



یادگیری ماشین

پاییز و زمستان ۱۴۰۴

استاد: علی شریفی زارچی

گرددآورندگان: گرامی‌راد - پیمایی - بهزاداصل - کوچک‌نیا - رستمیان

مهلت ارسال: ۲۶ آذر

شبکه های عصبی

تمرین سوم

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.
- درنظر داشته باشید ددلاین این تمرین شب قبل از آزمون میان‌ترم می‌باشد و امکان استفاده از تاخیر مجاز برای این تمرین وجود ندارد.
- همکاری و هم‌فکری شما در انجام تمرین مانعی ندارد اما پاسخ‌های هر کس حتماً باید توسط خود او نوشته شده باشد.
- 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.
- لطفاً تصویری واضح از پاسخ‌های سوالات نظری بارگذاری کنید. در غیر این صورت پاسخ شما تصحیح نخواهد شد.

سوالات نظری (۱۰۰ نمره)

۱. (۸ نمره)

بخش اول: گزاره‌های زیر صحیح است یا خیر؟ مختصر توضیح دهید.

- آ. اگر به تعداد کافی پرسپترون و لایه داشته باشیم، می‌توانیم هر تابع پیوسته‌ای را با دقت دلخواه تقریب بزنیم.
- ب. تابع هزینه در یک شبکه عصبی عمیق، به دلیل وجود توابع فعال‌سازی غیرخطی و ساختار چند لایه، همواره یک تابع غیرمحدب است.
- پ. یک شبکه عصبی که تنها یک لایه دارد و از تابع فعال‌سازی Sigmoid برای خروجی، به همراه تابع هزینه Cross-entropy برای دسته‌بندی دودویی استفاده می‌کند، کاملاً با مدل Logistic regression یکسان است.
- ت. پارامترهای گاما (γ) و بتا (β) در لایه نرمال‌سازی دسته‌ای، هاپرپارامتر هستند که باید قبل از شروع آموزش به صورت دستی تنظیم شوند و در حین Backpropagation به‌روزرسانی نمی‌شوند.

بخش دوم: به سوالات زیر پاسخ کوتاه دهید.

- آ. چرا مقداردهی اولیه صفر برای تمام وزن‌های یک شبکه عصبی ایده بدی است و باعث چه مشکلی در فرآیند یادگیری می‌شود؟
- ب. چرا توابع فعال‌سازی در لایه‌های مخفی یک شبکه عصبی باید غیرخطی باشند؟ اگر همه آن‌ها خطی باشند چه اتفاقی می‌افتد؟
- پ. مکانیزم بهینه‌ساز «مومنتوم» (Momentum) را با روش «کاهش شیب تصادفی» (SGD) استاندارد مقایسه کنید. مومنتوم برای حل چه مشکلی در فرآیند بهینه‌سازی معرفی شده است؟
- ت. روش «حذف تصادفی»^۱ چگونه کار می‌کند و نقش اصلی آن به عنوان یک روش تنظیم‌گری^۲ چیست؟ چرا این روش فقط در مرحله آموزش فعال است و در مرحله آزمون باید غیرفعال شود؟

Dropout^۱
Regularization^۲

۲. (۱۲ نمره) یک شبکه عصبی با ۲ ورودی، ۲ نورون در لایه مخفی^۳ با تابع فعال‌سازی ReLU و یک نورون خروجی با تابع فعال‌سازی Sigmoid داده شده است. پارامترها به صورت زیر هستند:

$$W^{(1)} = \begin{bmatrix} 1 & -1 \\ 0.5 & 0.5 \end{bmatrix}, \quad b^{(1)} = \begin{bmatrix} 0 \\ 0.5 \end{bmatrix}, \quad W^{(2)} = [1 \quad -1/5], \quad b^{(2)} = 0.2$$

برای ورودی $x = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$ و برجسب $y = 1$ ، موارد زیر را انجام دهید:

الف. یک مرحله پیش‌روی کامل^۴ را اجرا کرده و مقادیر $h, a^{(1)}, z^{(2)}$ و \hat{y} را به دست آورید. (در اینجا h خروجی فعال‌شده لایه مخفی است).

