



## یادگیری ماشین

پاییز ۱۴۰۳

استاد: علی شریفی زارچی

مسئول تمرین: علی شاه حیدر

مهلت ارسال نهایی: ۹ آذر

## فصل سوم

تمرین سوم

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روزهای مشخص شده است.
- در طول ترم، برای هر تمرین می‌توانید تا ۵ روز تأخیر مجاز داشته باشید و در مجموع حداکثر ۱۵ روز تأخیر مجاز خواهید داشت. توجه داشته باشید که تأخیر در تمرین‌های عملی و تئوری به صورت جداگانه محاسبه می‌شود و مجموع تأخیر هر دو نباید بیشتر از ۱۵ روز شود. پس از اتمام زمان مجاز، دو روز اضافی برای آپلود غیرمجاز در نظر گرفته شده است که در این بازه به ازای هر ساعت تأخیر، ۲ درصد از نمره تمرین کسر خواهد شد.
- اگر بخش عملی یا تئوری تمرین را قبل از مهلت ارسال امتیازی آپلود کنید، ۲۰ درصد نمره اضافی به آن بخش تعلق خواهد گرفت و پس از آن، ویدئویی تحت عنوان راهنمایی برای حل تمرین منتشر خواهد شد.
- حتماً تمرین‌ها را بر اساس موارد ذکر شده در صورت سوالات حل کنید. در صورت وجود هرگونه ابهام، آن را در صفحه تمرین در سایت کوئرا مطرح کنید و به پاسخ‌هایی که از سوی دستیار آموزشی مربوطه ارائه می‌شود، توجه کنید.
- 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.
- فایل پاسخ‌های سوالات نظری را در قالب یک فایل pdf به فرمت  $HW3\_T\_ [STD\_ID].pdf$  آماده کنید و برای سوالات عملی، هریک را در یک فایل zip جداگانه قرار دهید و فایل zip اول را به فرمت  $HW3\_P1\_ [STD\_ID].zip$  و فایل zip دوم را به فرمت  $HW3\_P2\_ [STD\_ID].zip$  و فایل سوم را نیز به فرمت  $HW3\_P3\_ [STD\_ID].zip$  نامگذاری کرده و هرکدام را به صورت جداگانه آپلود کنید.
- گردآورندگان تمرین: کیارش جولایی، علیرضا صابونچی، امیرعلی لقمانی، آسمانه نافع، دنیا نوابی، فاطمه شیری

## سوالات نظری (۱۰۰ نمره)

حل.  
سوال ۱  
(الف)

پاسخ برابر ۱۴ است. برای ورودی‌ها ۴ حالت  $\{(0,0), (0,1), (1,0), (1,1)\}$  داریم. برای هر یک از این ۴ حالت، یک perceptron داریم که فقط یکی از این ۴ تا را با برجسب ۱ خروجی می‌دهد و بقیه را برجسب ۰ می‌دهد. مشابهاً هم برای هر کدام فقط یک perceptron هست که آن یکی را برجسب ۰ می‌دهد و بقیه را برجسب ۱ می‌دهد. لذا تا اینجا ۸ تا داریم. سپس یک perceptron وجود دارد که ورودی دوم را نادیده می‌گیرد و دو نقطه‌ای را که ورودی اولشان برابر ۱ است برجسب ۱ می‌دهد و دو نقطه دیگر را برجسب ۰ می‌زند. همچنین یک perceptron وجود دارد که برعکسش را انجام می‌دهد. این هم دو تای دیگر اضافه می‌کند. در نهایت یک جفت کاملاً عین هم از perceptron ها وجود دارد که ورودی اول را نادیده می‌گیرد و مانند حالت قبل عمل می‌کند، فقط این بار از ورودی دوم استفاده می‌کند. همچنین یک حالت بدیهی وجود دارد که همه داده‌ها برجسب ۱ می‌گیرند یا همه برجسب ۰ می‌گیرند که نماینده مرزهای خطی هستند که همه نقاط روی یک سمت آن قرار دارند.

(ب)

غلط است. با اینکه تابع  $\tanh$  برعکس sigmoid میانگین صفر است، اما باز هم مشکل ناپدید شدن گرادیان را دارد.

(ج)

صحیح است. می‌دانیم که یک رابطه منطقی را می‌توان به صورت جمع  $p_i$  هایی نوشت که هر یک به صورت ضرب  $a_j$  ها و نه آن‌ها نوشته می‌شوند. نورون‌های ورودی را  $a_j$  ها تعریف می‌کنیم که در صورتی که خودشان به کار رفته باشند، با وزن یک و در صورتی که not نشان به کار رفته باشد با وزن ۱- مشخص می‌کنیم و اگر آن ورودی به کار نرفته باشد، وزن آن را صفر قرار می‌دهیم. بایاس هر جمله هم به اندازه یکی کمتر از تعداد نورون‌ها است که در ساخت جمله ضربی کاربرد دارند البته به صورت منفی. در لایه اول ضرایب یعنی  $p_i$  ها را می‌سازیم. در لایه دوم هم از هر  $p_i$  با وزن یک به نورون خروجی متصل می‌کنیم و بایاس نمی‌گذاریم تا در صورتی که حداقل یکی از  $p_i$  ها یک شد، مقدار خروجی مثبت شود.

(د)

غلط است. افزایش عمق و عرض ظرفیت مدل را افزایش می‌دهد، اما بدون تنظیم منظم‌سازی مناسب یا داده‌های آموزشی کافی، ممکن است منجر به بیش‌برازش یا چالش‌های بهینه‌سازی شود.

(ه)

غلط است. اگرچه یک شبکه سطحی و پهن می‌تواند طبق قضیه تقریب عمومی توابع را تقریب بزند، ممکن است نیاز به تعداد بسیار بیشتری نورون نسبت به یک شبکه عمیق برای نمایش کارآمد توابع پیچیده داشته باشد.

(و)

مشکل ناپدید شدن گرادینان زمانی رخ می‌دهد که گرادینان‌ها در لایه‌های ابتدایی شبکه‌های عمیق بسیار کوچک می‌شوند و یادگیری شبکه متوقف می‌شود. توابع فعال‌سازی مانند ReLU با صفر کردن مقادیر منفی ورودی، این مشکل را کاهش می‌دهند و از کاهش بیش از حد گرادینان جلوگیری می‌کنند.

