



به نام خدا

دانشگاه تهران

دانشکده علوم و فناوری های میان رشته ای

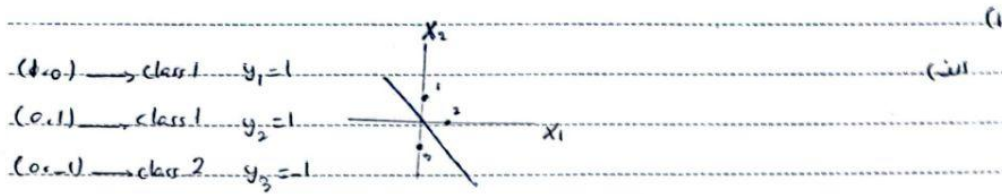
# Neural Networks

تمرین دوم

|             |                    |
|-------------|--------------------|
| OMID MORADI | نام و نام خانوادگی |
|             | شماره دانشجویی     |
|             | تاریخ ارسال گزارش  |

# سوال اول

(الف)



$$w_1 x_1 + w_2 x_2 + b = 0 \quad \text{بافتراض } w_1 = w_2 = 1 \quad \rightarrow \quad x_1 + x_2 = 0 \quad \rightarrow \quad \begin{cases} w_1 = 1 \\ w_2 = 1 \\ b = 0 \end{cases}$$

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

حساب:

$$D_1 = \frac{1}{\sqrt{2}}, \quad D_2 = \frac{1}{\sqrt{2}}, \quad D_3 = \frac{1}{\sqrt{2}}$$

$$\text{margin} = 2 \cdot \frac{1}{\sqrt{2}} = \sqrt{2}$$

(ب)

$$\begin{cases} \min_{w,b} \frac{1}{2} \|w\|^2 \\ y_i (w \cdot x_i + b) \geq 1, \forall i \end{cases}$$

dual form:  $\max_d \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) \quad \begin{cases} \alpha_i \geq 0 \\ \sum_{i=1}^n \alpha_i y_i = 0 \end{cases}$

$n=3, K = x_i \cdot x_j$

$$K = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & -1 \\ 0 & -1 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} x_1 x_1 & x_1 x_2 & x_1 x_3 \\ x_2 x_1 & x_2 x_2 & x_2 x_3 \\ x_3 x_1 & x_3 x_2 & x_3 x_3 \end{bmatrix}$$

ماتریس K:

$$\begin{array}{l|l} i=1, j=1 & \alpha_1 \alpha_1 y_1 y_1 (x_1 \cdot x_1) = \alpha_1^2 \\ i=1, j=2 & \alpha_1 \alpha_2 y_1 y_2 (x_1 \cdot x_2) = 0 \\ i=1, j=3 & \alpha_1 \alpha_3 y_1 y_3 (x_1 \cdot x_3) = 0 \end{array} \quad \begin{array}{l|l} i=2, j=2 & \alpha_2 \alpha_2 y_2 y_2 (x_2 \cdot x_2) = \alpha_2^2 \\ i=2, j=3 & \alpha_2 \alpha_3 y_2 y_3 (x_2 \cdot x_3) = -\alpha_2 \alpha_3 \\ i=3, j=3 & \alpha_3 \alpha_3 y_3 y_3 (x_3 \cdot x_3) = \alpha_3^2 \end{array}$$

$$\rightarrow \frac{1}{2} \sum_{i=1}^3 \sum_{j=1}^3 \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) = \frac{1}{2} (\alpha_1^2 + \alpha_2^2 + \alpha_3^2 + 2\alpha_2 \alpha_3)$$

$$\max_{\alpha_1, \alpha_2, \alpha_3} \alpha_1 + \alpha_2 + \alpha_3 - \frac{1}{2} (\alpha_1^2 + \alpha_2^2 + \alpha_3^2 + 2\alpha_2 \alpha_3)$$

$$\text{قید: } \alpha_1 + \alpha_2 + \alpha_3 = 1 \quad \text{و} \quad \alpha_1 y_1 + \alpha_2 y_2 + \alpha_3 y_3 = 0 \Rightarrow \boxed{\alpha_1 + \alpha_2 = \alpha_3}$$