ب. با استفاده از تابع هزینه^۵ Cross-Entropy که در زیر تعریف شده است، گرادیان‌های هزینه نسبت به تمام پارامترها را محاسبه کرده و مقادیر عددی آن‌ها را به دست آورید.

$$\mathcal{L}(\hat{y}, y) = -[y \ln \hat{y} + (1 - y) \ln(1 - \hat{y})]$$

$$\text{گرادیان‌های مورد نیاز: } \frac{\partial \mathcal{L}}{\partial b^{(2)}}, \frac{\partial \mathcal{L}}{\partial W^{(2)}}, \frac{\partial \mathcal{L}}{\partial b^{(1)}}, \frac{\partial \mathcal{L}}{\partial W^{(1)}}$$

پ. (۱) تعداد کل پارامترهای قابل آموزش برای یک شبکه عصبی تماماً متصل^۶ با n نورون در لایه ورودی، یک لایه مخفی با k نورون، و m نورون در لایه خروجی را به صورت یک فرمول کلی بنویسید. (۲) با استفاده از فرمول به دست آمده، تعداد کل پارامترهای قابل آموزش برای شبکه تعریف شده در این سوال را محاسبه کنید.

ت. از لحاظ مفهومی توضیح دهید که افزودن یک لایه مخفی جدید به یک شبکه عصبی، چه تأثیری بر ظرفیت مدل می‌گذارد و چرا این کار احتمال بیش‌برازش^۷ را افزایش می‌دهد؟

۳. (۲۵ نمره) تابع سافت‌مکس (Softmax) این خاصیت مطلوب را دارد که یک توزیع احتمال خروجی می‌دهد و اغلب به عنوان تابع فعال‌سازی در بسیاری از شبکه‌های دسته‌بندی چندکلاسه استفاده می‌شود. یک شبکه عصبی ۲-لایه برای دسته‌بندی K -کلاسه، با استفاده از فعال‌سازی سافت‌مکس و تابع هزینه Cross-Entropy مطابق تعریف زیر در نظر بگیرید:

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$a^{[1]} = \text{LeakyReLU}(z^{[1]}, \alpha = 0.01)$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$\bar{y} = \text{softmax}(z^{[2]})$$

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

که در آن مدل ورودی x به شکل $1 \times D_x$ و برجسب One-hot $y \in \{0, 1\}^K$ را می‌گیرد. فرض کنید لایه پنهان دارای D_a گره است، یعنی $z^{[1]}$ یک بردار با اندازه $1 \times D_a$ است. به یاد بیاورید که تابع سافت‌مکس به صورت زیر محاسبه می‌شود:

$$\hat{y} = \begin{bmatrix} \exp(z_1^{[2]})/Z \\ \vdots \\ \exp(z_K^{[2]})/Z \end{bmatrix}, \quad Z = \sum_{j=1}^K \exp(z_j^{[2]})$$

Hidden Layer^۳

Pass Forward^۴

Loss Function^۵

Fully Connected Neural Network^۶

Overfit^۷

آ. شکل (ابعاد) $W^{[۲]}$ و $b^{[۲]}$ چیست؟ اگر بر روی m مثال برداری سازی می کردیم، یعنی از یک دسته^۸ نمونه $X \in \mathbb{R}^{D_x \times m}$ به عنوان ورودی استفاده می کردیم، شکل خروجی لایه پنهان چه می بود؟

ب. مشتق $\frac{\partial \hat{y}_k}{\partial z_k^{[۲]}}$ چیست؟ پاسخ خود را بر حسب مؤلفه (های) \hat{y} ساده کنید.

پ. به ازای $i \neq k$ ، مشتق $\frac{\partial \hat{y}_k}{\partial z_i^{[۲]}}$ چیست؟ پاسخ خود را بر حسب مؤلفه (های) \hat{y} ساده کنید.

ت. فرض کنید بر حسب y درایه k -ام خود مقدار ۱ و در سایر درایه ها مقدار ۰ دارد. $\frac{\partial L}{\partial z_i^{[۲]}}$ چیست؟ پاسخ خود را بر حسب \hat{y}_i ساده کنید. نکته: هر دو حالت $i = k$ و $i \neq k$ را در نظر بگیرید.

ث. مشتق $\frac{\partial z^{[۲]}}{\partial a^{[۱]}}$ چیست؟ این نتیجه را با نماد δ_1 نشان دهید.

ج. مشتق $\frac{\partial a^{[۱]}}{\partial z^{[۱]}}$ چیست؟ این نتیجه را با نماد δ_2 نشان دهید. می توانید از نماد آکولاد (تابع چندضابطه ای) استفاده کنید.

