

محمدحسین مازندرانیان - ۸۳۰۴۰۲۰۶۶ - تمرین شماره ۲ شبکه های عصبی

$$X = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \\ 1 & 1 \end{pmatrix} \rightarrow K = X^T X = \begin{pmatrix} 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{pmatrix}$$

سوال (۱)

$$\alpha_{val} \rightarrow \max \alpha_1 + \alpha_2 + \alpha_3 + \alpha_4 - \frac{1}{2} \left(\alpha_1^2 K_{11} + \alpha_2^2 K_{22} + \alpha_3^2 K_{33} + \alpha_4^2 K_{44} + 2\alpha_1\alpha_2 K_{12} + 2\alpha_1\alpha_3 K_{13} + 2\alpha_1\alpha_4 K_{14} + 2\alpha_2\alpha_3 K_{23} + 2\alpha_2\alpha_4 K_{24} + 2\alpha_3\alpha_4 K_{34} \right)$$

$$\alpha_1 + \alpha_2 - \alpha_3 - \alpha_4 = 0 \quad \alpha_i \geq 0 \xrightarrow{\text{فرض کنیم}} \alpha_1 = 0.5, \alpha_2 = 0.5, \alpha_3 = \alpha_4 = 0$$

$$w = \sum_{i=1}^n \alpha_i y_i x_i \rightarrow w = 0.5 \cdot (0) + 0.5 \cdot (1) - 0.5 \cdot (1) - 0.5 \cdot (1) = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$b = y_k - w^T x_k \rightarrow b = 1 - (-3 \cdot 1 + -3 \cdot 0) = 4$$

$$w^T x + b = 0 \rightarrow -3x_1 - 3x_2 + 4 = 0$$

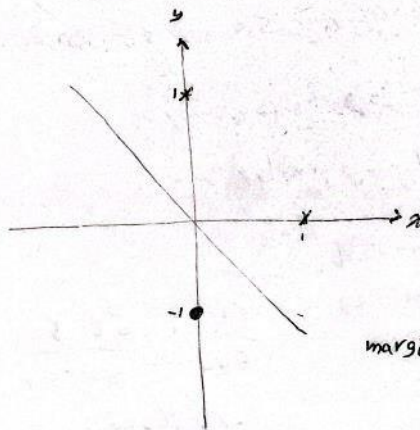
$$y_i (w^T x_i + b) = 1 \rightarrow \text{نقاط } (1,0) \text{ و } (0,1)$$

$$\alpha_1 = \frac{1 - \epsilon \alpha_4 - 3\alpha_2}{1}$$

$$\alpha_2 = \frac{1 - \epsilon \alpha_4 - 3\alpha_1}{1}$$

$$\alpha_3 = \frac{1 - 3\alpha_1 - 3\alpha_2 - \epsilon \alpha_4}{1}$$

$$\alpha_4 = \frac{1 - \epsilon \alpha_1 - \epsilon \alpha_2 - 3\alpha_3}{3}$$



حل بصورت سفردن

$$x_1 + x_2 = 0$$

وزن $(1,1)$ of bias

$$\text{margin} = \frac{|w_1 x_1 + w_2 x_2 + b|}{\sqrt{w_1^2 + w_2^2}} = 0.5 \sqrt{2}$$

سوال (۲)