KKT:

$$\frac{\partial L}{\partial \alpha_1} = 1 - \alpha_1 = 0 \Rightarrow \alpha_1 = 1$$

$$\frac{\partial L}{\partial \alpha_2} = 1 - \alpha_2 + \alpha_3 = 0 \Rightarrow \alpha_2 + \alpha_3 = 1 \quad \rightarrow \quad \begin{cases} \alpha_1 = 1 \\ \alpha_2 = 0 \\ \alpha_3 = 1 \end{cases}$$

$$\frac{\partial L}{\partial \alpha_3} = 1 - \alpha_3 + \alpha_2 = 0 \Rightarrow \alpha_2 - \alpha_3 = -1$$

$$\text{نتیجه: } \alpha_1 + \alpha_2 = \alpha_3 \Rightarrow 1 + 0 = 1 \quad \checkmark \checkmark$$

نتیجه:

$$w = \sum_{i=1}^n \alpha_i y_i x_i = (1)(1)(1, 0) + (0)(1)(0, 1) + (1)(-1)(0, -1) = (1, 1)$$

نتیجه: در هر دو خطی که می‌بینیم، است. آری

پس، اگر این نتایج را با هم مقایسه کنیم، می‌توانیم ببینیم که کدام نقطه از نقاط KKT، شرایط را برآورده می‌کند و در این صورت

$$1) \quad y_i (w x_i + b) = 1$$

با توجه به خطی که می‌بینیم، نقطه ۱ بر خط قرار دارد، پس

$$2) \quad \alpha_i = 0$$

$$x_{in} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$w_1 = \begin{bmatrix} 0 & 1 & 2 \\ 0 & -1 & 0 \end{bmatrix}$$

$$b = 1$$

$$z_1 = w_1^T x + b = \begin{bmatrix} 0 & 1 & 2 \\ 0 & -1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 1 = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

$$z'_1 = G(z_1) = \begin{bmatrix} \frac{1}{1+e^{-1}} \\ \frac{1}{1+e^{-2}} \\ \frac{1}{1+e^{-3}} \end{bmatrix} = \begin{bmatrix} 0.73 \\ 0.88 \\ 0.95 \end{bmatrix}$$

خروجی:

$$z_{out} = w_2^T z'_1 + b = \begin{bmatrix} 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0.73 \\ 0.88 \\ 0.95 \end{bmatrix} + 1 = 2.68$$

$$z'_{out} = G(z_{out}) = \frac{1}{1+e^{-2.68}} = 0.936$$

Error:

$$MSE = \frac{1}{2} (y_{true} - z'_{out})^2 = \frac{1}{2} (1 - 0.936)^2 = 0.002$$

backpropagation:

$$\delta_{out} = \frac{\partial E}{\partial z_{out}} = (a_{out} - y_{true}) \cdot G'(z_{out}) \quad \leftarrow \text{وزن ورودی}$$

$$G'(z_{out}) = a_{out} (1 - a_{out})$$

$$\delta_{out} = (a_{out} - y_{true}) \cdot a_{out} (1 - a_{out})$$

$$= (0.936 - 1) \cdot 0.936 (1 - 0.936) = -0.0038$$

بهتر ساز وزن در لایه مخفی:

$$w_2 = w_2 - \eta \cdot \delta_{out} \cdot a_{hidden} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} - 0.5 (-0.0038) \begin{bmatrix} 0.73 \\ 0.88 \\ 0.95 \end{bmatrix}$$

$$= \begin{bmatrix} 1.0014 \\ 0.0017 \\ 1.0018 \end{bmatrix}$$

بهتر ساز وزن در لایه خروجی:

$$b_{out} = b_{out} - \eta \cdot \delta_{out} = 1 - 0.5 (-0.0038) = 1.0019$$

در لایه مخفی:

(نصب کردن)

$$\delta_{hidden_1} = \delta_{out} \cdot w_2 \cdot a_{hidden} (1 - a_{hidden})$$

$$= -0.0038 \times 1 \times 0.73 (1 - 0.73) = -0.00075$$

(نودون):

$$\delta_{hidden_2} = -0.0038 \times 0.88 (1 - 0.88) = 0$$

(نودون):

$$\delta_{hidden_1} = -0.0038 \times 1 \times 0.99 (1 - 0.99) = -0.00017$$

$$\delta_{hidden} = \begin{bmatrix} -0.00075 \\ 0 \\ -0.00017 \end{bmatrix}$$

آپرستون:

$$w_1 = w_1 - \eta \delta_{hidden} \cdot \alpha^T = \begin{bmatrix} 0 & 1 & 2 \\ 0 & -1 & 0 \end{bmatrix} - 0.5 \begin{bmatrix} -0.00075 \\ 0 \\ -0.00017 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0.00037 & 1 & 2 \\ 0 & -1 & 0 \end{bmatrix}$$