(ز)

غلط است. در حالی که Stochastic Gradient Descent (SGD) در مقایسه با Mini-Batch Gradient Descent به روزرسانی‌ها را بیشتر (پس از هر نمونه) انجام می‌دهد، این همگرایی سریع‌تر را تضمین نمی‌کند. نوین بودن به روزرسانی‌های فردی در SGD می‌تواند منجر به نوسانات بیشتر شود که ممکن است همگرایی را کاهش دهد یا باعث شود مدل یک راه‌حل بهینه را از دست بدهد. Mini-batch Gradient Descent، با اندازه دسته‌ای مناسب، با کاهش نوین تعادل ایجاد می‌کند و در عین حال امکان به روزرسانی‌های مکرر را فراهم می‌کند، که اغلب منجر به همگرایی پایدارتر و سریع‌تر در مقایسه با SGD خالص می‌شود. نرخ همگرایی همچنین به عواملی مانند نرخ یادگیری، اندازه دسته‌ای و مسئله بهینه‌سازی خاص بستگی دارد.

(ی)

معادلات لازم به صورت زیر هستند:

$$g_i = \nabla_{x_i} y(x_{i-1}) = 1/2 x_{i-1}^3 - 0.3 x_{i-1}^2 - 4 x_{i-1} - 0.8 \quad (1)$$

$$m_i = \mu m_{i-1} + g_i \quad (2)$$

$$x_i = x_{i-1} - \gamma m_i \quad (3)$$

مقادیر  $m_i$  از ممان و اطلاعات تکرارهای قبلی استفاده می‌کنند. برای تکرار اول،  $m_0 = 0$ . بنابراین، تکرار اول همان جواب مسئله بدون ممان را خواهد داشت.

$$\begin{aligned} g_1 &= 1/2 \times (-2/8000)^3 - 0.3 \times (-2/8000)^2 - 4 \times (-2/8000) - 0.8 = -18/2944, \\ m_1 &= 0.7 \times m_0 + g_1 = 0.7 \times 0 + (-18/2944) = -18/2944, \\ x_1 &= x_0 - 0.05 \times m_1 = -2/8000 - 0.05 \times (-18/2944) = -1/8853, \\ y(x_1) &= 0.3 \times (-1/8853)^4 - 0.1 \times (-1/8853)^3 - 2 \times (-1/8853)^2 - 0.8 \times (-1/8853) \\ &= -1/1404 \end{aligned}$$

مرحله ۱:

$$\begin{aligned}
 g_2 &= 1/2 \times (-1/8853)^2 - 0/3 \times (-1/8853)^2 - 4 \times (-1/8853) - 0/8 = -2/3661, \\
 m_2 &= 0/7 \times m_1 + g_2 = 0/7 \times (-18/2944) + (-2/3661) = -15/1722, \\
 x_2 &= x_1 - 0/05 \times m_2 = -1/8853 - 0/05 \times (-15/1722) = -1/1267, \\
 y(x_2) &= 0/3 \times (-1/1267)^4 - 0/1 \times (-1/1267)^3 - 2 \times (-1/1267)^2 - 0/8 \times (-1/1267) \\
 &= -1/0110.
 \end{aligned}$$

مرحله ۲:

حل.  
سوال ۲  
(الف)

$$z_1 : f_{net} = \max \{ |x_1 + w_{z1,1}|, |x_2 + w_{z1,2}|, |x_3 + w_{z1,3}| \} \Rightarrow w_{z1,1} = -3 \quad w_{z1,2} = -5 \quad w_{z1,3} = 2$$

$$f_{act} = \exp\left(\frac{-net}{\theta}\right) \Rightarrow \theta = 4$$

$$f_{out} = f_{act}$$

$$z_2 : f_{net} = \sqrt{(x_1 + w_{z2,1})^2 + (x_3 + w_{z2,3})^2} \Rightarrow w_{z2,1} = 3 \quad w_{z2,3} = -4$$

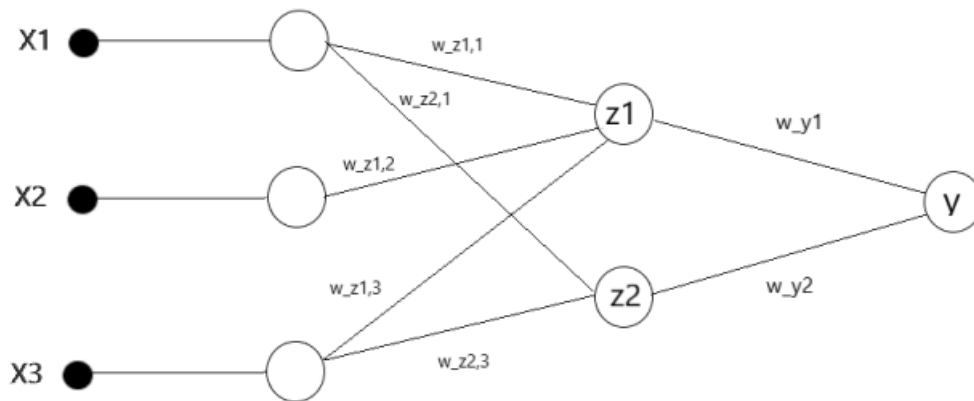
$$f_{act} = \begin{cases} 0, & \text{net} > \theta \\ 1, & \text{otherwise} \end{cases} \quad \theta = 4$$

$$f_{out} = f_{act}$$

$$y : f_{net} = w_{y1}z_1 + w_{y2}z_2 + \sigma_1 \Rightarrow w_{y1} = 6 \quad w_{y2} = 2 \quad \sigma_1 = 4$$

$$f_{act} = f_{net}$$

$$f_{out} = act$$



ب) برای طراحی این شبکه مسئله را به دو بخش تقسیم می‌کنیم.  
بخش اول (طراحی تابع  $x_1^{x_2}$ ):

در این مرحله از توالی نورون‌هایی با تابع فعالسازی لگاریتمی استفاده می‌کنیم. سپس با یک نورون با تابع فعالسازی نمایی برای ایجاد خروجی مدنظر استفاده می‌کنیم.