$$\begin{aligned}
 (الف) \quad z_1 &= w_1 \cdot x = \begin{pmatrix} 0 \\ 0 \end{pmatrix} & z_2 &= w_2 \cdot u_1 = \begin{pmatrix} 0.5 \\ -0.5 \end{pmatrix} \\
 a_1 &= \delta(z_1) = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix} & a_2 &= \delta(z_2) = \begin{pmatrix} 0.4225 \\ 0.4775 \end{pmatrix} \\
 z_3 &= w_3 \cdot a_1 = \begin{pmatrix} -0.4225 \\ 0.4775 \end{pmatrix} & \delta_3 &= a_3 - y = \begin{pmatrix} 0.4225 \\ 0.4775 \end{pmatrix} - \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} -0.5775 \\ 0.4775 \end{pmatrix} \\
 a_3 &= \begin{pmatrix} 0.4225 \\ 0.4775 \end{pmatrix} & \delta^2 &= (w_3^T \cdot \delta_3) \odot \delta'(z_2) \\
 \delta'(z_2) &= \delta(z_2) \odot (1 - \delta(z_2)) = \begin{pmatrix} 0.4225 \\ 0.4775 \end{pmatrix} \odot \begin{pmatrix} 0.5775 \\ 0.5225 \end{pmatrix} = \begin{pmatrix} 0.2441 \\ 0.2501 \end{pmatrix} \\
 \delta^2 &= \begin{pmatrix} -1 & 0 \\ 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} -0.5775 \\ 0.4775 \end{pmatrix} \odot \begin{pmatrix} 0.2441 \\ 0.2501 \end{pmatrix} = \begin{pmatrix} 0.1245 \\ 0.1250 \end{pmatrix} \\
 \Delta w_3 &= -x \cdot \delta_3 \cdot a_1^T = -0.5 \cdot \begin{pmatrix} 0.1245 \\ 0.1250 \end{pmatrix} \cdot \begin{pmatrix} 0.5 & 0.5 \end{pmatrix} = \begin{pmatrix} -0.06225 \\ -0.0625 \end{pmatrix} \\
 w_{3new} &= w_3 + \Delta w_3 = \begin{pmatrix} 0.93775 & -0.06225 \\ -0.06225 & -1.0625 \end{pmatrix}
 \end{aligned}$$

نسباً) توابع فعال ساز در شبکه های عمیق چند لایه برای ایجاد غیر خطی بودن استفاده می شوند. این غیر خطی بودن به شبکه اجازه می دهد تا روابط پیچیده تر را یاد بگیرد و قادر به یادگیری روابط غیر خطی بین ورودی ها و خروجی ها باشد.

(ج) $\text{vanishing/exploding gradients}$: در میان تاج سگوبر بردن تغییر بزرگ یا کوچ به عنوان می تواند باعث لغزش یا دیگر در شبکه های عمیق می شود به عنوان راه حل می توانیم از توابع فعال ساز مانند ReLU استفاده کنیم.

(د) خروجی های غیر متوازن: تابع Sigmoid خروجی های در بازه $(0, 1)$ تولید می کند که باعث متغیر شدن نرخ خروجی ها در حدهای صفر و یک می شود. به عنوان راه حل می توانیم از توابع فعال ساز استفاده کنیم.

$$\begin{aligned}
 \text{forward} \rightarrow z &= w_{11} \cdot x_1 + w_{01} \cdot x_2 + b_1 \\
 a &= \text{ReLU}(z) = \max(0, z) \\
 \text{backward} \rightarrow \frac{\partial L}{\partial z} &= \begin{cases} 0 & \text{if } z \leq 0 \\ 1 & \text{if } z > 0 \end{cases} \\
 \frac{\partial L}{\partial w_{11}} &= \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w_{11}}
 \end{aligned}$$

در این مرحله وزن ها را بر روی وزن های می کنیم.

سوال ۳)

ج. بررسی راه حل (نیازی به پاسخ طولانی نیست، تنها رساندن مفهوم کافی ست)

۱. در مورد مفاهیم Grid Search و Random Search تحقیق کنید و هر کدام را مختصراً توضیح دهید.

۲. در مورد کرنل های مختلف مانند Linear, RBF, Polynomial تحقیق کنید و یک یا دو مورد از مهم ترین پارامترهای هر کدام را مختصر توضیح دهید.

۳. روش های one vs all و one vs rest را مختصراً توضیح داده و باهم مقایسه کنید. آیا در این مسئله نیازی به استفاده از آنها داریم؟

Grid Search

در این روش، فضای جستجو به یک شبکه منظم تقسیم می شود، و الگوریتم تمام نقاط شبکه را بررسی می کند.

ویژگی ها:

- ساختار منظم: فضای جستجو به طور سیستماتیک به نقاط مشخص تقسیم می شود.

- جامعیت: تمامی ترکیب های ممکن از پارامترها مورد بررسی قرار می گیرند.

- مزایا:

- ساده و قابل درک است.

- تضمین می کند که هیچ نقطه ای از فضای جستجو از قلم نمی افتد.

- معایب:

- بسیار زمان بر و محاسباتی سنگین است، به ویژه در مسائل با تعداد زیاد پارامترها (به دلیل افزایش نمایی نقاط شبکه).

کاربرد:

معمولاً زمانی استفاده می شود که فضای جستجو کوچک و تعداد پارامترها محدود باشد.