چ. مشتق $\frac{\partial L}{\partial z^{[۲]}}$ را با δ نشان دهید. $\frac{\partial L}{\partial b^{[۱]}}$ و $\frac{\partial L}{\partial W^{[۱]}}$ چیست؟ می توانید از نمادهای بخش های قبلی استفاده مجدد کنید. نکته: به اشکال (ابعاد) دقت کنید.

ح. برای جلوگیری از مشکلات پایداری عددی، می توان از یک ترفند هنگام پیاده سازی تابع سافت مکس استفاده کرد. فرض کنید $m = \max_{j=1}^K z_j^{[۲]}$ حداکثر مقادیر $z_i^{[۲]}$ باشد، سپس \hat{y}_i به صورت زیر محاسبه می شود:

$$\hat{y}_i = \frac{\exp(z_i^{[۲]} - m)}{\sum_{j=1}^K \exp(z_j^{[۲]} - m)}.$$

مشکل عددی در محاسبه اولیه سافت مکس چیست؟ چرا فرمول تغییر یافته به حل آن مشکل کمک می کند؟
خ. آیا اصلاح انتقال-حداکثر (max-shift) در بخش آخر، مشتق نسبت به $z^{[۲]}$ را تغییر می دهد؟

۴. (۱۲ نمره) یک شبکه عصبی برای پیاده سازی تابع XOR برای چهار ورودی دودویی $x_1, x_2, x_3, x_4 \in \{0, 1\}$ طراحی شده است. این شبکه از توابع فعال سازی زیر استفاده می کند:

$$\text{ReLU}(z) = \max(0, z) \quad , \quad \sigma(z) = \frac{1}{1 + e^{-z}}$$

خروجی به عنوان ۱ در نظر گرفته می شود اگر $y > 0.5$ و در غیر این صورت ۰ است. تابع هزینه، خطای مربع است:

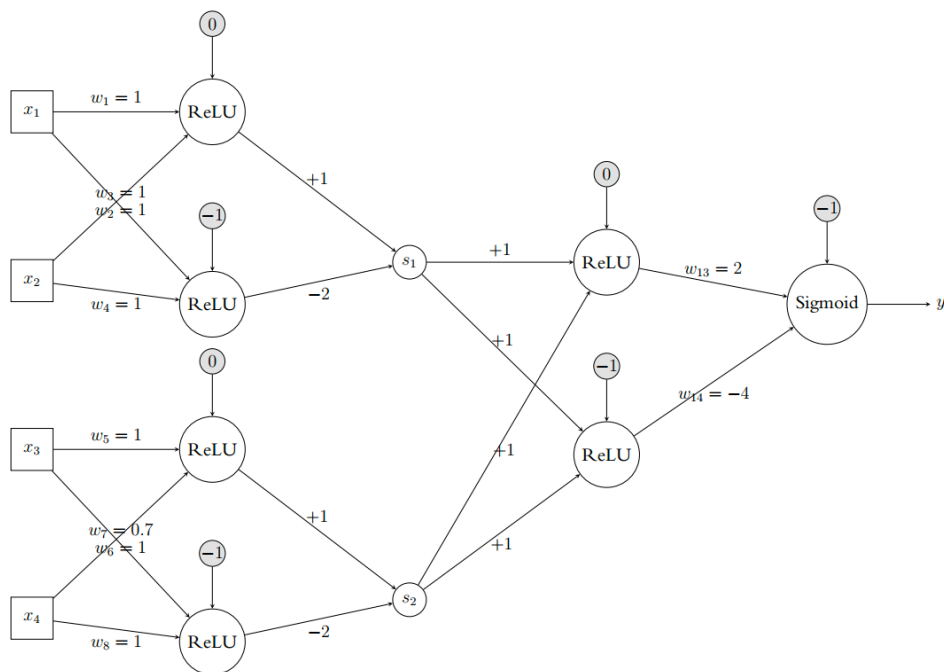
$$L(x) = (y_{\text{هدف}} - y)^2.$$

ساختار شبکه

$$\begin{aligned} h_1 &= \text{ReLU}(x_1 + x_2), \\ h_2 &= \text{ReLU}(x_1 + x_2 - 1), \\ h_3 &= \text{ReLU}(x_3 + x_4), \\ h_4 &= \text{ReLU}(0.7x_3 + x_4 - 1). \\ s_1 &= h_1 - 2h_2, \\ s_2 &= h_3 - 2h_4. \\ h_5 &= \text{ReLU}(s_1 + s_2), \\ h_6 &= \text{ReLU}(s_1 + s_2 - 1). \\ y &= \sigma(2h_5 - 4h_6 - 1). \end{aligned}$$