بخش دوم (طراحی تابع ماکسیمم):

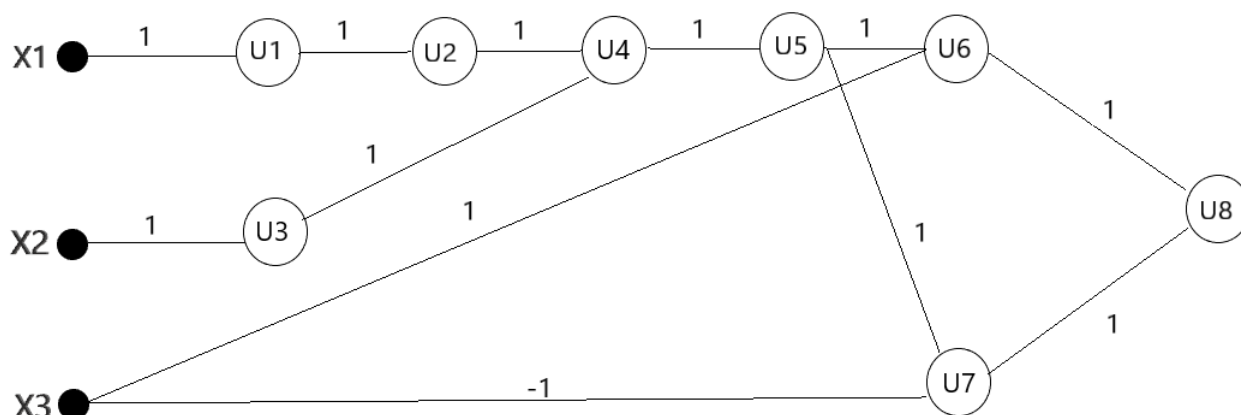
برای طراحی این بخش از عبارت زیر استفاده می‌کنیم:

$$\max(x_1, x_2) = \frac{x_1 + x_2 + |x_1 - x_2|}{2}$$

به کمک ترکیب کردن بخش‌های فوق به شبکه نهایی می‌رسیم.

$$U_{1,2,3} \begin{cases} f_{net}(x) = x \\ f_{act}(x) = \ln(x) \\ f_{out}(x) = x \end{cases} \quad U_4 \begin{cases} f_{net}(x, y) = x + y \\ f_{act}(x) = e^x \\ f_{out}(x) = x \end{cases} \quad U_5 \begin{cases} f_{net}(x) = x \\ f_{act}(x) = e^x \\ f_{out}(x) = x \end{cases}$$

$$U_6 \begin{cases} f_{net}(x, y) = x + y \\ f_{act}(x) = x \\ f_{out}(x) = x \end{cases} \quad U_7 \begin{cases} f_{net}(x, y) = x - y \\ f_{act}(x) = |x| \\ f_{out}(x) = x \end{cases} \quad U_8 \begin{cases} f_{net}(x, y) = x + y \\ f_{act}(x) = \frac{x}{y} \\ f_{out}(x) = x \end{cases}$$



حل.

سوال ۳

۱.  $W_1 \in R^{D_{a1} \times D_x}, b_1 \in R^{D_{a1} \times 1}, W_2 \in R^{1 \times D_{a1}}, b_2 \in R^{1 \times 1}$ . شکل وزن‌ها/بایاس‌ها بعد از برداری‌سازی یکسان خواهد بود.  $X \in R^{D_x \times m}, Y \in R^{m \times 1}$  بعد از برداری‌سازی.

$$\delta_1^{(i)} = y^{(i)} \hat{y}^{(i)} - (1 - y^{(i)}) / (1 - \hat{y}^{(i)}). \delta J / \delta \hat{y} = -\frac{1}{m} \sum_i \delta_1^{(i)} \quad 2.$$

$$\delta_2^{(i)} = \sigma(z_2)(1 - \sigma(z_2)) \quad 3.$$

$$\delta_3^{(i)} = W_2 \quad 4.$$

$$\delta_4^{(i)} = 0 \text{ اگر } z_1 < 0, 1 \text{ اگر } z_1 \geq 0 \quad 5.$$

$$\delta_5^{(i)} = x^{(i)T} \quad 6.$$

$$\delta_6 = -\frac{1}{m} \sum_i \delta_1^{(i)} * \delta_2^{(i)} * (\delta_3^{(i)} \circ \delta_4^{(i)}) * \delta_5^{(i)} \quad 7.$$

حل.

سوال ۴

**Pass Forward (آ**

لایه ۱ (ورودی به مخفی ۱)

محاسبه ورودی وزن‌دار  $z^{(1)}$ :

$$z^{(1)} = W^{(1)}x + b^{(1)} = \begin{bmatrix} a & -a \\ b & a \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix}$$

محاسبه خروجی فعال سازی  $a^{(1)}$  با استفاده از تابع **Sigmoid**:

$$a^{(1)} = \sigma(z^{(1)}) = \begin{bmatrix} \sigma(a) \\ \sigma(b) \end{bmatrix} = \begin{bmatrix} \frac{1}{1+e^{-a}} \\ \frac{1}{1+e^{-b}} \end{bmatrix}$$

لایه ۲ (مخفی ۱ به مخفی ۲)

محاسبه ورودی وزن دار  $z^{(2)}$ :

$$z^{(2)} = W^{(2)} \cdot a^{(1)} + b^{(2)} = \begin{bmatrix} a & b \\ -b & -a \end{bmatrix} \cdot \begin{bmatrix} \sigma(a) \\ \sigma(b) \end{bmatrix} = \begin{bmatrix} a \cdot \sigma(a) + b \cdot \sigma(b) \\ -b \cdot \sigma(a) - a \cdot \sigma(b) \end{bmatrix}$$