Random Search

در این روش، نقاطی به صورت تصادفی در فضای جستجو انتخاب و بررسی می شوند.

ویژگی‌ها:

- عدم ساختار منظم: نقاط به‌طور کاملاً تصادفی انتخاب می‌شوند.
- انعطاف‌پذیری: برخلاف Search Grid، نیازی به تقسیم‌بندی سیستماتیک فضای جستجو نیست.
- مزایا:

- در مسائل با فضای جستجوی بزرگ‌تر، کارآمدتر از Search Grid است.
- امکان کشف نقاط غیرمنتظره‌ای از فضای جستجو را دارد.

معایب:

- ممکن است برخی نقاط از فضای جستجو نادیده گرفته شوند.
- نتایج وابسته به شانس و تعداد نمونه‌های تصادفی است.

کاربرد:

در مسائل پیچیده با فضای جستجوی بزرگ، زمانی که منابع محاسباتی محدود باشند.

کرنل خطی (Linear Kernel)

این کرنل ساده‌ترین نوع است و فرض می‌کند داده‌ها در فضای اصلی به‌صورت خطی جداپذیر هستند.

فرمول:

$$y \cdot x = K(x, y)$$

پارامتر مهم:

• C (پارامتر منظم‌سازی)

- تنظیم تعادل بین حداکثرسازی حاشیه و کاهش خطا در داده‌های آموزشی.
- مقدار بزرگ‌تر: جریمه بیشتر برای خطاها (احتمال بیش‌برازش).
- مقدار کوچک‌تر: حاشیه بزرگ‌تر (احتمال کم‌برازش).
-

کاربرد:

- برای مسائل ساده و داده‌های خطی مناسب است.

کرنل چندجمله‌ای (Polynomial Kernel)

این کرنل برای مسائل پیچیده‌تر با روابط غیرخطی استفاده می‌شود.

فرمول:

$$d(c + y \cdot x) = K(x, y)$$

پارامترهای مهم:

۱. (Degree): درجه یا

- تعیین می‌کند تا چه حد ویژگی‌های غیرخطی در داده مدل‌سازی شوند.
- مقدار کوچک‌تر (d=2, d=2, d=2) یا (d=3, d=3, d=3): مناسب برای روابط ساده.
- مقدار بزرگ‌تر: ممکن است پیچیدگی زیاد و بیش‌برازش ایجاد کند.

۲. c (Bias):

- مقدار ثابت که تأثیر تعادل بین درجات پایین و بالا را تنظیم می‌کند.
- مقدار مناسب معمولاً با آزمایش مشخص می‌شود.

کاربرد:

- برای داده‌هایی با الگوهای پیچیده‌تر که دارای روابط چندجمله‌ای هستند.

کرنل شعاعی پایه RBF

یکی از پرکاربردترین کرنل‌ها، به‌ویژه برای مسائل غیرخطی. این کرنل داده‌ها را به فضای ویژگی‌های بسیار بزرگ نگاشت می‌کند.

فرمول:

$$K(x, y) = \exp(-\gamma \|x - y\|^2)$$

پارامترهای مهم:

۱. گاما

- تعیین کننده ی گستره تأثیر هر نمونه.
- مقدار کوچک تر: حاشیه وسیع تر (انعطاف کمتر).
- مقدار بزرگ تر: حاشیه باریک تر (ممکن است بیش برآزش شود).

۲. C

- مشابه کرنل خطی، کنترل تعادل بین دقت در آموزش و تعمیم پذیری مدل.

کاربرد:

- برای داده هایی که کاملاً غیرخطی هستند و مرز تصمیم گیری پیچیده ای نیاز دارند.

۱. روش One-vs-Rest

در این روش، برای هر کلاس، یک دسته بند دوتایی آموزش داده می شود.

• ایده:

- یک کلاس به عنوان مثبت در نظر گرفته می شود، و باقی کلاس ها به عنوان منفی.
- به ازای هر کلاس، یک مدل دوتایی ایجاد می شود.

• نحوه عملکرد:

- تعداد دسته بندها برابر با تعداد کلاس ها (k) است.
- هنگام پیش بینی، نمونه به کلاسی نسبت داده می شود که مدل مربوطه، بالاترین احتمال یا اطمینان را داشته باشد.

مزایا:

- ساده و قابل فهم.
- محاسبات کمتر نسبت به روش All-vs-All.