بایس:

$$b_{hidden} = b_{hidden} - \eta \delta_{hidden} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} - 0.5 \begin{bmatrix} -0.00075 \\ 0 \\ -0.00017 \end{bmatrix} = \begin{bmatrix} 1.00037 \\ 1 \\ 0.000085 \end{bmatrix}$$

Scanned with CamScanner

(ب)

1. غیرخطی کردن خروجی‌ها
2. ایجاد مرزهای تصمیم‌گیری پیچیده
3. افزایش توانایی شبکه در یادگیری الگوهای پیچیده
4. کنترل و محدود کردن خروجی‌ها
5. مدل‌سازی مفاهیم غیرخطی در داده‌ها

(پ)

### 1. محو شدن گرادیان

یکی از مشکلات سیگموئید، پدیده محو شدن گرادیان است. در این حالت، اگر ورودی بسیار بزرگ یا بسیار کوچک باشد، مشتق سیگموئید تقریباً به صفر نزدیک می‌شود. این موضوع در شبکه‌های عمیق باعث می‌شود گرادیان در طی پس‌انتشار کاهش یابد و لایه‌های ابتدایی نتوانند به خوبی به‌روزرسانی شوند. برای



رفع این مشکل، استفاده از ReLU پیشنهاد می‌شود، زیرا در ناحیه مثبت، مشتق آن برابر با 1 است و از محو شدن گرادیان جلوگیری می‌کند. همچنین نرمال‌سازی ورودی‌ها (مثلاً استاندارد کردن داده‌ها) می‌تواند از ورود مقادیر بیش از حد بزرگ یا کوچک به تابع سیگموئید جلوگیری کند.

## 2. خروجی‌های غیر صفر مرکزی (Non-Zero Centered Outputs)

مشکل دیگر سیگموئید این است که خروجی‌های آن **غیر صفر مرکزی** هستند، یعنی مقادیر خروجی همیشه در بازه  $[0,1]$  قرار دارند. این موضوع باعث می‌شود گرادیان‌ها همیشه یا مثبت یا صفر باشند، که الگوریتم گرادیان نزولی را کند می‌کند. برای حل این مشکل، می‌توان از تابع تانژانت هایپربولیک (Tanh) استفاده کرد که خروجی آن در بازه  $[-1,1]$  قرار دارد و گرادیان‌ها را متقارن‌تر می‌کند. استفاده از روش‌های بهینه‌سازی پیشرفته مانند Adam نیز می‌تواند تأثیر این مشکل را کاهش دهد.

### برتری RELU نسبت به Sigmoid:

ReLU نسبت به سیگموئید مزایای زیادی دارد. این تابع محاسبات بسیار ساده‌تری دارد، زیرا فقط مقدار حداکثر  $\max(0, x)$  را محاسبه می‌کند، در حالی که سیگموئید به محاسبات نمایی نیاز دارد. علاوه بر این، ReLU به دلیل صفر کردن خروجی برای مقادیر منفی، یک خروجی پراکنده تولید می‌کند که باعث ساده‌تر شدن ساختار شبکه می‌شود. همچنین، در شبکه‌های عمیق، ReLU عملکرد بهتری دارد زیرا از محو شدن گرادیان جلوگیری کرده و سرعت یادگیری را افزایش می‌دهد.

بنابراین، ReLU به دلیل سادگی، کارایی بالا و توانایی حل مشکلاتی که سیگموئید ایجاد می‌کند، انتخاب بهتری برای شبکه‌های عمیق است.

(ت)

در مسیر پیشرو (Forward Pass)، تابع فعال‌ساز ReLU به صورت  $f(x) = \max(0, x)$  عمل می‌کند. اگر ورودی  $x$  مثبت باشد، خروجی همان مقدار ورودی است، و اگر ورودی منفی یا صفر باشد، خروجی صفر می‌شود. این رفتار باعث می‌شود که مقادیر مثبت بدون تغییر به لایه بعدی منتقل شوند و مقادیر منفی به صفر تبدیل شوند. بنابراین، برخی از نورون‌ها خاموش شده و در محاسبات بعدی تأثیر ندارند.

مثال عددی:

ورودی 2:

$$f(x) = \max(0, x) = (0, 2) = 2$$

ورودی 3:-

$$f(x) = \max(0, x) = (0, -3) = 0$$

در مسیر پسرو (Backward Pass)، مشتق تابع ReLU نسبت به ورودی محاسبه می‌شود.