محاسبه عناصر  $z^{(2)}$ :

$$z_1^{(2)} = a \cdot \sigma(a) + b \cdot \sigma(b)$$

$$z_2^{(2)} = -b \cdot \sigma(a) - a \cdot \sigma(b)$$

محاسبه خروجی فعال سازی  $a^{(2)}$  با استفاده از تابع **ReLU**:

$$a^{(2)} = \text{ReLU}(z^{(2)}) = \begin{bmatrix} \max(0, z_1^{(2)}) \\ \max(0, z_2^{(2)}) \end{bmatrix} = \begin{bmatrix} \max(0, a \cdot \sigma(a) + b \cdot \sigma(b)) \\ \max(0, -b \cdot \sigma(a) - a \cdot \sigma(b)) \end{bmatrix}$$

لایه ۳ (مخفی ۲ به مخفی ۳)

محاسبه ورودی وزن دار  $z^{(3)}$ :

$$z^{(3)} = W^{(3)}a^{(2)} + b^{(3)} = \begin{bmatrix} a & -b \\ -a & b \end{bmatrix} a^{(2)} + \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

محاسبه عناصر  $z^{(3)}$ :

$$z_1^{(3)} = a \cdot a_1^{(2)} - b \cdot a_2^{(2)}$$

$$z_2^{(3)} = -a \cdot a_1^{(2)} + b \cdot a_2^{(2)}$$

محاسبه خروجی فعال سازی  $a^{(3)}$  با استفاده از تابع **Softmax**:

$$e^{z_1^{(3)}}, \quad e^{z_2^{(3)}}$$

$$a_i^{(3)} = \frac{e^{z_i^{(3)}}}{e^{z_1^{(3)}} + e^{z_2^{(3)}}}$$

لایه ۴ (مخفی ۳ به خروجی)

محاسبه ورودی وزن دار  $z^{(۴)}$ :

$$z^{(۴)} = W^{(۴)}a^{(۳)} + b^{(۴)} = \begin{bmatrix} a & -b \\ ۰.۵a & a \end{bmatrix} a^{(۳)} + \begin{bmatrix} ۰ \\ ۰ \end{bmatrix}$$

محاسبه عناصر  $z^{(۴)}$ :

$$\begin{aligned} z_1^{(۴)} &= a \cdot a_1^{(۳)} - b \cdot a_2^{(۳)} \\ z_2^{(۴)} &= ۰.۵a \cdot a_1^{(۳)} + a \cdot a_2^{(۳)} \end{aligned}$$

از آنجا که تابع فعال سازی لایه خروجی Identity است، داریم:

$$a^{(۴)} = z^{(۴)}$$

## ب) محاسبه خطا

با استفاده از تابع خطای میانگین مربعات (MSE):

$$E = \frac{1}{۲} \left( (a_1^{(۴)} - t_1)^2 + (a_2^{(۴)} - t_2)^2 \right)$$

## ج) Backward Pass

لایه ۴ (خروجی)

محاسبه  $\delta^{(۴)}$ :

تابع خطا به صورت زیر است:

$$E = \frac{1}{۲} \sum_{k=1}^۲ (a_k^{(۴)} - t_k)^2$$

گرادیان خطا نسبت به خروجی شبکه:

$$\frac{\partial E}{\partial a_k^{(۴)}} = a_k^{(۴)} - t_k$$

$$\frac{\partial a_k^{(۴)}}{\partial z_k^{(۴)}} = ۱$$

$$\delta_k^{(۴)} = \frac{\partial E}{\partial z_k^{(۴)}} = \frac{\partial E}{\partial a_k^{(۴)}} \cdot \frac{\partial a_k^{(۴)}}{\partial z_k^{(۴)}} = (a_k^{(۴)} - t_k) \times ۱ = a_k^{(۴)} - t_k$$

—

لایه ۳ (مخفی ۳)

محاسبه  $\delta^{(۳)}$ :

$$z^{(۴)} = W^{(۴)} a^{(۳)} + b^{(۴)}$$

$$\frac{\partial z_i^{(۴)}}{\partial a_j^{(۳)}} = W_{ij}^{(۴)}$$

$$\frac{\partial E}{\partial a_j^{(۳)}} = \sum_{i=۱}^۲ \frac{\partial E}{\partial z_i^{(۴)}} \cdot \frac{\partial z_i^{(۴)}}{\partial a_j^{(۳)}} = \sum_{i=۱}^۲ \delta_i^{(۴)} W_{ij}^{(۴)} = \sum_k (a_k^{(۴)} - t_k) \cdot W_{ki}^{(۴)}$$

$$\frac{\partial a_l^{(۳)}}{\partial z_k^{(۳)}} = a_l^{(۳)} (\delta_{lk} - a_k^{(۳)})$$

که در آن  $\delta_{lk}$  دلتای کرونکر است که اگر  $l = k$  باشد برابر ۱ و در غیر این صورت برابر ۰ است.

$$\frac{\partial a_i^{(۳)}}{\partial z_j^{(۳)}} = \begin{cases} a_i^{(۳)} (1 - a_i^{(۳)}) & \text{if } i = j \\ -a_i^{(۳)} \cdot a_j^{(۳)} & \text{if } i \neq j \end{cases}$$

بنابراین، گرادیان خطا نسبت به  $z_k^{(۳)}$ :

$$\delta_k^{(۳)} = \frac{\partial E}{\partial z_k^{(۳)}} = \sum_{j=۱}^۲ \frac{\partial E}{\partial a_j^{(۳)}} \cdot \frac{\partial a_j^{(۳)}}{\partial z_k^{(۳)}}$$

$$\delta_k^{(۳)} = \sum_{j=۱}^۲ \left( \frac{\partial E}{\partial a_j^{(۳)}} \cdot \begin{cases} a_j^{(۳)} (1 - a_j^{(۳)}) & \text{if } j = k \\ -a_j^{(۳)} a_k^{(۳)} & \text{if } j \neq k \end{cases} \right)$$

به دو بخش تقسیم می‌شود:

$$\delta_k^{(۳)} \text{ (for } j = k) = \frac{\partial E}{\partial a_k^{(۳)}} \cdot a_k^{(۳)} (1 - a_k^{(۳)})$$

$$\delta_k^{(۳)} \text{ (for } j \neq k) = \sum_{j \neq k} \left( \frac{\partial E}{\partial a_j^{(۳)}} \cdot (-a_j^{(۳)} a_k^{(۳)}) \right)$$

با ترکیب دو بخش:

$$\delta_k^{(۳)} = \sum_{l=۱}^۲ \frac{\partial E}{\partial a_l^{(۳)}} \cdot \frac{\partial a_l^{(۳)}}{\partial z_k^{(۳)}} = \sum_{l=۱}^۲ \left( \frac{\partial E}{\partial a_l^{(۳)}} \cdot a_l^{(۳)} (\delta_{lk} - a_k^{(۳)}) \right)$$

$$\delta_k^{(۳)} = \frac{\partial E}{\partial a_k^{(۳)}} \cdot a_k^{(۳)} (1 - a_k^{(۳)}) - a_k^{(۳)} \sum_{j \neq k} \left( \frac{\partial E}{\partial a_j^{(۳)}} a_j^{(۳)} \right)$$

(تا بالای اینجا اگر نوشته شود نمره کامل گرفته می‌شود.)