معایب:

- در مسائل نامتوازن ممکن است عملکرد خوبی نداشته باشد.
- نیازمند مقایسه نتایج برای تعیین کلاس است.

۲. روش One-vs-One

در این روش، برای هر جفت از کلاس‌ها، یک دسته‌بند دوتایی آموزشی داده می‌شود.

• ایده:

- هر بار فقط دو کلاس در نظر گرفته می‌شوند.
- تمام ترکیب‌های ممکن از کلاس‌ها بررسی می‌شود.

• نحوه عملکرد:

- تعداد دسته‌بندها برابر با $k(k-1)/2$ (تعداد جفت‌های کلاس‌ها) است.
- هنگام پیش‌بینی، هر مدل دوتایی یک رأی می‌دهد، و کلاس با بیشترین رأی انتخاب می‌شود.

مزایا:

- مناسب برای مسائل نامتوازن.
- می‌تواند مرزهای تصمیم‌گیری دقیق‌تری ارائه دهد.

معایب:

- تعداد مدل‌های بیشتری نیاز است، که محاسبات بیشتری می‌طلبد.
- پیچیدگی افزایش می‌یابد.

در مسئله‌ی مطرح شده چون دسته‌بندی فقط دو کلاسه است نیازی به استفاده از این روش‌ها نیست.

```
print(grid.best_score_)
print(grid.best_estimator_)

0.7854545454545454
SVC(C=0.1, degree=2, gamma=0.1, kernel='linear', max_iter=1000)

[ ] best_model = grid.best_estimator_
y_pred = best_model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy}")

Accuracy: 0.7681159420289855
```


سوال ۴)

۵. یکی از مراحل پیش‌پردازش معمولاً نرمال سازی یا استاندارد سازی است، ابتدا این دو روش را با هم مقایسه کنید سپس بدون استفاده از توابع آماده، آنها را روی داده‌ها اعمال کنید و نتایج را با قسمت های قبل مقایسه کنید.

نرمال سازی مقیاس داده‌ها را به یک بازه مشخص، معمولاً $[0,1]$ یا $[-1,1]$ تبدیل می‌کند.

• فرمول رایج (Min-Max Scaling):

$$\frac{\min X - X}{\min X_{max} - X} = 'X$$

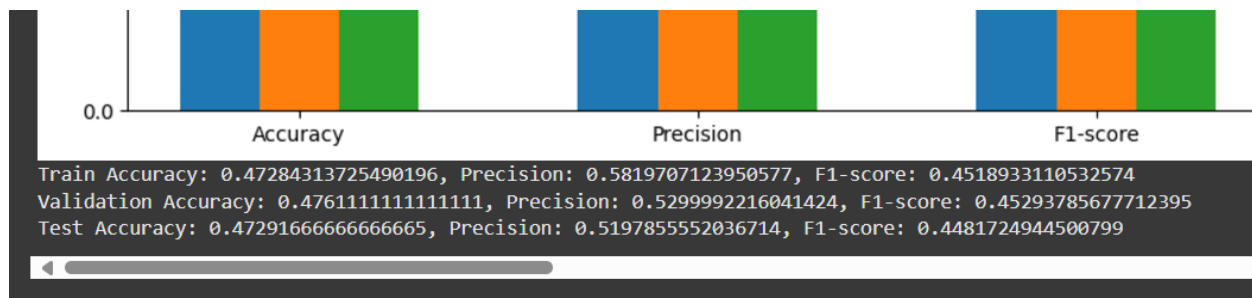
استاندارد سازی داده‌ها را به گونه‌ای تغییر می‌دهد که میانگین صفر و انحراف معیار واحد داشته باشند.