برای ورودی‌های مثبت ( $x > 0$ )، مشتق برابر با 1 است، که به این معنی است که گرادیان به طور کامل به لایه‌های قبلی منتقل می‌شود و یادگیری ادامه پیدا می‌کند. اما برای ورودی‌های غیرمثبت ( $x \leq 0$ )، مشتق برابر صفر است، و گرادیانی به وزن‌های مرتبط با این نورون‌ها منتقل نمی‌شود. این باعث می‌شود که این نورون‌ها به اصطلاح "خاموش" شوند و دیگر یاد نگیرند.

این رفتار مزایایی دارد، مانند جلوگیری از محو شدن گرادیان و تسریع یادگیری، زیرا مشتق ReLU برای مقادیر مثبت ثابت است. اما مشکل اصلی این است که نورون‌هایی که ورودی غیرمثبت دریافت می‌کنند ممکن است برای همیشه خاموش باقی بمانند (مشکل نورون‌های مرده). برای رفع این مشکل می‌توان از نسخه‌های اصلاح‌شده ReLU مانند Leaky ReLU استفاده کرد، که برای مقادیر منفی یک گرادیان کوچک (مثلاً  $0.01x$ ) در نظر می‌گیرد. این کار باعث می‌شود نورون‌های خاموش همچنان امکان یادگیری داشته باشند. گزینه دیگر ReLU6 است که خروجی ReLU را در بازه  $[0, 6]$  محدود می‌کند و کنترل بیشتری روی گرادیان فراهم می‌کند.

به طور کلی، ReLU در مسیر پیشرو و پسرو برای مقادیر مثبت بسیار مؤثر عمل می‌کند و در شبکه‌های عمیق به دلیل جلوگیری از محو شدن گرادیان محبوب است.

برای ورودی 2 گرادیان:

$$f'(x) = 1$$

برای ورودی 3- گرادیان:

$$f'(x) = 0$$

## سوال سوم

(ب)

بررسی مقادیر گمشده: (Missing Values)

- **دلیل:** مقادیر گم‌شده می‌توانند باعث خطا در مدل‌سازی شوند. روش‌های مختلفی مانند حذف ردیف‌ها، پر کردن مقادیر با میانگین، میانه، یا مدل‌سازی آماری وجود دارند.

#### **حذف یا مدیریت مقادیر پرت: (Outliers)**

- **دلیل:** مقادیر پرت می‌توانند تأثیر زیادی بر الگوریتم‌های حساس به مقیاس داشته باشند. حذف یا تعدیل این مقادیر باعث می‌شود مدل پایدارتر شود.

#### **مقیاس‌بندی داده‌ها: (Scaling)**

- **دلیل:** در الگوریتم‌هایی که از فاصله برای محاسبات استفاده می‌کنند (مانند رگرسیون یا الگوریتم‌های مبتنی بر درخت)، مقیاس متغیرها اهمیت زیادی دارد.

#### **رمزگذاری متغیرهای دسته‌ای: (Encoding)**

- **دلیل:** داده‌های دسته‌ای باید به فرم عددی تبدیل شوند تا الگوریتم‌های یادگیری ماشین بتوانند با آن‌ها کار کنند.

#### **حذف یا ترکیب ویژگی‌های غیرمفید: (Feature Selection/Engineering)**

- **دلیل:** برخی ویژگی‌ها ممکن است همبستگی بالایی با یکدیگر داشته باشند یا اطلاعات کمی برای مدل‌سازی ارائه دهند.

#### **تبدیل داده‌های نامتقارن یا غیرخطی:**

- **دلیل:** داده‌های نامتقارن می‌توانند نتایج مدل را تحت تأثیر قرار دهند. استفاده از تبدیلات لگاریتمی یا جذر می‌تواند توزیع داده‌ها را متعادل کند.

(پ)

(1)

- **Grid Search:** جستجوی سازمان‌یافته روی تمام ترکیب‌های ممکن از مقادیر پارامترهای مشخص شده. دقیق است اما زمان‌بر.
- **Random Search:** نمونه‌برداری تصادفی از ترکیبات پارامترها در فضای جستجو. سریع‌تر است و برای فضاهای بزرگ‌تر مؤثرتر.



(2)

- **Linear Kernel**: برای داده‌های خطی؛ پارامتر مهم: نیاز به پارامتر اضافی خاصی ندارد.

- **RBF Kernel**: مناسب برای داده‌های غیرخطی؛ پارامترها:

- **Gamma**: تأثیرگذاری نقاط دور و نزدیک.

