ب) فرض کنید ویژگی‌ها مستقل هستند (یعنی سطرهاى ماتریس داده‌ها مستقل هستند). ثابت کنید که پارامترهای بهینه از آموزش رگرسیون بر روی همه ویژگی‌ها با پارامترهای بهینه حاصل از آموزش روی هر ویژگی به‌طور مستقل یکسان است.

پ) فرض کنید  $y = \underline{w}^T \underline{x} + w$  و رگرسیون را فقط بر روی ویژگی  $j$  انجام دهیم.  $w$  و  $w_j$  را بدست آورید. حل.

الف) تخمین‌گر OLS برای بردار وزن به صورت معادله زیر بیان می‌شود:

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

در سوال مورد نظر، رگرسیون خطی ساده با یک متغیر پیش‌بینی‌کننده بررسی شده است. در این حالت، ماتریس  $X$  شامل یک بردار ستونی  $X_j$  است و وزن  $w_j$  برای متغیر پیش‌بینی‌کننده  $X_j$  به صورت زیر قابل محاسبه است:

$$w_j = \frac{X_j y}{X_j X_j^T}$$

ب) در این حالت، ماتریس کوواریانس  $X^T X$  یک ماتریس قطری است، چرا که عناصر غیر قطری که معرف کوواریانس بین ویژگی‌های مختلف هستند، به دلیل فرض استقلال صفر هستند. بنابراین، معکوس  $X^T X$  برابر است با معکوس هر عنصر قطری.

$$X = \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_N \end{bmatrix}$$

$$X^T = \begin{bmatrix} 1 & & & \\ \sigma_1^2 & & & \\ & \ddots & & \\ & & \sigma_N^2 & \\ & & & 1 \end{bmatrix}$$

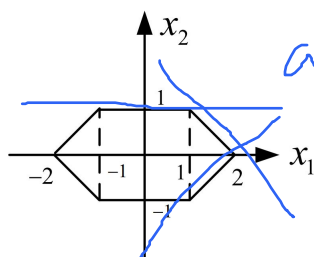
وزن برای هر ویژگی  $w_j$  می‌تواند به صورت مستقل با استفاده از فرمول بخش قبلی محاسبه شود. این به دلیل فرض استقلال است که نشان می‌دهد اطلاعات مشترکی بین ویژگی‌ها وجود ندارد که بر محاسبه وزن‌ها تاثیر بگذارد.

(پ)

$$w_0 = \bar{y} = E\{y\}, \quad w_j = \frac{X_j(y - w_0)}{X_j X_j^T}$$

در این بخش،  $w_0$  برابر با میانگین  $y$  است و  $w_j$  به صورت فوق محاسبه می‌شود.

۳. (۲۰ نمره) یک شبکه عصبی با دو گره در لایه ورودی و دولایه مخفی و تابع فعالسازی پله داریم. وزن‌ها و بایاس‌های این شبکه را به گونه‌ای تعیین کنید که در ناحیه داخل ۶ ضلعی و روی اضلاع خروجی شبکه ۱ و در باقی نواحی ۰ باشد.



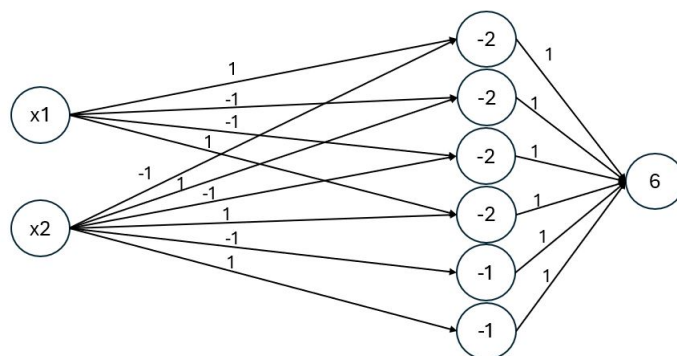
$$w_0 + w_1 x_1 + w_2 x_2 \geq 1$$

حل.

