



تمرین شماره ۳

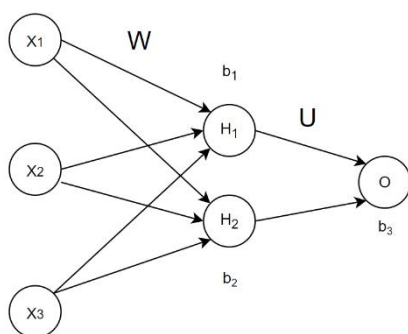
- قبل از شروع تمرین، فایل مربوط به قوانین حل و تحویل تمرین‌ها را مطالعه کنید.
- سؤالات و مشکلات خود را درباره این تمرین می‌توانید در گروه تلگرامی درس یا با طراحان این تمرین مطرح کنید.
- نوشتن گزارش کامل و تفسیر نتایج اجباری است. جزئیاتی مانند روش‌های مورد استفاده، تاثیر هر روش در نتیجه نهایی و بهبود حاصل شده به همراه ارائه معیارهای ارزیابی خواسته شده در گزارش ضروری است. با هر تغییر و هر بهبود، تغییر مقادیر معیارهای ارزیابی نیز ذکر شود. ضمناً برای گزارش سؤالات عملی، نیاز به فایل دیگری نیست و در همان محیط کدزنی، در قالب بلاک‌های مارک‌داون، توضیحات مورد نیاز را قرار دهید.
- پاسخ سؤالات تئوری را در یک فایل پی‌دی‌اف^۱ با عنوان گزارش قرار دهید.
- برای سؤالات عملی، نوتبوک‌های داده شده را تکمیل نمایید.
- طراحان این تمرین: [امیرحسین ایزدی](#) – [امیررضا افتخاری](#)

5 + 35 نمره

سؤالات تئوری

سؤال ۱ (۱۴ نمره)

یکی از اساسی‌ترین اجزا در یک شبکه ژرف، انتشار به عقب^۲ است. شبکه عصبی پیش‌رونده^۳ با معماری زیر را در نظر بگیرید:



- لایه اول: لایه تماماً متصل^۴ با تابع فعالساز سیلو^۵
- لایه دوم: لایه تماماً متصل با تابع فعالساز سیگموئید^۶

^۱ PDF

^۲ Backpropagation

^۳ Feedforward

^۴ Fully Connected

^۵ SiLU

^۶ Sigmoid

اگر وزن‌های ابتدایی و ورودی شبکه مقادیر زیر باشند و همچنین تابع خسارت^۷، تابع Cross-Entropy و بهینه‌ساز ما در این مسئله گرادیان کاهشی^۸ به همراه مومتوم^۹ و پارامترهای η و γ به ترتیب ۰.۸ و ۰.۹ باشند، به موارد زیر پاسخ دهید.

نمونه‌های ورودی:

وزن‌های ابتدایی شبکه:

$$\begin{aligned} w_{11} &= 0.4 & b_1 &= 0 \\ w_{12} &= -0.3 & b_2 &= -0.4 \\ w_{21} &= 0.2 & b_3 &= 0.1 \\ w_{22} &= 0.2 & u_1 &= 0.3 \\ w_{31} &= -0.5 & u_2 &= -0.1 \\ w_{32} &= 0.1 \end{aligned}$$

$$x^{(1)} = \begin{bmatrix} 1 \\ 0 \\ 0.5 \end{bmatrix}, y^{(1)} = 1$$

$$x^{(2)} = \begin{bmatrix} 0.1 \\ 0.5 \\ 1 \end{bmatrix}, y^{(2)} = 0$$

$$x^{(3)} = \begin{bmatrix} 0 \\ 1 \\ 0.7 \end{bmatrix}, y^{(3)} = 0$$

الف) گراف محاسباتی^{۱۰} را رسم کنید.

ب) عملیات انتشار رو به جلو^{۱۱} و انتشار به عقب را برای دو بار بروزرسانی وزن‌ها انجام دهید و وزن‌های جدید به همراه میزان تابع خطا را پس از بروزرسانی وزن‌ها گزارش کنید. توجه داشته باشید تمامی محاسبات شما باید بصورت ماتریسی انجام شود در غیر این صورت نمره‌ای از این سوال کسب نخواهید کرد.

سؤال ۲ (۱۲ نمره)

به سوالات زیر پاسخ دهید.

الف) در شکل زیر مسیرهای بهینه‌سازی برای روش‌های Momentum و Nestrov-Momentum و GD^{۱۲} و RMSProp برای یک تابع مرتبه ۲ از نقطه شروع $[-2, 2]$ رسم شده است. نمودار آبی مربوط به روش GD است.