- **C**: جریمه خطای دسته‌بندی.

- **Polynomial Kernel**: برای مسائل با روابط چندجمله‌ای؛ پارامترها:

- **degree**: درجه چندجمله‌ای.

- **C**: میزان تحمل خطا.

(3)

- **One-vs-Rest (OvR)**: هر کلاس را در مقابل سایر کلاس‌ها مقایسه می‌کند. مزیت: ساده‌تر و

سریع‌تر برای تعداد کلاس زیاد.

- **One-vs-One (OvO)**: تمام جفت کلاس‌ها را مقایسه می‌کند. مزیت: معمولاً برای داده‌های

کوچک‌تر بهتر عمل می‌کند.

**در این مسئله:** اگر دسته‌بندی چندکلاسی باشد، بسته به تعداد کلاس‌ها یکی از این روش‌ها لازم

است. OvR برای تعداد کلاس‌های زیاد مناسب‌تر است. در نتیجه برای این مثال OvR مناسب نیست.

(ث)

مدل SVM با استفاده از روش GridSearch برای بهینه‌سازی پارامترها به نتایج قابل توجهی دست یافت.

بهترین پارامترهای به‌دست‌آمده شامل C برابر با 1، degree برابر با 2، gamma با مقدار "scale" و کرنل

"rbf" هستند. امتیاز اعتبارسنجی متقاطع (Cross-Validation) مدل 77.89 درصد گزارش شده است.

دقت مدل روی داده‌های تست برابر با 82.23 درصد است که عملکرد مطلوب آن را نشان می‌دهد.

ماتریس درهم‌ریختگی نشان می‌دهد که مدل توانسته است 56.58 درصد از نمونه‌های منفی واقعی را

به‌درستی منفی پیش‌بینی کند و 25.66 درصد از نمونه‌های مثبت واقعی را نیز به‌درستی به‌عنوان مثبت

شناسایی کند. خطاهای False Positive و False Negative به ترتیب برابر با 8.55 درصد و 9.21

درصد هستند. این مقادیر نشان می‌دهند که مدل در ایجاد تعادل بین پیش‌بینی‌های مثبت و منفی عملکرد مناسبی داشته است.

به‌طور کلی، مدل عملکرد خوبی روی داده‌ها ارائه داده است و برای بهبود بیشتر می‌توان روش‌های پیش‌پردازش داده را تقویت کرد یا از داده‌های آموزشی بیشتری استفاده کرد.

## سوال چهارم

(ب)

جلوگیری از بیش‌برازش:

مدل باید روی داده‌های دیده‌نشده آزموده شود تا بتوان گفت عملکرد آن به داده‌های خاص وابسته نیست.

تنظیم هایپرپارامترها:

استفاده از مجموعه اعتبارسنجی کمک می‌کند بهترین تنظیمات را برای مدل انتخاب کنیم.

ارزیابی بی‌طرفانه:

مجموعه آزمون به‌عنوان داده‌های "کاملاً نادیده" استفاده می‌شود تا ارزیابی منصفانه‌ای از مدل ارائه دهد.

(ت)

**روند بهبود:** افزایش تعداد نوروها و اضافه کردن لایه‌های بیشتر باعث بهبود دقت مدل شده است.

**مشکل احتمالی:** دقت کلی هنوز پایین است. این موضوع ممکن است ناشی از:

1. نیاز به تنظیم بهتر هایپرپارامترها (مانند نرخ یادگیری یا تعداد دوره‌ها).

2. حجم کم داده‌های آموزشی یا کیفیت پایین داده‌ها.

3. مدل ساده‌ای که هنوز برای این داده‌ها کافی نیست

|   | Model                             | Test Accuracy |
|---|-----------------------------------|---------------|
| 0 | 1 Hidden Layer (16 Neurons)       | 0.16          |
| 1 | 1 Hidden Layer (128 Neurons)      | 0.12          |
| 2 | 2 Hidden Layers (64, 128 Neurons) | 0.08          |

(ث)

نرمال سازی داده ها را به یک بازه ثابت، معمولاً  $[0, 1]$  یا  $[-1, 1]$ ، مقیاس بندی می کند. این روش در مواقعی که مقیاس داده ها اهمیت دارد یا الگوریتم از فاصله محور استفاده می کند (مانند KNN و SVM) بسیار مفید است. با این حال، نرمال سازی به مقادیر پرت حساس است، زیرا این مقادیر می توانند به شدت بر محدوده داده ها تاثیر بگذارند.