باتوجه به شکل داده شده می‌توانیم نامساوی‌های زیر را بنویسیم:

$$x_1 - x_2 \geq -2, x_2 - x_1 \geq -2, -x_2 \geq -1, x_2 \geq -1, -x_1 - x_2 \geq -2, x_1 + x_2 \geq -2$$

با توجه به اینکه تابع فعال ساز هم پله می‌باشد، می‌توان جواب زیر را برای این سوال در نظر گرفت:



۴. (۲۰ نمره) در یک مسئله دسته‌بندی دو کلاسه (binary classification) از رگرسیون لاجستیک با تابع هزینه cross entropy استفاده کرده‌ایم. تابع هزینه برای یک نقطه از داده‌ها به صورت زیر تعریف می‌شود:

$$L = - \sum_{i=1}^n y_i \log(\hat{y}_i)$$

که در آن  $y$  بردار برجسب‌های واقعی و  $\hat{y}$  احتمالات پیش‌بینی شده با استفاده از تابع softmax به صورت زیر است:

$$\hat{y}_i = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}$$

الف) مشتق تابع softmax را نسبت به  $z_k$  برای دو حالت  $k = i$  و  $k \neq i$  به دست آورید.

ب) با استفاده از مشتق تابع softmax، مشتق تابع هزینه cross entropy را نسبت به  $z_k$  محاسبه نمایید. حل.

الف) برای محاسبه مشتق  $\hat{y}_i$  نسبت به  $z_k$  دو حالت بررسی می‌کنیم:

حالت ۱:  $i = k$  در این حالت، هدف ما محاسبه مشتق  $\hat{y}_i$  نسبت به  $z_i$  (یعنی  $\frac{\partial \hat{y}_i}{\partial z_i}$ ) است.

$$\hat{y}_i = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}$$

$$\frac{\partial \hat{y}_i}{\partial z_i} = \frac{\partial}{\partial z_i} \left( \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}} \right)$$

با مشتق‌گیری از کسر داریم:

$$\frac{\partial \hat{y}_i}{\partial z_i} = \frac{e^{z_i} \left( \sum_{j=1}^n e^{z_j} \right) - e^{z_i} e^{z_i}}{\left( \sum_{j=1}^n e^{z_j} \right)^2}$$

که ساده می‌شود به:

$$\frac{\partial \hat{y}_i}{\partial z_i} = \hat{y}_i (1 - \hat{y}_i)$$

حالت ۲:  $i \neq k$  در این حالت، هدف محاسبه مشتق  $\hat{y}_i$  نسبت به  $z_k$  است (یعنی  $\frac{\partial \hat{y}_i}{\partial z_k}$  برای  $i \neq k$ ).

$$\hat{y}_i = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}$$

با مشتق‌گیری از تابع کسری، داریم:

$$\frac{\partial \hat{y}_i}{\partial z_k} = \frac{0 \cdot \left( \sum_{j=1}^n e^{z_j} \right) - e^{z_i} e^{z_k}}{\left( \sum_{j=1}^n e^{z_j} \right)^2}$$

که نتیجه می‌دهد:

$$\frac{\partial \hat{y}_i}{\partial z_k} = -\hat{y}_i \hat{y}_k$$

ب) برای محاسبه مشتق تابع هزینه Cross-Entropy نسبت به  $z_k$ ، از قاعده مشتق زنجیره‌ای استفاده می‌کنیم. قاعده مشتق زنجیره‌ای:

$$\frac{\partial L}{\partial z_k} = \sum_{i=1}^n \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial z_k}$$

با توجه به تعریف تابع هزینه:

$$L = - \sum_{i=1}^n y_i \log(\hat{y}_i)$$

مشتق نسبت به  $\hat{y}_i$  به این صورت است:

$$\frac{\partial L}{\partial \hat{y}_i} = -\frac{y_i}{\hat{y}_i}$$

بنابراین، می‌توانیم رابطه زیر را بنویسیم:

$$\frac{\partial L}{\partial z_k} = \sum_{i=1}^n \left( -y_i \frac{1}{\hat{y}_i} \right) \cdot \frac{\partial \hat{y}_i}{\partial z_k}$$