در ادامه می‌توان نوشت:

$$\sum_{j \neq k} \left( \frac{\partial E}{\partial a_j^{(\mathbf{r})}} a_j^{(\mathbf{r})} \right) = \left( \sum_{j=1}^{\mathbf{r}} \frac{\partial E}{\partial a_j^{(\mathbf{r})}} a_j^{(\mathbf{r})} \right) - \frac{\partial E}{\partial a_k^{(\mathbf{r})}} a_k^{(\mathbf{r})}$$

$$S = \sum_{j=1}^{\mathbf{r}} \frac{\partial E}{\partial a_j^{(\mathbf{r})}} a_j^{(\mathbf{r})}$$

$$\delta_k^{(\mathbf{r})} = \frac{\partial E}{\partial a_k^{(\mathbf{r})}} a_k^{(\mathbf{r})} (1 - a_k^{(\mathbf{r})}) - a_k^{(\mathbf{r})} \left( S - \frac{\partial E}{\partial a_k^{(\mathbf{r})}} a_k^{(\mathbf{r})} \right)$$

$$\delta_k^{(\mathbf{r})} = \frac{\partial E}{\partial a_k^{(\mathbf{r})}} a_k^{(\mathbf{r})} (1 - a_k^{(\mathbf{r})}) - a_k^{(\mathbf{r})} S + \frac{\partial E}{\partial a_k^{(\mathbf{r})}} a_k^{(\mathbf{r})} \mathbf{r}$$

$$\delta_k^{(\mathbf{r})} = \frac{\partial E}{\partial a_k^{(\mathbf{r})}} a_k^{(\mathbf{r})} - a_k^{(\mathbf{r})} S$$

$$\delta_k^{(\mathbf{r})} = a_k^{(\mathbf{r})} \left( \frac{\partial E}{\partial a_k^{(\mathbf{r})}} - S \right)$$

$$\delta^{(\mathbf{r})} = a^{(\mathbf{r})} \odot \left( \frac{\partial E}{\partial a^{(\mathbf{r})}} - S \mathbf{1} \right)$$

لایه ۲ (مخفی ۲)

محاسبه  $\delta^{(2)}$ :

$$\frac{\partial E}{\partial a_j^{(\mathbf{r})}} = \sum_{k=1}^{\mathbf{r}} \delta_k^{(\mathbf{r})} \cdot \frac{\partial z_k^{(\mathbf{r})}}{\partial a_j^{(\mathbf{r})}}$$

$$z^{(\mathbf{r})} = W^{(\mathbf{r})} a^{(\mathbf{r})} + b^{(\mathbf{r})}$$

$$\frac{\partial z_k^{(\mathbf{r})}}{\partial a_j^{(\mathbf{r})}} = W_{kj}^{(\mathbf{r})}$$

$$\frac{\partial E}{\partial a_j^{(\mathbf{r})}} = \sum_{k=1}^{\mathbf{r}} \delta_k^{(\mathbf{r})} W_{kj}^{(\mathbf{r})}$$

$$\frac{\partial a_j^{(\mathbf{r})}}{\partial z_j^{(\mathbf{r})}} = \begin{cases} 1 & z_j^{(\mathbf{r})} > 0 \\ 0 & \text{o.w.} \end{cases}$$

$$\delta_j^{(\mathbf{r})} = \frac{\partial E}{\partial a_j^{(\mathbf{r})}} \cdot \frac{\partial a_j^{(\mathbf{r})}}{\partial z_j^{(\mathbf{r})}} = \begin{cases} \frac{\partial E}{\partial a_j^{(\mathbf{r})}} & z_j^{(\mathbf{r})} > 0 \\ 0 & \text{o.w.} \end{cases}$$



لایه ۱ (مخفی ۱)

محاسبه  $\delta^{(1)}$ :

$$\frac{\partial E}{\partial a_j^{(1)}} = \sum_{k=1}^{\mathfrak{Y}} \delta_k^{(\mathfrak{Y})} \cdot \frac{\partial z_k^{(\mathfrak{Y})}}{\partial a_j^{(1)}}$$

$$z^{(\mathfrak{Y})} = W^{(\mathfrak{Y})} a^{(1)} + b^{(\mathfrak{Y})}$$

$$\frac{\partial z_k^{(\mathfrak{Y})}}{\partial a_j^{(1)}} = W_{kj}^{(\mathfrak{Y})}$$

$$\frac{\partial E}{\partial a_j^{(1)}} = \sum_{k=1}^{\mathfrak{Y}} \delta_k^{(\mathfrak{Y})} W_{kj}^{(\mathfrak{Y})}$$

$$\frac{\partial a_j^{(1)}}{\partial z_j^{(1)}} = a_j^{(1)} (1 - a_j^{(1)})$$

$$\delta_j^{(1)} = \frac{\partial E}{\partial a_j^{(1)}} \cdot \frac{\partial a_j^{(1)}}{\partial z_j^{(1)}} = \left( \sum_{k=1}^{\mathfrak{Y}} \delta_k^{(\mathfrak{Y})} W_{kj}^{(\mathfrak{Y})} \right) \cdot a_j^{(1)} (1 - a_j^{(1)})$$