استاندارد سازی، برعکس، داده ها را به گونه ای تغییر می دهد که میانگین آن ها صفر و واریانس یک شود. این روش برای داده هایی با توزیع گاوسی مناسب است و تاثیر مقادیر پرت را کاهش می دهد. استاندارد سازی در الگوریتم هایی که به توزیع داده ها حساس هستند (مانند رگرسیون خطی و شبکه های عصبی) اغلب بهتر عمل می کند.

**نتایج کد:**

Test Accuracy with Normalization: 0.20499999821186066

Test Accuracy with Standardization: 0.41499999165534973

استاندارد سازی معمولاً برای داده های تصویری مانند miniMNIST که شامل مقادیر عددی گسترده ای است، عملکرد بهتری ارائه می دهد. دلیل این موضوع این است که استاندارد سازی ویژگی های هر نمونه را حول میانگین صفر و انحراف معیار یک متمرکز می کند، که باعث می شود شبکه عصبی بهتر بتواند اطلاعات مهم را استخراج کند و از اثرات مقادیر پرت جلوگیری کند.

در مقابل، نرمال سازی ممکن است برای داده هایی که به طور طبیعی در یک بازه محدود نیستند یا دارای مقادیر پرت قابل توجهی هستند، کارایی کمتری داشته باشد، زیرا بازه ثابت  $[0, 1]$  نمی تواند تفاوت های پیچیده بین مقادیر داده را به خوبی مدیریت کند.

(ج)

با نرخ یادگیری  $1e-9$ :

|   | Model                             | Test Accuracy |
|---|-----------------------------------|---------------|
| 0 | 1 Hidden Layer (16 Neurons)       | 0.155         |
| 1 | 1 Hidden Layer (128 Neurons)      | 0.090         |
| 2 | 2 Hidden Layers (64, 128 Neurons) | 0.135         |

با نرخ یادگیری  $1e-5$ :

|   | Model                             | Test Accuracy |
|---|-----------------------------------|---------------|
| 0 | 1 Hidden Layer (16 Neurons)       | 0.045         |
| 1 | 1 Hidden Layer (128 Neurons)      | 0.135         |
| 2 | 2 Hidden Layers (64, 128 Neurons) | 0.045         |

نتایج با نرخ یادگیری  $1e-5$ :

1. مدل با 16 نورون در یک لایه مخفی: دقت بسیار پایین (4.5٪). این نشان می‌دهد که با نرخ یادگیری  $1e-5$ ، مدل نتوانسته به خوبی به حداقل خطا نزدیک شود و عملاً یادگیری موثری رخ نداده است.

2. مدل با 128 نورون در یک لایه مخفی: دقت کمی بهتر (13.5٪)، که نشان می‌دهد تعداد بیشتر نورون‌ها به مدل کمک کرده است ویژگی‌های بیشتری از داده را بیاموزد، اما نرخ یادگیری هنوز نتوانسته مدل را بهینه کند.

3. مدل با دو لایه مخفی (64 و 128 نورون): دقت بسیار پایین (4.5٪). این نتیجه نشان می‌دهد که با وجود معماری پیچیده‌تر، نرخ یادگیری ناکافی بوده و مدل عملاً به درستی آموزش ندیده است.

## نتایج با نرخ یادگیری $1e-9$ :

1. مدل با 16 نورون در یک لایه مخفی: دقت 14.5٪. این مدل با نرخ یادگیری بسیار پایین عملکرد بهتری نسبت به  $1e-5$  داشته است، زیرا احتمالاً توانسته تغییرات بسیار کوچک در وزن‌ها را به دقت مدیریت کند.
2. مدل با 128 نورون در یک لایه مخفی: دقت 10٪. در اینجا عملکرد ضعیف‌تر نسبت به مدل اول نشان می‌دهد که نرخ یادگیری بسیار پایین ممکن است از رسیدن به یک حداقل مطلوب جلوگیری کرده باشد.
3. مدل با دو لایه مخفی (64 و 128 نورون): دقت 8.5٪. معماری پیچیده‌تر در اینجا حتی با نرخ یادگیری بسیار پایین عملکرد ضعیف‌تری داشته است، احتمالاً به دلیل نیاز به نرخ یادگیری مناسب‌تر برای یادگیری موثر در لایه‌های عمیق‌تر.