اکنون از مشتق تابع Softmax که در بخش الف به دست آوردیم، استفاده می‌کنیم: برای  $i = k$  داریم:

$$\frac{\partial \hat{y}_k}{\partial z_k} = \hat{y}_k (1 - \hat{y}_k)$$

برای  $i \neq k$  داریم:

$$\frac{\partial \hat{y}_i}{\partial z_k} = -\hat{y}_i \hat{y}_k$$

اکنون می‌توانیم دو قسمت را با هم ترکیب کنیم:

$$\frac{\partial L}{\partial z_k} = -y_k (1 - \hat{y}_k) + \sum_{i \neq k} y_i \hat{y}_k \quad (۱)$$

از آنجایی که  $\sum_{i \neq k} y_i = 1 - y_k$ ، معادله به شکل زیر ساده می‌شود:

$$\frac{\partial L}{\partial z_k} = -y_k (1 - \hat{y}_k) + \hat{y}_k (1 - y_k) \quad (۲)$$

که در نهایت ساده می‌شود به:

$$\frac{\partial L}{\partial z_k} = \hat{y}_k - y_k \quad (۳)$$

نتیجه نهایی: این نشان می‌دهد که مشتق تابع هزینه Cross-Entropy به سادگی برابر است با تفاوت احتمال پیش‌بینی شده ( $\hat{y}_k$ ) و برچسب واقعی ( $y_k$ ):

$$\frac{\partial L}{\partial z_k} = \hat{y}_k - y_k \quad (۴)$$

۵. (۲۰ نمره) به سوالات زیر در مورد رگرسیون Ridge پاسخ دهید:

الف) نشان دهید که به ازای مقادیر  $\lambda > 0$  واریانس ضرایب Ridge از واریانس ضرایب رگرسیون خطی کوچکتر است:

$$\text{Var}(\hat{\beta}^{LS}) > \text{Var}(\hat{\beta}^{Ridge}(\lambda))$$

ب) نشان دهید به ازای مقادیر  $\lambda > 0$  رابطه زیر برقرار است:

$$\text{tr} \left\{ \text{Var}[\hat{Y}(\lambda)] \right\} = \sigma^2 \sum_{j=1}^p (D_x)_{jj}^2 [(D_x)_{jj}^2 + \lambda]^{-2}$$

در این رابطه  $D_x$  یک ماتریس قطری است که شامل مقادیر تکین  $X$  است.

حل. الف) از تجزیه SVD ماتریس  $X$  استفاده می‌کنیم.

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

$$\beta_{LS} = (X^T X)^{-1} X^T y \Rightarrow \text{Cov}(\hat{\beta}_{LS}) = E \left[ (\hat{\beta}_{LS} - E(\hat{\beta})) (\hat{\beta}_{LS} - E(\hat{\beta}))^T \right]$$

ثابت

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

$$\text{Cov}(\beta_{ls}) = \text{Cov}((X^T X)^{-1} X^T \epsilon)$$

$$\text{Cov}(\beta_{ls}) = E \left[ \left( (X^T X)^{-1} X^T \epsilon - E \left[ (X^T X)^{-1} X^T \epsilon \right] \right) \left( (X^T X)^{-1} X^T \epsilon - E \left[ (X^T X)^{-1} X^T \epsilon \right] \right)^T \right]$$

$$\text{Cov}(\beta_{ls}) = E \left[ ((X^T X)^{-1} X^T \epsilon) (\epsilon^T X ((X^T X)^{-1})^T) \right]$$

$$\xrightarrow{X=UDV^T} \text{cov}(\vec{\beta}_{ls}) = E[(UD^T V^T)^{-1} V D U^T \epsilon \epsilon^T U D V^T (V D^T V^T)^{-1}] \rightarrow E[\epsilon \epsilon^T] = \sigma^2 I$$

$$\Rightarrow \sigma^2 ((V D^T V^T)^{-1})^T = \sigma^2 \sum_{j=1}^p D_{jj}^{-2} V_j V_j^T$$

