



تمرین سری اول - درس یادگیری ژرف

فاطمه براتی ۴۰۲۳۲۴۵۵۰۳

۱- الگوریتم Perceptron را که در کلاس دیدیم، برای طبقه‌بندی داده‌های زیر اعمال کنید (تا زمان همگرا شدن به‌روزرسانی بردار وزن را ادامه دهید).

$$(4,3,3) \in N \quad (2,-2,3) \in P \quad (1,0,-3) \in P \quad (4,2,2) \in N$$

همانطور که می‌بینید چهار داده داریم که دو داده متعلق به کلاس منفی (N) و دو داده متعلق به کلاس مثبت (P) هستند. فرض کنید بردار وزن اولیه در الگوریتم Perceptron برابر $(w_0, w_1, w_2, w_3) = (1,0,0,0)$ باشد، که w_0 بیانگر بایاس است.

الف- مراحل الگوریتم و مقدار بردار وزن را در هر مرحله از بروزرسانی الگوریتم به‌طور کامل بنویسید.

$A: (4, 3, 3) \in N$ $B: (2, -2, 3) \in P$
 $C: (4, 2, 2) \in N$ $D: (1, 0, -3) \in P$

$(w_0, w_1, w_2, w_3) = (1, 0, 0, 0)$
 \downarrow
 b

-1

Net Input
 $y_{in} = w \cdot x = w_0 + \sum_{i=1}^3 w_i \cdot x_i$

output
 $y = f(y_{in}) = \begin{cases} 1 & \text{if } w \cdot x > 0 \\ 0 & \text{if } w \cdot x = 0 \\ -1 & \text{if } w \cdot x < 0 \end{cases}$

از این جا که y و t متفاوت اند پس w باید update شود.

Update: $\alpha = \text{learning rate} = 1$
 $\Delta w_i = \alpha t x_i \quad i = 1, 2, 3$
 $\Delta b = \alpha t$
 $w_{i_{new}} = w_{i_{old}} + \Delta w_i$
 $b_{new} = b_{old} + \Delta b$

(الف)

epoch	Inputs			Target (t)	y_{in}	y	weight changes				weights			
	x_1	x_2	x_3				Δb	Δw_1	Δw_2	Δw_3	b	w_1	w_2	w_3
	4	3	3	-1	1	1	-1	-4	-3	-3	0	-4	-3	-3
	2	-2	3	1	-11	-1	1	2	-2	3	1	-2	-5	0
	4	2	2	-1	-17	-1	0	0	0	0	1	-2	-5	0
	1	0	-3	1	-1	-1	1	1	0	-3	2	-1	-5	-3

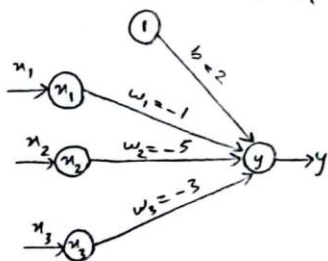
چون داتن برده w با update نشدند، پس به iteration دیگر می‌رویم.

epoch 2	Inputs			Target (t)	y_{in}	y	weight changes				weights			
	x_1	x_2	x_3				Δb	Δw_1	Δw_2	Δw_3	b	w_1	w_2	w_3
	4	3	3	-1	-26	-1	0	0	0	0	2	-1	-5	-3
	2	-2	3	1	1	1	0	0	0	0	2	-1	-5	-3
	4	2	2	-1	-18	-1	0	0	0	0	2	-1	-5	-3
	1	0	-3	1	10	1	0	0	0	0	2	-1	-5	-3

چون داتن برده w با update نشدند، y و t مطابقت داشت، پس نیاز به به ابله داتن

Final weights = $(2, -1, -5, -3)$
 $b \quad w_1 \quad w_2 \quad w_3$

iteration مانیت



ب- چه موقع مطمئن می‌شوید که الگوریتم Perceptron در این مسئله همگرا شده است؟ پس از چند مرحله روزرسانی این اتفاق رخ داده است؟

ب) زمانی که در چند مرحله y برابر شود و وزن ها $update$ نشوند و تغییر ندهند
برابر به منبرست، ممکن شویم که الگوریتم Perceptron در این مسئله همگرا شده است. در این مسئله پس از ۳ مرحله به روزرسانی، الگوریتم Perceptron همگرا نشود.

ج- کد پایتون الگوریتم Perceptron بالا را در پایتون بنویسید. برای پیاده سازی این الگوریتم تنها به محاسبات ساده نیاز دارید و می‌توانید برای راحتی بیشتر از Numpy یا PyTorch نیز استفاده کنید.