Batch^۸

نمودار شبکه



شکل ۱: نمودار شبکه عصبی مورد استفاده در سوال

- آ. تمام بردارهای ورودی $x = (x_1, x_2, x_3, x_4)$ که در آن مدل پیش‌بینی نادرست انجام می‌دهد (یعنی $y_{\text{مدل}} \neq y_{\text{هدف}}$) را شناسایی کنید.
- ب. یک گام SGD روی وزن $w_7 = 0.7$ با استفاده از بردار ورودی $(x_1, x_2, x_3, x_4) = (0, 0, 1, 1)$ انجام دهید. نرخ یادگیری $\eta = 0.1$ را فرض کنید. مقدار جدید w_7 را تعیین کنید.

۵. (۲۰ نمره) ماتریس طراحی X (که سطر i -ام آن نقطه نمونه \mathbf{x}_i^T است) و یک بردار n -تایی از برچسب‌ها $\mathbf{y} = [y_1, \dots, y_n]^T$ به شما داده شده است. برای سادگی، فرض کنید X سفید (whitened) شده است، به طوری که $X^T X = nI$. جمله بایاس یا بعد اضافی اضافه نکنید؛ برای ورودی صفر، خروجی همیشه صفر است. فرض کنید \mathbf{x}_{*i} ستون i -ام ماتریس X را نشان می‌دهد.

آ. نشان دهید که تابع هزینه برای کمترین مربعات با تنظیم گر ℓ_1 ، (Lasso) یعنی:

$$J_1(\mathbf{w}) \triangleq \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_1 \quad (\lambda > 0 \text{ در آن})$$

می‌تواند به صورت زیر بازنویسی شود:

$$J_1(\mathbf{w}) = \|\mathbf{y}\|^2 + \sum_{i=1}^d f(\mathbf{x}_{*i}, w_i)$$

که در آن $f(\cdot, \cdot)$ یک تابع مناسب است که آرگومان اول آن یک بردار و آرگومان دوم آن یک اسکالر است.

ب. با استفاده از پاسخ خود به بخش (آ)، شرایط لازم و کافی برای مولفه i -ام بهینه‌ساز \mathbf{w}^* تابع $J_1(\cdot)$ را استخراج کنید به طوری که هر یک از سه خاصیت زیر را برآورده کند: $w_i^* > 0$ ، $w_i^* = 0$ و $w_i^* < 0$.

پ. برای بهینه‌ساز $\mathbf{w}^\#$ تابع هزینه کمترین مربعات با تنظیم گر ℓ_2 ، (Ridge) یعنی:

$$J_2(\mathbf{w}) \triangleq \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_2^2 \quad (\lambda > 0 \text{ در آن})$$

یک شرط لازم و کافی برای $w_i^\# = 0$ استخراج کنید، که در آن $w_i^\#$ مولفه i -ام بردار $\mathbf{w}^\#$ است.

ت. یک بردار را **تُنک**^۹ می‌نامیم اگر بیشتر مولفه‌های آن صفر باشند. بر اساس پاسخ‌های خود به بخش‌های (ب) و (پ)، کدام یک از بردارهای \mathbf{w}^* و $\mathbf{w}^\#$ به احتمال زیاد **تُنک**‌تر است؟ چرا؟

۶. (۱۰ نمره) تابع هزینه یک مدل $J(\theta)$ دارای حداقل‌های محلی است و می‌خواهیم از الگوریتم کاهش شیب تصادفی (SGD) برای بهینه‌سازی پارامترهای آن (θ) استفاده کنیم. فرض کنید از یک نرخ یادگیری با کاهش نمایی^{۱۰} به شکل زیر استفاده می‌کنیم:

$$\alpha_t = \alpha_0 e^{-kt}$$

که در آن t شماره گام (epoch) است. نرخ یادگیری اولیه $\alpha_0 = 0.1$ و ضریب کاهش $k = 0.05$ داده شده است.

آ. قاعده به‌روزرسانی پارامترها (θ) را با استفاده از SGD بنویسید.