محاسبه گرادیان‌ها نسبت به وزن‌ها و بایاس‌ها

لایه ۴:

$$\frac{\partial E}{\partial W_{ij}^{(\mathfrak{F})}} = \delta_i^{(\mathfrak{F})} \cdot a_j^{(\mathfrak{F})}$$

$$\frac{\partial E}{\partial b_i^{(\mathfrak{F})}} = \delta_i^{(\mathfrak{F})}$$

لایه ۳:

$$\frac{\partial E}{\partial W_{ij}^{(\mathfrak{R})}} = \delta_i^{(\mathfrak{R})} \cdot a_j^{(\mathfrak{R})}$$

$$\frac{\partial E}{\partial b_i^{(\mathfrak{R})}} = \delta_i^{(\mathfrak{R})}$$

لایه ۲:

$$\frac{\partial E}{\partial W_{ij}^{(\mathfrak{Y})}} = \delta_i^{(\mathfrak{Y})} \cdot a_j^{(1)}$$

$$\frac{\partial E}{\partial b_i^{(r)}} = \delta_i^{(r)}$$

لایه ۱:

$$\frac{\partial E}{\partial W_{ij}^{(1)}} = \delta_i^{(1)} \cdot x_j$$

$$\frac{\partial E}{\partial b_i^{(1)}} = \delta_i^{(1)}$$

### د) به روزرسانی وزن ها و بایاس ها

با استفاده از نرخ یادگیری  $\eta$ ، وزن ها و بایاس ها را به روزرسانی می کنیم:

$$W_{ij}^{(l)} \leftarrow W_{ij}^{(l)} - \eta \frac{\partial E}{\partial W_{ij}^{(l)}}$$

$$b_i^{(l)} \leftarrow b_i^{(l)} - \eta \frac{\partial E}{\partial b_i^{(l)}}$$

حل.  
سوال ۵

الف) هدف از این طراحی، ایجاد یک شبکه با استفاده از واحدهای خطی آستانه‌ای (TLUs) است که بتواند تعیین کند آیا یک نقطه در نواحی رنگی مشخص شده در تصویر قرار دارد یا خیر. این نواحی به شکل لوزی و متقارن نسبت به مبدا هستند. در تصویر چهار ناحیه لوزی شکل وجود دارد که مراکز آنها در مختصات  $(2, 2)$ ،  $(-2, 2)$ ،  $(2, -2)$  و  $(-2, -2)$  قرار دارند. این لوزی ها می توانند به صورت مجموعه ای از نابرابری های خطی توصیف شوند که مرزهای آنها را تعریف می کنند. برای هر لوزی، می توان مرزهای آن را با استفاده از معادله زیر نمایش داد:

$$|x - x_c| + |y - y_c| \leq d$$

که در آن:

$(x_c, y_c)$  مرکز لوزی است،  $d$  فاصله مرکز تا رأس لوزی (نصف طول قطر) است.

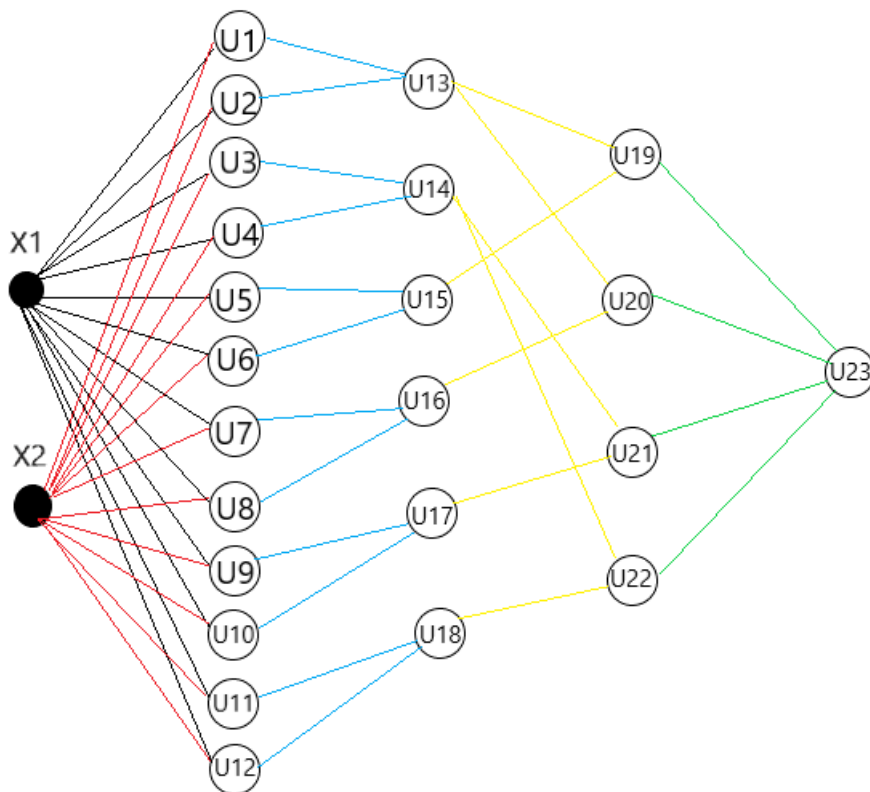
برای هر لوزی، چهار نابرابری خطی برای تعیین مرزها وجود دارد. هر واحد TLU می تواند یک نابرابری خطی را نشان دهد. به عنوان مثال، برای یک لوزی، نابرابری های زیر را می توان استفاده کرد:

$$\begin{aligned} x + y &\leq b_1 \\ x - y &\leq b_2 \\ -x + y &\leq b_3 \\ -x - y &\leq b_4 \end{aligned}$$

در مجموع، برای هر لوزی چهار TLU نیاز است که هر کدام یک مرز را نمایش می دهند. اگر یک نقطه تمام این نابرابری ها را برآورده کند، یعنی در داخل لوزی قرار دارد. برای تعیین اینکه آیا نقطه ای در هر یک

از لوزی‌ها قرار دارد، خروجی‌های مربوط به هر لوزی باید با یک تابع OR ترکیب شوند. اگر خروجی یکی از لوزی‌ها ۱ باشد، نتیجه نهایی نیز ۱ خواهد بود، یعنی نقطه در یکی از نواحی رنگی قرار دارد. ساختار نهایی شبکه به این صورت است:

- **لایه ورودی:** شامل دو ورودی  $x$  و  $y$  که مختصات نقطه هستند.
- **لایه مخفی اول:** شامل ۱۲ واحد TLU (۴ واحد برای هر لوزی) که شرایط مرزی لوزی‌ها را بررسی می‌کنند. (توجه کنید تعدادی از خطوط بین نواحی مشترک هستند).
- **لایه مخفی دوم:** شامل ۶ واحد TLU که قرار گرفتن ورودی بین جفت‌های خطوط موازی شکل را بررسی می‌کند. (۶ جفت خط موازی داریم).
- **لایه مخفی سوم:** شامل ۴ واحد TLU که قرار گرفتن ورودی درون هر یک از نواحی را بررسی می‌کند.
- **لایه خروجی:** ترکیب نهایی نتایج لایه مخفی برای تعیین اینکه آیا نقطه در داخل یکی از لوزی‌ها قرار دارد یا خیر.

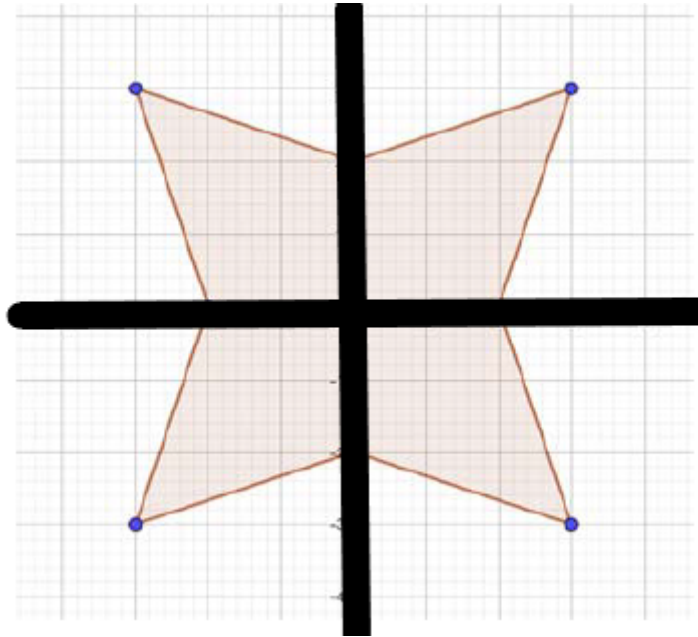


$$U_1 \begin{cases} W_{U_1, X_1} = -1 \\ W_{U_1, X_2} = -1 \\ \theta_1 = -1 \end{cases} \quad U_2 \begin{cases} W_{U_2, X_1} = -1 \\ W_{U_2, X_2} = -1 \\ \theta_2 = 1 \end{cases} \quad U_3 \begin{cases} W_{U_3, X_1} = -1 \\ W_{U_3, X_2} = 1 \\ \theta_3 = -1 \end{cases} \quad U_4 \begin{cases} W_{U_4, X_1} = -1 \\ W_{U_4, X_2} = 1 \\ \theta_4 = 1 \end{cases}$$

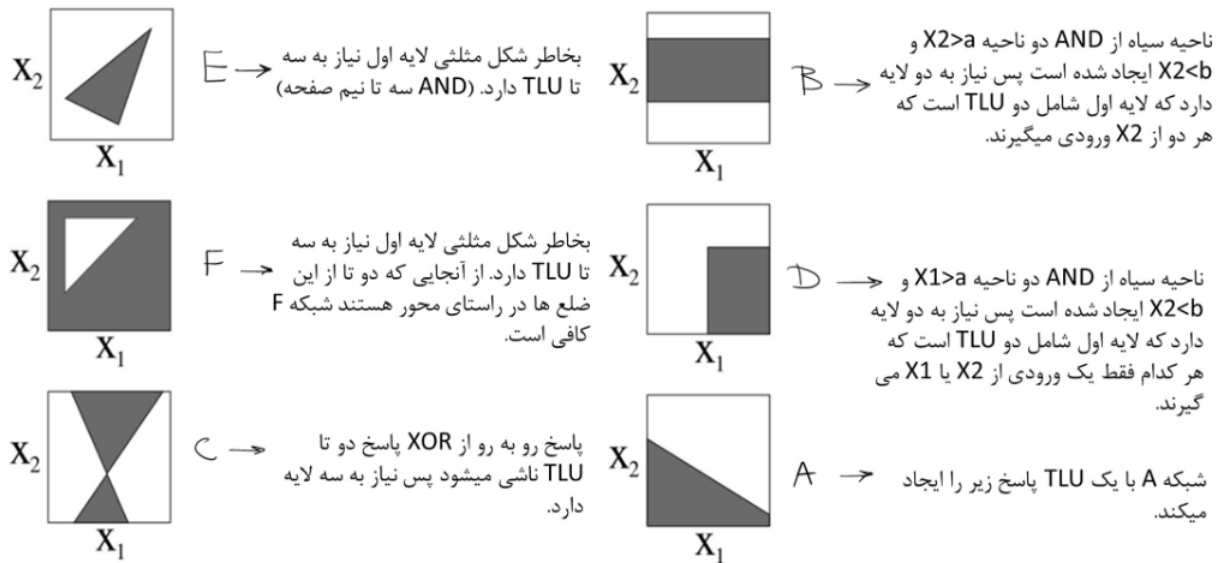
$$U_5 \begin{cases} W_{U_5, X_1} = -1 \\ W_{U_5, X_2} = 1 \\ \theta_5 = 3 \end{cases} \quad U_6 \begin{cases} W_{U_6, X_1} = 1 \\ W_{U_6, X_2} = -1 \\ \theta_6 = -5 \end{cases} \quad U_7 \begin{cases} W_{U_7, X_1} = 1 \\ W_{U_7, X_2} = -1 \\ \theta_7 = 3 \end{cases} \quad U_8 \begin{cases} W_{U_8, X_1} = -1 \\ W_{U_8, X_2} = 1 \\ \theta_8 = -5 \end{cases}$$