⁷ Loss Function

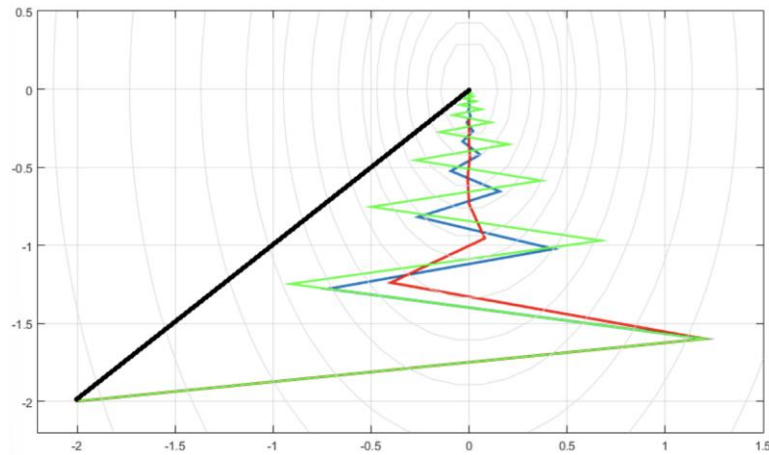
⁸ Gradient Descent

⁹ Momentum

¹⁰ Computational Graph

¹¹ Forward Propagation

¹² Gradient Descent



- با ذکر دلیل کافی توضیح دهید سه نمودار قرمز، مشکی و سبز هر کدام مربوط به کدام یک از سه روش دیگر است؟ می‌توانید برای بررسی دقیق‌تر نحوه‌ی عملکرد این الگوریتم‌ها، به این [لینک](#) نیز مراجعه کنید.
- مزایا و معایب سه روش دیگر را بیان کنید و بگویید چگونه این روش‌ها مشکلات GD را حل می‌کنند؟
- الگوریتم AdaGrad را در نظر بگیرید. ضمن مقایسه آن با RMSProp بطور جامع، توصیف کنید اگر مسیر این الگوریتم را در شکل بالا قصد داشتیم حدس بزنیم، چگونه بود؟

ب) با توجه به گرادین محاسبه شده در یک نقطه، بهینه ساز Adam سه مرحله مجزا دارد: اول به روزرسانی میانگین متحرک؛ دوم اعمال تصحیح بایاس و سوم به‌روزرسانی پارامترها. میانگین متحرک مربع گرادین‌ها را بصورت بازگشتی که در ادامه داده شده است در نظر بگیرید:

$$s_t = \beta_2 s_{t-1} + (1 - \beta_2) g_t^2$$

- عبارت s_t را فقط بر حسب گرادین‌های g_0, g_1, \dots, g_t بنویسید.
- با توجه به عبارت داده شده در قسمت قبل، $E[s_t]$ را بر حسب $E[g_t^2]$ و β_2 بنویسید. می‌توانید فرض کنید که g_i ها از هم مستقل و هم‌توزیع هستند. فرمول‌های زیر ممکن است مفید باشند:

$$\sum_{i=0}^{n-1} (a + id) = \frac{n}{2} (2a + (n-1)d)$$

$$\sum_{i=0}^{n-1} ar^i = a \frac{r^n - 1}{r - 1}$$

- با استفاده از نتیجه گیری خود در قسمت قبل، توضیح دهید که اگر مرحله تصحیح بایاس را انجام نمی‌دادید، چه اتفاقی می‌افتاد؟
- الگوریتم Adam را با RMSProp مقایسه کنید. هر یک چه مزیت و معایبی نسبت به همدیگر دارند؟

سؤال ۳ (۸ نمره)

به سوالات زیر در رابطه با نرمالسازی دسته‌ای^{۱۳} پاسخ دهید.

الف) فرمول زیر را در نظر بگیرید و توضیح دهید اگر به جای y ، \hat{x} را در نظر بگیریم چه مشکلی بوجود می‌آید؟

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

ب) توضیح دهید زمانی که سائز دسته‌ها کوچک باشد، استفاده از نرمالسازی دسته‌ای چه تأثیری در فرآیند آموزش دارد؟ استفاده از نرمالسازی دسته‌ای با اندازه دسته ۱ و مقدار اولیه شیفتم نرمالسازی دسته‌ای برابر با صفر، چه مشکلی دارد؟

ج) شما در حال اعمال نرمالسازی دسته‌ای یک لایه کاملاً متصل با اندازه ورودی ۱۰ و اندازه خروجی ۲۰ هستید. با احتساب پارامترهای نرمالسازی دسته‌ای، این لایه چند پارامتر آموزشی دارد؟

د) نرمالسازی لایه‌ای^{۱۴} تکنیک نرمالسازی دیگری است که برای غلبه بر معایب نرمالسازی دسته‌ای طراحی شده است. فرمول نرمالسازی لایه‌ای به صورت زیر است:

$$\hat{x} = \frac{x - \mu_L}{\sqrt{\sigma_L^2 + \epsilon}} \quad y = \gamma \hat{x} + \beta$$

$$\mu_L = \frac{1}{d} \sum_{j=1}^d x_j \quad \sigma_L = \sqrt{\frac{1}{d} \sum_{j=1}^d (x_j - \mu_L)^2}$$

با توجه به تعریف نرمالسازی لایه‌ای، بزرگ‌ترین تفاوت بین نرمالسازی دسته‌ای و لایه‌ای چیست؟ در چه شرایطی نرمالسازی دسته‌ای نسبت به نرمالسازی لایه‌ای ترجیح داده می‌شود؟ و برعکس آن چگونه؟ چرا؟

¹³ Batch Normalization

¹⁴ Layer Normalization

سؤال ۴ (۶ نمره)

در مورد دراپاوت^{۱۵} به موارد زیر پاسخ دهید.

- الف) توضیح دهید چرا دراپ اوت مانند منظم‌ساز عمل می‌کند؟
- ب) آیا می‌توان گفت دراپاوت عملکردی شبیه یادگیری جمعی^{۱۶} دارد؟
- ج) یک شبکه عصبی با N نود را در نظر بگیرید که هر کدام از نودها می‌توانند در طول آموزش بصورت مستقل با احتمال $0 < p < 1$ حذف شوند. تعداد کل مدل‌های منحصر به فردی که می‌توان با اعمال دراپاوت تحقق بخشید، چقدر است؟
- د) استفاده از دراپاوت هنگام آموزش و آزمایش چه تفاوتی با هم دارد؟ چرا؟

15 + 65 نمره

سوالات عملی

سوالات عملی در نوتبوک‌های پیوست شده به تمرین قابل مشاهده هست.

سالم و موفق باشید.

¹⁵ Dropout

¹⁶ Ensemble Learning