(چ)

## بیش‌برازش: (Overfitting)

زمانی رخ می‌دهد که مدل به‌طور دقیق داده‌های آموزشی را یاد می‌گیرد اما در تعمیم‌دهی به داده‌های جدید و نادیده (مانند داده‌های آزمون) عملکرد ضعیفی دارد. این معمولاً زمانی رخ می‌دهد که مدل بیش از حد پیچیده باشد یا داده‌های آموزشی ناکافی باشند.

نشانه‌ها:

- دقت زیاد روی مجموعه آموزش، اما دقت پایین روی مجموعه اعتبارسنجی یا آزمون.
- افزایش مقدار Validation Loss با کاهش مقدار Training Loss.

## • راهکارها برای پیشگیری و حل:

### 1. رگولاراسیون: (Regularization)

- اضافه کردن جریمه‌هایی به تابع هزینه، مانند  $L1$  یا  $L2$ .

### 2. Dropout:

- به صورت تصادفی غیرفعال کردن درصدی از نورون‌ها در طول آموزش برای کاهش پیچیدگی مدل.

## کم‌برازش: (Underfitting)

زمانی رخ می‌دهد که مدل حتی روی داده‌های آموزشی نیز عملکرد ضعیفی دارد و نتوانسته الگوهای موجود در داده‌ها را به درستی یاد بگیرد. این معمولاً زمانی رخ می‌دهد که مدل بیش از حد ساده باشد یا زمان کافی برای آموزش نداشته باشد.

**نشانه‌ها:**

- دقت پایین روی مجموعه آموزش.
- عدم کاهش مقدار Training Loss حتی پس از تعداد زیادی اپوک.

• راهکارها برای پیشگیری و حل:

1. افزایش پیچیدگی مدل:

- اضافه کردن تعداد نوروں‌ها یا لایه‌ها.

2. کاهش نرخ یادگیری: (Learning Rate)

- کاهش نرخ یادگیری برای بهبود جزئیات یادگیری مدل.

تحلیل نتایج مراحل قبل با این مفاهیم:

بیش برآزش برای دو حالت با 128 نوروں و دولایه

کم برآزش برای حالت با 16 نوروں

**سوال پنجم**

**(الف)**

برای یافتن مرز تصمیم گیری باید شکل که توانیم با هم اشتراک دهیم را پیدا کنیم :

class 1 = class 2 :

$$x - y + 1 = -x - y + 2 \Rightarrow x = \frac{1}{2}$$

class 1 = class 3 :

$$x - y + 1 = -3 + y - x \Rightarrow x - y = -2$$

class 1 = class 4 :  $x - y + 1 = x + y \Rightarrow y = \frac{1}{2}$

class 2 = class 3 :

$$-x - y + 2 = -x + y - 3 \Rightarrow y = \frac{5}{2}$$

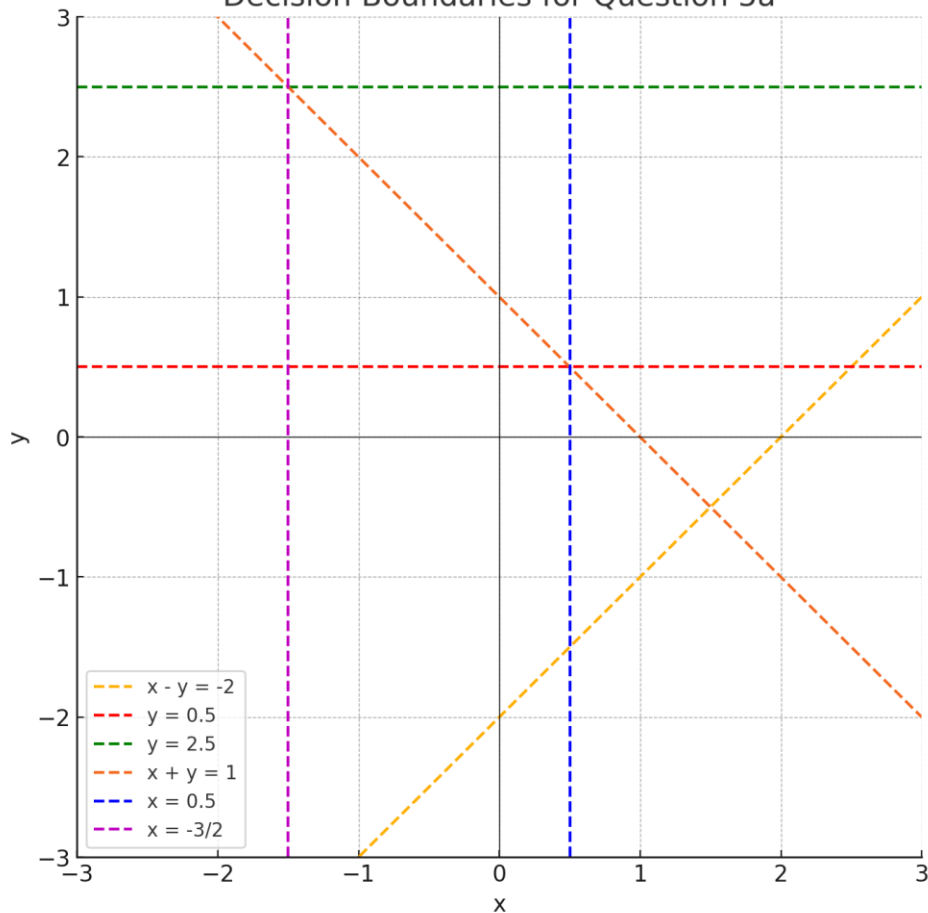
class 2 = class 4 :