$$\begin{aligned}
U_9 & \begin{cases} W_{U_9, X_1} = 1 \\ W_{U_9, X_2} = 1 \\ \theta_9 = 3 \end{cases} & U_{10} & \begin{cases} W_{U_{10}, X_1} = -1 \\ W_{U_{10}, X_2} = -1 \\ \theta_{10} = -5 \end{cases} & U_{11} & \begin{cases} W_{U_{11}, X_1} = -1 \\ W_{U_{11}, X_2} = -1 \\ \theta_{11} = 3 \end{cases} & U_{12} & \begin{cases} W_{U_{12}, X_1} = 1 \\ W_{U_{12}, X_2} = 1 \\ \theta_{12} = -5 \end{cases} \\
U_{13} & \begin{cases} W_{U_{13}, U_1} = 2 \\ W_{U_{13}, U_2} = 2 \\ \theta_{13} = 3 \end{cases} & U_{14} & \begin{cases} W_{U_{14}, U_3} = 2 \\ W_{U_{14}, U_4} = 2 \\ \theta_{14} = 3 \end{cases} & U_{15} & \begin{cases} W_{U_{15}, U_5} = 2 \\ W_{U_{15}, U_6} = 2 \\ \theta_{15} = 3 \end{cases} & U_{16} & \begin{cases} W_{U_{16}, U_7} = 2 \\ W_{U_{16}, U_8} = 2 \\ \theta_{16} = 3 \end{cases} \\
U_{17} & \begin{cases} W_{U_{17}, U_9} = 2 \\ W_{U_{17}, U_{10}} = 2 \\ \theta_{17} = 3 \end{cases} & U_{18} & \begin{cases} W_{U_{18}, U_{11}} = 2 \\ W_{U_{18}, U_{12}} = 2 \\ \theta_{18} = 3 \end{cases} & U_{19} & \begin{cases} W_{U_{19}, U_{13}} = 2 \\ W_{U_{19}, U_{15}} = 2 \\ \theta_{19} = 3 \end{cases} & U_{20} & \begin{cases} W_{U_{20}, U_{13}} = 2 \\ W_{U_{20}, U_{16}} = 2 \\ \theta_{20} = 3 \end{cases} \\
U_{21} & \begin{cases} W_{U_{21}, U_{14}} = 2 \\ W_{U_{21}, U_{17}} = 2 \\ \theta_{21} = 3 \end{cases} & U_{22} & \begin{cases} W_{U_{22}, U_{14}} = 2 \\ W_{U_{22}, U_{18}} = 2 \\ \theta_{22} = 3 \end{cases} & U_{23} & \begin{cases} W_{U_{23}, U_{19}} = 2 \\ W_{U_{23}, U_{20}} = 2 \\ W_{U_{23}, U_{21}} = 2 \\ W_{U_{23}, U_{22}} = 2 \\ \theta_{23} = 1 \end{cases}
\end{aligned}$$

ب) شکل این بخش را می‌توان به چهار ناحیه زیر تقسیم کرد.



هر بخش از تقاطع چهار خط بدست می‌آید. از آنجاییکه تعدادی از خطوط مشترک هستند، در واقع ده خط متفاوت داریم. بنابراین می‌توانیم با تغییر ضرایب و حد آستانه‌ها این طبقه‌بندی را انجام دهیم.  
(ج)



## حل. سوال ۶

(الف) با استفاده از استقرا این قضیه را اثبات می‌کنیم. نورون‌های لایه اول  $v_i^{(1)}$  را در نظر بگیرید. درایه یا سطر  $i$ ام را برای پارامترهای نورون یعنی  $W_i^{(1)}$  و  $b_i^{(1)}$  در نظر بگیرید. حال اگر تصویر ورودی نسبت به لایه اول شبکه را تشکیل دهیم، هر نورون یک ابرصفحه به فرم  $W_i^{(1)}x + b_i^{(1)} = 0$  در فضا تشکیل می‌دهد. در نهایت در فضا تعدادی ابرصفحه خواهیم داشت که ناحیه‌های مختلف را تشکیل می‌دهند. هر polytope تشکیل شده مخصوص یک activation pattern خواهد بود.

حال فرض کنید لایه‌های ۱ تا  $d-1$  فضا را به polytope ها تقسیم کرده‌اند. نورون  $v_i^{(d)}$  و ناحیه  $R_i$  را در نظر بگیرید. در ناحیه  $R_i$  یک activation pattern ثابت داریم. اگر دقت کنید متوجه می‌شوید که ورودی‌های نورون در واقع ترکیب خطی‌ای از تعدادی از ورودی‌ها به فرم  $\sum_j \lambda_j x_j + b$  می‌باشد. چرا که تابع ReLU یا خود ورودی را که ترکیب خطی‌ای از ورودی‌هاست خروجی می‌دهد و یا صفر خروجی می‌دهد. با قرار دادن این رابطه برابر ۰ به یک ابرصفحه دیگر می‌رسیم که البته تنها در ناحیه  $R_i$  صحت دارد. حال یا این صفحه در صورت برخورد با ناحیه آن را به دو ناحیه تقسیم می‌کند و یا تغییری در pattern خروجی این نورون نمی‌دهد. بنابراین همچنان فضا به convex polytope ها تقسیم می‌شود و قضیه در نتیجه آن اثبات می‌شود.

(ب) مشابه بخش قبل با استقرا ثابت می‌کنیم. طبق بخش الف، لایه اول شبکه فضا را به  $r(k, m)$  ناحیه تقسیم می‌کند که طبق راهنمایی تعداد این ناحیه‌ها کمتر از  $O(k^m)$  است. حال فرض کنید تعداد نواحی ایجاد شده توسط لایه ۱- $n$  ام کمتر از  $O(k^{m(n-1)})$  است. حال از آنجایی که هر ناحیه دارای pattern activation ثابت است طبق بخش الف هر کدام از این ناحیه‌ها می‌توانند به حداکثر  $r(k, m)$  ناحیه دیگر تقسیم شوند. یعنی تعداد نواحی ما حداکثر  $O(k^{m(n-1)}) \times r(k, m)$  است که خود این مقدار کمتر از  $O(k^{mn})$  می‌باشد.