کد نوشته شده به نام Perceptron Learning Algorithm.py در پوشه code موجود است.

۲- یکی از الگوریتم‌های بهینه‌سازی مورد بررسی در کلاس الگوریتم بهینه‌سازی Momentum بود.
الف- با توجه به مباحث کلاس و نیز با جستجو در اینترنت دلیل برتری این روش نسبت به الگوریتم پایه Gradient descent را توضیح دهید.

Momentum اغلب بهتر از Gradient Descent در نظر گرفته می‌شود زیرا می‌تواند برخی از محدودیت‌های GD را رفع کند.
علل برتری Momentum:

۱. همگرایی سریع‌تر در موارد خاص:

- Momentum به تسریع همگرایی از طریق مناطق flat (مسطح) یا با انحنای کم کمک می‌کند. هنگامی که گرادیان نوین دارد یا تابع هدف دارای نواحی flat است، Momentum به فرآیند بهینه‌سازی اجازه می‌دهد تا به طور مؤثرتری از این مناطق عبور کند.

- در مقابل، GD می‌تواند در مناطق flat گیر کند یا به دلیل شیب‌های نوین نوسان کند.

۲. غلبه بر مینییم‌های محلی:

- Momentum الگوریتم بهینه‌سازی را قادر می‌سازد تا به طور مؤثرتری از مینییم‌های محلی فرار کند. با جمع‌آوری گرادیان‌های گذشته، اینرسی ایجاد می‌کند و می‌تواند از مینییم‌های محلی کم عمق فرار کند.

- GD، بدون Momentum، ممکن است به مینییم محلی کمتر از حد مطلوب همگرا شود.

۳. هموارسازی مسیر بهینه‌سازی:

- Momentum با در نظر گرفتن گرادیان‌های تاریخی مسیر بهینه‌سازی را هموار می‌کند. نوسانات را کم می‌کند و مسیری پایدارتر به سمت مینییم فراهم می‌کند.

- GD می‌تواند رفتار نامنظم از خود نشان دهد، به خصوص زمانی که گرادیان‌ها دارای نوین هستند.

۴. تجسم:

- تصور کنید توپی را از یک تپه می‌غلطانید: Momentum به توپ اجازه می‌دهد تا هنگام غلتیدن در سراسیمگی سرعت بگیرد و احتمال گیر کردن آن در یک دره کم عمق را کاهش دهد.

- GD، بدون Momentum، مانند یک توپ که مستقیماً بر اساس گرادیان می‌غلطد رفتار می‌کند، که می‌تواند به مسیرهای غیربهبوده منجر شود.

۵. اجرای عملی:

- در عمل، بسیاری از الگوریتم‌های بهینه‌سازی مانند Adam و RMSProp را برای بهبود همگرایی ترکیب می‌کنند.

- محققان Momentum را در انجام وظایف مختلف یادگیری ماشین موثر می‌دانند.

ب- تفاوت روش Nesterov Momentum با روش Momentum چیست و کدامیک نسبت به دیگری برتری دارد؟
در این قسمت نیز با توجه به مباحث کلاس و جستجو در اینترنت پاسخ دهید.

۱. Momentum:

- الگوریتم:

- در GD مبتنی بر Momentum، از به‌روزرسانی وزن قبلی m و گرادیان فعلی g برای به‌روزرسانی پارامترها p استفاده می‌کنیم.

- قانون به‌روزرسانی این است: $v = \beta m - \eta g$ ، که در آن β یک ثابت است، η نرخ یادگیری است، و v نشان‌دهنده سرعت است.

- مقدار پارامتر جدید به صورت زیر بدست می‌آید: $p_{new} = p + v$.

- رفتار:

- Momentum شیب‌های گذشته را جمع می‌کند و به مسیرهای بهینه‌سازی هموارتر منجر می‌شود.

- می‌تواند به طور موثرتری از مینیمم‌های محلی فرار کند.

- به دلیل Momentum انباشته شده سریعتر حرکت می‌کند.

- تجسم:

- تصور کنید توپی را از یک تپه می‌غلطانید: Momentum به توپ اجازه می‌دهد تا هنگام غلتیدن در سراسیمگی سرعت بگیرد و احتمال گیر کردن آن در یک دره کم عمق را کاهش دهد.

۲. Nesterov Momentum (NAG):

- الگوریتم:

- به جای محاسبه گرادیان برای پارامترهای فعلی W ، Momentum نستروف گرادیان ها را برای $(W - \beta * v_{t-1})$ محاسبه می کند، که در آن v_{t-1} سرعت قبلی است.