$$-x - y + 2 = x + y \Rightarrow x + y = 1$$

class 3 = class 4 :

$$-x + y + 3 = x + y \Rightarrow x = \frac{3}{2}$$

Decision Boundaries for Question 5a





(ب)

- اگر داده‌ها به صورت خطی جداسدنی باشند، یک پرسپترون با دو گره ورودی کفایت می‌کند.
- اگر داده‌ها به صورت غیرخطی جداسدنی باشند، نیاز به افزودن ویژگی‌هایی است تا داده‌ها در فضای ویژگی جدید جداسدنی شوند.

(پ)

(1)

برای حل این بخش، باید شبکه‌ای طراحی شود که ورودی‌های  $X_1, X_2, \dots, X_n$  را دریافت کرده و تابع اکثریت را محاسبه کند. تابع اکثریت به این صورت است که اگر تعداد ورودی‌هایی که مقدار 1 دارند از نصف تعداد کل ورودی‌ها ( $n/2$ ) بیشتر باشد، خروجی برابر با 1 است. در غیر این صورت، خروجی برابر با 0 خواهد بود.

این شبکه یک لایه ورودی دارد که شامل  $n$  نورون است، و وظیفه آن دریافت ورودی‌های  $X_i$  است. وزن هر ورودی برابر با 1 در نظر گرفته می‌شود تا مقدار هر ورودی مستقیماً به جمع نهایی اضافه شود. برای تصمیم‌گیری، در لایه خروجی یک بایاس ( $b$ ) با مقدار  $-n/2$  تعریف می‌شود. خروجی شبکه از معادله

$$f\left(\sum_{i=1}^n w_i X_i\right) + b$$

به دست می‌آید، که  $f$  یک تابع فعال (مانند سیگموئید یا Heaviside) است.

(2)

تمامی مفادیر با توجه به فرمول  $\sum w_i x_i + b$  محاسبه شده.

گره  $\alpha$

| مقدار | وزن          |
|-------|--------------|
| 0     | $w_1 \alpha$ |
| 1     | $w_2 \alpha$ |
| 0     | $w \alpha$   |

گره  $\beta$

| مقدار | وزن |
|-------|-----|
|-------|-----|

|              |    |
|--------------|----|
| $w_{1\beta}$ | 1  |
| $w_{2\beta}$ | 0  |
| $w_\beta$    | -4 |

گره  $\gamma$

|               |       |
|---------------|-------|
| وزن           | مقدار |
| $w_{1\gamma}$ | 0     |
| $w_{2\gamma}$ | 1     |
| $w_\gamma$    | -4    |

گره  $\delta$

|                    |       |
|--------------------|-------|
| وزن                | مقدار |
| $w_{\alpha\delta}$ | 1     |
| $w_{\beta\delta}$  | 1/2   |
| $w_\delta$         | -1/2  |

گره A

|                |       |
|----------------|-------|
| وزن            | مقدار |
| $w_{\delta A}$ | 1     |
| $w_{\gamma A}$ | 1     |
| $w_A$          | -1    |

گره B

|                |       |
|----------------|-------|
| وزن            | مقدار |
| $w_{\alpha B}$ | 1     |
| $w_{\beta B}$  | 1     |
| $w_{\gamma B}$ | -1    |
| $w_B$          | -1.5  |

گره C

|                |       |
|----------------|-------|
| وزن            | مقدار |
| $w_{\alpha C}$ | -1    |
| $w_{\beta C}$  | 1     |
| $w_C$          | -1    |

