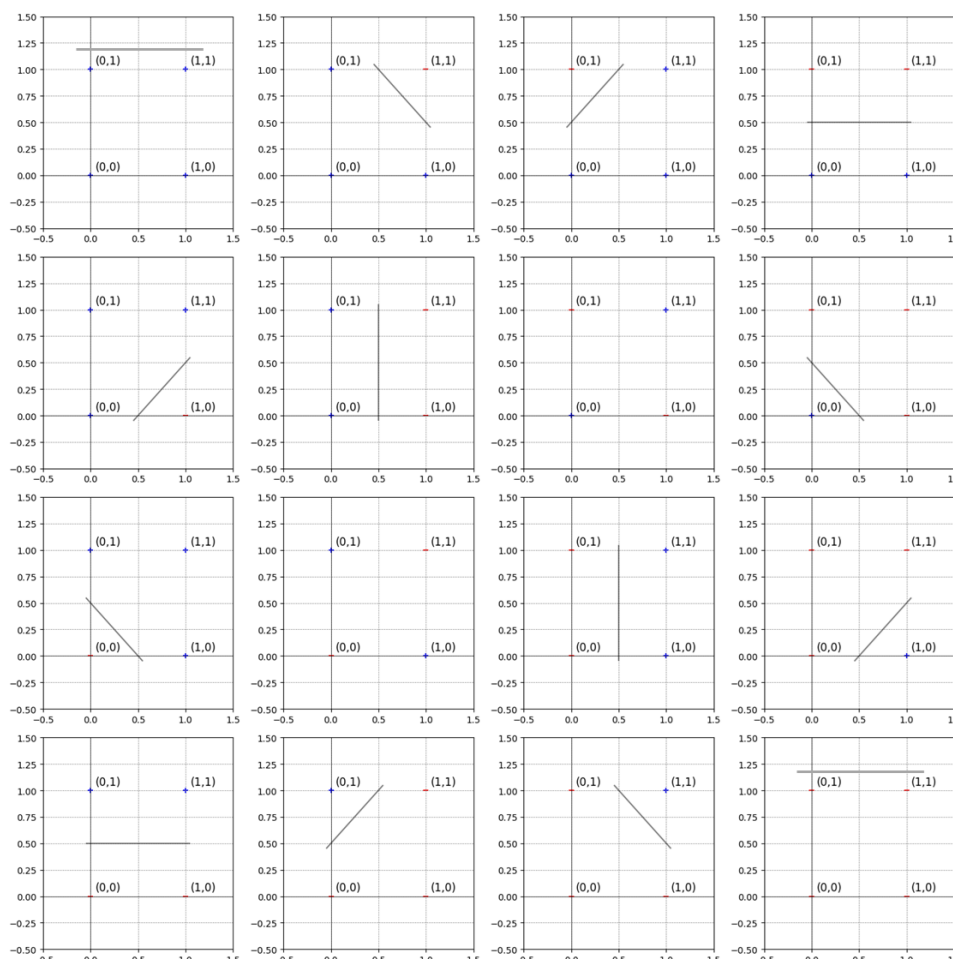


## تمرین سوم

### پاسخ مسئله ۱.

#### الف) ۱۴ تابع.

در شکل زیر، تمامی ۱۶ حالت ممکن برای خروجی یک تابع boolean بر روی دو ورودی را مشاهده می کنید. نقاط موجود در صفحه، نشان دهنده ورودی های تابع هستند (به طور مثال نقطه  $(0, 0)$  در صفحه نشان دهنده صفر بودن هر دو ورودی تابع است). همچنین اگر یک نقطه با علامت "+" مشخص شده باشد، یعنی خروجی تابع برای آن ورودی برابر یک بوده و اگر با علامت "-" مشخص شده باشد، یعنی خروجی تابع برای آن ورودی برابر صفر است.