ب. به صورت مفهومی توضیح دهید که چرا استفاده از نرخ یادگیری متغیر با کاهش^{۱۱} می‌تواند هم به سرعت همگرایی و هم به پایداری فرآیند آموزش کمک کند.

پ. اگر الگوریتم در یک ناحیه با شیب تند از تابع هزینه قرار گیرد، رفتار آن را با نرخ یادگیری ثابت در مقایسه با نرخ یادگیری دارای کاهش تحلیل و مقایسه کنید. این دو رویکرد چه تفاوتی در مسیر رسیدن به نقطه بهینه دارند؟

^۹ sparse
^{۱۰} exponential decay
^{۱۱} decaying learning rate

ت. یکی از شرایط کلاسیک برای تضمین همگرایی SGD (تحت فروض مشخص)، شرایط رابینز-مونرو^{۱۲} برای نرخ یادگیری α_t است:

$$\sum_{t=1}^{\infty} \alpha_t = \infty \quad \text{و} \quad \sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

بررسی کنید که آیا نرخ یادگیری با کاهش نمایی ($\alpha_t = \alpha \cdot e^{-kt}$) در این شرایط صدق می‌کند یا خیر. برای تحلیل می‌توانید از تقریب انتگرال برای سری‌ها استفاده کنید. نتیجه‌ای که به دست می‌آورید چه مفهومی در مورد تضمین همگرایی این نوع نرخ یادگیری دارد؟

۷. (۱۳ نمره) فرض کنید یک لایه نرمال‌سازی دسته‌ای^{۱۳} یک‌بعدی روی یک Mini-batch با اندازه m اعمال می‌شود. ورودی‌های اسکالر برای یک نورون خاص را با x_1, \dots, x_m نمایش می‌دهیم. لایه BN آماره‌ها و خروجی‌ها را به صورت زیر محاسبه می‌کند:

$$\begin{aligned} \mu_B &= \frac{1}{m} \sum_{i=1}^m x_i \\ \sigma_B^2 &= \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2 \\ \hat{x}_i &= \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \\ y_i &= \gamma \hat{x}_i + \beta \end{aligned}$$

که در آن γ و β پارامترهای قابل یادگیری هستند و ϵ یک عدد کوچک ثابت برای پایداری عددی است. تابع هزینه L تابعی از خروجی‌های y_1, \dots, y_m است و گرادیان نسبت به خروجی‌ها $g_i := \frac{\partial L}{\partial y_i}$ می‌باشد.

آ. با استفاده از قاعده زنجیره‌ای، فرمول بسته‌ای برای گرادیان‌های $\frac{\partial L}{\partial \beta}$ و $\frac{\partial L}{\partial \gamma}$ بر حسب g_i و \hat{x}_i به دست آورید.

ب. گرادیان $\frac{\partial L}{\partial x_i}$ را برای هر i محاسبه کنید. توجه داشته باشید که μ_B و σ_B^2 توابعی از تمام x_j ها هستند. پاسخ نهایی باید تنها بر حسب $\sigma_B^2, \mu_B, x_i, \gamma, \epsilon, g_j$ باشد. هدف رسیدن به رابطه زیر است:

$$\frac{\partial L}{\partial x_i} = \frac{\gamma}{\sqrt{\sigma_B^2 + \epsilon}} \left(g_i - \frac{1}{m} \sum_{j=1}^m g_j - \frac{\hat{x}_i}{m} \sum_{j=1}^m g_j \hat{x}_j \right)$$

راهنمایی: ابتدا $\frac{\partial L}{\partial \hat{x}_i}$ را بر حسب g_i و γ بنویسید، سپس گرادیان‌ها را از طریق \hat{x}_i, μ_B و σ_B^2 به عقب منتقل کنید.

پ. با استفاده از رابطه‌ای که در قسمت (ب) به دست آوردید، به صورت تحلیلی ثابت کنید که:

$$\sum_{i=1}^m \frac{\partial L}{\partial x_i} = 0$$

در یک یا دو جمله توضیح دهید که این ویژگی چه مفهوم شهودی‌ای در مورد رفتار لایه BN نسبت به شیفت دادن تمام ورودی‌های دسته با یک مقدار ثابت (مثلاً $x_i \leftarrow x_i + c$) دارد.

Robbins-Monro conditions^{۱۲}
Batch Normalization (BN)^{۱۳}