• فرمول:

$$\frac{X - \mu}{\sigma} = 'X$$

```
Epoch 1/10
797/797 ————— 2s 3ms/step - accuracy: 0.6358 - loss: 1.3827 - val_accuracy: 0.5
Epoch 2/10
797/797 ————— 1s 2ms/step - accuracy: 0.8782 - loss: 0.4761 - val_accuracy: 0.5
Epoch 3/10
797/797 ————— 1s 2ms/step - accuracy: 0.8960 - loss: 0.3850 - val_accuracy: 0.5
Epoch 4/10
797/797 ————— 1s 2ms/step - accuracy: 0.9025 - loss: 0.3513 - val_accuracy: 0.5
Epoch 5/10
797/797 ————— 3s 2ms/step - accuracy: 0.9118 - loss: 0.3162 - val_accuracy: 0.5
Epoch 6/10
797/797 ————— 3s 2ms/step - accuracy: 0.9143 - loss: 0.3023 - val_accuracy: 0.5
Epoch 7/10
797/797 ————— 2s 2ms/step - accuracy: 0.9193 - loss: 0.2868 - val_accuracy: 0.5
Epoch 8/10
797/797 ————— 1s 2ms/step - accuracy: 0.9206 - loss: 0.2752 - val_accuracy: 0.5
Epoch 9/10
797/797 ————— 1s 2ms/step - accuracy: 0.9247 - loss: 0.2663 - val_accuracy: 0.5
Epoch 10/10
797/797 ————— 2s 2ms/step - accuracy: 0.9286 - loss: 0.2549 - val_accuracy: 0.5
313/313 ————— 1s 3ms/step - accuracy: 0.9247 - loss: 0.2707
Test accuracy: 0.9343000054359436
```

نقته مدل با normalization



دقت مدل بدون normalization

و. نرخ یادگیری را به مقادیر $1e-5$ و $9e-1$ تغییر دهید و مشاهدات خود را تفسیر کنید.

متأسفانه من تغییری بین این مقادیر مشاهده نکردم، شاید اگر تعداد epoch ها بیشتر باشند میتوانیم تغییرات محسوسی داشته باشیم.

ز. در مورد بیش برآزش و کم برآزش در شبکه‌های عصبی توضیح دهید و برای پیشگیری و حل این مشکلات راهکارهایی پیشنهاد کنید (برای هر کدام دو مورد). نتایج قسمت‌های قبل را با توجه به این مفاهیم تحلیل کنید و بیان کنید در هر مرحله با کدام یک از این مشکلات مواجه بودید.

بیش‌برآزش زمانی رخ می‌دهد که مدل به‌قدری پیچیده باشد که الگوهای جزئی یا نویز موجود در داده‌های آموزشی را نیز یاد بگیرد. در نتیجه، مدل در داده‌های آموزشی عملکرد بسیار خوبی دارد اما در داده‌های جدید (آزمایشی) عملکرد ضعیفی نشان می‌دهد. کم‌برآزش زمانی رخ می‌دهد که مدل حتی در داده‌های آموزشی هم عملکرد خوبی ندارد و قادر به یادگیری الگوهای اساسی نیست.

از آنجایی که مدل بدون استفاده از normalization روی داده‌های آموزشی هم دقت خوبی نداشتند به نظر می‌آید دچار کم برآزش شده ایم ولی با استفاده از normalization این مشکل برطرف شده است.

سوال ۵

الف) $g_1 = g_2 \rightarrow x = \frac{1}{4}$

$g_2 = g_4 \rightarrow x = y = 1$

$g_1 = g_3 \rightarrow x = y = -2$

$g_2 = g_4 \rightarrow x = y = \frac{3}{4}$

$g_1 = g_4 \rightarrow y = \frac{1}{4}$

$g_2 = g_3 \rightarrow y = \frac{5}{4}$

شکل ۳-۴ ← ~~شکل ۸-۴~~ نمودار

ب) شکل ۱ و ۲ به کمره

وزن هر فرد ۱

Bias $-\frac{n}{4}$

ج) قسمت ۱ ←

$y = g\left(\sum_{i=1}^n x_i - \frac{n}{4}\right)$ تابع نوسازی

$g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$ ← تابع پله

<u>a</u>	<u>B</u>	<u>y</u>	<u>b</u>
$w_{1a} = 0$	$w_{1B} = 1$	$w_{1y} = 0$	$w_{ab} = 1$
$w_{2a} = 1$	$w_{2B} = 0$	$w_{2y} = 1$	$w_{Ba} = 1$
$w_a = 0$	$w_B = 1$	$w_y = -3$	$w_b = 4$

<u>A</u>	<u>B</u>	<u>C</u>
$w_{BA} = 1$	$w_{AB} = 1$	$w_{AC} = -1$
$w_{yA} = 1$	$w_{BB} = 1$	$w_{BC} = 1$
$w_A = -2$	$w_{yB} = -1$	$w_{Cs} = -1$
	$w_B = 1$	