- قانون به روز رسانی این است: $p_{new} = p + \beta v - \eta g$

- رفتار:

- Nesterov Momentum رویکردی «gamble-correct» دارد:

- ابتدا در جهت سرعت (مثل قمار) قدم می گذارد.

- سپس بردار سرعت را بر اساس مکان جدید تصحیح می کند.

- هدف آن این است که به جای دنبال کردن سرعت به مقدار واقعی برآورد شده نزدیکتر شود.

- به بهبود سرعت همگرایی SGD کمک می کند.

به طور خلاصه، Nesterov Momentum جهت به روز رسانی را با در نظر گرفتن گرادیان در موقعیت "آینده" تنظیم می کند که منجر به رفتار همگرایی بهتر می شود. در حالی که هر دو روش گرادیان نزولی را بهبود می دهند، Nesterov Momentum مزایای بیشتری از نظر سرعت همگرایی و فرار از مینیمم های محلی ارائه می کند.

۳- یک شبکه CNN دارای سه لایه کانولوشنی، سه لایه fully connected، و یک لایه max pooling است. همچنین یک لایه فعالسازی بعد از هر لایه کانولوشنی و هر لایه fully connected قرار گرفته است. ورودی شبکه تصاویر رنگی با ابعاد 200×200 است. در مورد این شبکه داریم:

ترتیب لایه ها به این صورت است:

Conv1 >> Activation1 >> Conv2 >> Activation2 >> Max-pooling >> Conv3 >> Activation3 >> FC1 >> Activation4 >> FC2 >> Activation5 >> FC3 >> Activation6

که Conv بیانگر لایه کانولوشنی و FC بیانگر لایه Fully connected است. لایه های کانولوشنی اول تا سوم به ترتیب دارای 10، 4، و 20 فیلتر هستند. تمام فیلترها را با ابعاد 3×3 فرض کنید. بدیهی است عدد ۳ بیانگر عرض و ارتفاع فیلتر است و عمق فیلتر به عمق داده ورودی به لایه مرتبط است. padding برای لایه های کانولوشنی اول و دوم برابر ۱ و برای لایه کانولوشنی سوم برابر صفر است. Stride برای لایه کانولوشنی اول برابر ۲ و برای لایه های کانولوشنی دوم و سوم برابر ۱ است. شبکه فقط دارای یک لایه max pooling است و فیلترهای این لایه 2×2 هستند و stride آن نیز برابر ۲ است. تابع فعالسازی همه لایه ها sigmoid است. تعداد نورون های لایه های FC1، FC2، و FC3 را به ترتیب برابر ۲۰۰، ۱۰۰، و ۱ فرض کنید. به موارد زیر پاسخ دهید:

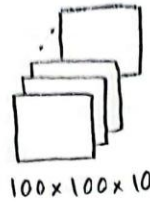
ابعاد داده سه بعدی در خروجی تمام لایه ها را بدست آورید. یعنی بیان کنید در خروجی هر لایه، داده با چه ابعادی داریم؟ (به ازای یک تصویر رنگی 200×200 به عنوان ورودی شبکه) اگر فرض کنیم هر فیلتر و نیز هر نورون علاوه بر وزن هایش یک مقدار بایاس نیز دارد، تعداد کل پارامترهای این شبکه (تعداد کل وزن ها و بایاس ها) را بدست آورید.

$$\# \text{ of output pixels} = \frac{n + 2p - f}{s} + 1$$

-14



Conv Layer 1
10 filters
filter size = 3
stride = 2
padding = 1



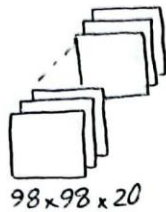
$$\text{parameters} = 10 \times (3 \times 3 \times 3 + 1) = 280$$

Conv Layer 2
4 filters
filter size = 3
stride = 1
padding = 1



$$\text{parameters} = 4 \times (3 \times 3 \times 10 + 1) = 364$$

Conv Layer 3
20 filters
filter size = 3
stride = 1
padding = 0



$$\text{parameters} = 20 \times (3 \times 3 \times 4 + 1) = 1820$$

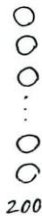
Max Pooling Layer

1 filter
filter size = 2
stride = 2

No additional parameters
(only downsampling)

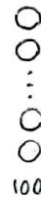


FC Layer 1
Neurons = 200



$$\text{parameters} = 49 \times 49 \times 20 \times 200 + 200 = 968200$$

FC Layer 2
Neurons = 100



$$\text{parameters} = 200 \times 100 + 100 = 20100$$

FC Layer 3

Neurons = 1



$$\text{parameters} = 100 \times 1 + 1 = 101$$

Total parameters in the entire network:

$$280 + 364 + 1820 + 968200 + 20100 + 101 = 990865$$