$$\beta_{ridge} = (X^T X + \lambda I)^{-1} X^T y \rightarrow \text{cov}(\vec{\beta}_{ridge}) = \sigma^2 (V D^T V^T + \lambda I)^{-1} V D^T V^T ((V D^T V^T + \lambda I)^{-1})^T$$

$$= \sigma^2 \sum_{j=1}^p \frac{D_{jj}^2}{(D_{jj}^2 + \lambda)^2} V_j V_j^T$$

با توجه به اینکه واریانس هریک از درایه‌های بردار  $\beta$  برابر با درایه متناظر روی قطر  $cov(\hat{\beta})$  است داریم:

$$var(\beta_{jls}) = \frac{\sigma^2}{D_{jj}^2}$$

$$var(\beta_{jridge}) = \sigma^2 \frac{D_{jj}^2}{(D_{jj}^2 + \lambda)^2}$$

$$\xrightarrow{\lambda \rightarrow \infty} (D_{jj}^2 + \lambda)^2 = D_{jj}^4 + \lambda^2 + 2\lambda D_{jj}^2 > D_{jj}^2$$

(ب)

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

$$\hat{Y} = X\beta + \epsilon$$

$$X = U D V^T$$

$$cov(\epsilon\epsilon^T) = \sigma^2 I$$

$$var[\hat{Y}(\lambda)] = E[(\hat{Y} - E\hat{Y})(\hat{Y} - E\hat{Y})^T] \Rightarrow E[UD(D^2 + \lambda I)^{-1} DU^T \epsilon \epsilon^T U D ((D^2 + \lambda I)^{-1})^T DU^T]$$

$$= \sigma^2 U D (D^2 + \lambda I)^{-1} D^2 ((D^2 + \lambda I)^{-1})^T D U^T$$

$$(D^2 + \lambda I)^{-1} = \begin{bmatrix} \frac{1}{(D_1^2 + \lambda)} & 0 & 0 & 0 \\ 0 & \frac{1}{(D_2^2 + \lambda)} & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & \frac{1}{(D_P^2 + \lambda)} \end{bmatrix}_{P \times P}$$

$$D^2 = \begin{bmatrix} D_1^2 & 0 & 0 & 0 \\ 0 & D_2^2 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & D_P^2 \end{bmatrix}_{P \times P}$$

$$var(\hat{Y}) = \sigma^2 U \begin{bmatrix} \frac{D_1^2}{(D_1^2 + \lambda)^2} & 0 & 0 & 0 \\ 0 & \frac{D_2^2}{(D_2^2 + \lambda)^2} & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & \frac{D_P^2}{(D_P^2 + \lambda)^2} \end{bmatrix}_{P \times P} U^T$$

$$\rightarrow tr(var(\hat{Y})) = \sigma^2 \sum_{j=1}^P \frac{D_j^2}{(D_j^2 + \lambda)^2}$$

۱. (۱۰۰ نمره) برای حل سوالات به *notebook* های ضمیمه شده مراجعه کنید.

(۱) (۶۰ نمره) برای پاسخ به تمرین عملی اول ابتدا فایل نوت‌بوک قرار گرفته را باز کنید و سپس مراحل را مطابق آنچه که از شما خواسته شده انجام دهید. در نهایت، مقادیر پیش‌بینی شده برای دیتاست *InsuranceData\_test.csv* را در یک فایل به نام *submission.csv* که شامل یک ستون به نام *charges* می‌باشد، ذخیره کنید. فایل خروجی و فایل نوت‌بوک را در یک فایل *zip* قرار دهید و آن را به فرمت *HW۱\_P۱\_[STD\_ID].zip* نامگذاری کرده و آپلود کنید. توجه بفرمایید این سوال دارای داوری خودکار می‌باشد و ۱۵ نمره از ۶۰ نمره به این قسمت تعلق دارد.

(۲) (۴۰ نمره) برای پاسخ به تمرین عملی دوم تنها کافی است نوت‌بوک *Perceptron.ipynb* را تکمیل کرده و سپس مطابق با فرمت ذکر شده آپلود کنید.

حل.  
*notebook* های حل شده ضمیمه شده‌اند.