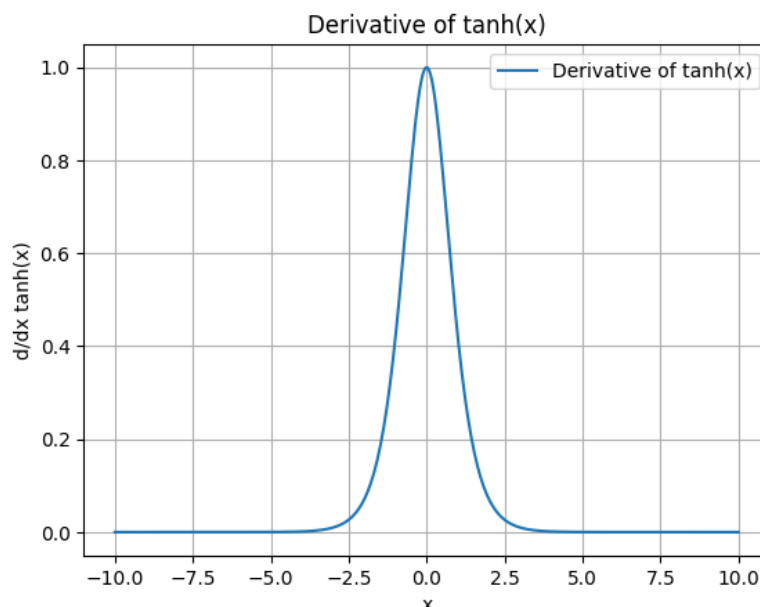


شکل ۱: حالات ممکن برای یک تابع boolean با ۲ ورودی

همانطور که مشاهده می شود، تنها در دو حالتی که کلاس های مختلف به شکل قطری قرار گرفته باشند مجموعه داده حاصل خطی جدایی پذیر نخواهد بود. پس پاسخ مسئله ۱۴ می باشد.

ب) نادرست.

در شکل زیر نمودار  $\frac{d(\tanh x)}{dx} = 1 - \tanh^2 x$  را مشاهده می کنید. به وضوح خروجی این مشتق همواره بین صفر و یک بوده و به خصوص برای  $x$  هایی که قدر مطلقشان از یک بزرگ تر است، بسیار نزدیک به صفر می باشد. بنابراین در صورت استفاده از این تابع فعال سازی در ساختار یک شبکه عصبی، هنگام ضرب شدن مشتقات جزئی در فرآیند پس انتشار، اندازه گرادیان بسیار کوچک می شود و ناپدید شدن گرادیان رخ می دهد.



شکل ۲: نمودار  $\frac{d(\tanh x)}{dx}$

ج) درست.

براساس قضیه تقریب پذیری (Universal Approximation Theorem)، شبکه های عصبی با یک لایه مخفی، تعداد کافی نورون و یک تابع فعال سازی غیر خطی (در اینجا RELU) قادر به تقریب زنی هر تابع پیوسته ای بر روی مجموعه های محدود هستند. اگرچه توابع منطقی در اصل گسسته و غیر پیوسته هستند، اما می توان آن ها را به صورت ترکیبی از توابع پیوسته مدل سازی کرد. در واقع هر تابع منطقی را می توان به شکل SOP (or تعدادی عبارت به شکل and تعدادی متغیر دودویی) نوشت، از طرفی می دانیم که برای متغیر های دودویی  $x_1, x_2$  داریم:

$$x_1 \wedge x_2 = x_1 x_2$$

$$x_1 \vee x_2 = \frac{|x_1 + x_2| + |x_1 - x_2|}{2}$$

که عبارات بالا به راحتی قابل تعمیم برای  $n$  متغیر دودویی می باشند.

د) نادرست.

عرض و عمق، دو مورد از هایپرپارامتر های یک شبکه عصبی هستند و باید مقدار مناسب آن ها را پیدا کرد (به عبارت دیگر باید tune شوند). افزایش بیش از اندازه این هایپرپارامتر ها، پیچیدگی شبکه عصبی را زیاد می کند و می تواند باعث بیش برآزش، افزایش خطای تعمیم (generalization error) و در نتیجه بدتر شدن عملکرد شبکه عصبی شود.

## ه) نادرست.

شبکه‌های عصبی عمیق دارای چندین لایه مخفی هستند که به آن‌ها امکان می‌دهد ویژگی‌های پیچیده‌تر و انتزاعی‌تری را از داده‌ها استخراج کنند. این ویژگی باعث می‌شود که شبکه‌های عمیق بتوانند توابع پیچیده را با تعداد کمتری نورون نسبت به شبکه‌های سطحی مدل‌سازی کنند.

اگرچه شبکه‌های عصبی سطحی و بسیار پهن قادر به تقریب زدن خروجی شبکه‌های عمیق و به دست آوردن عملکرد تقریباً مشابه هستند (با توجه به Universal Approximation Theorem)، اما این امر به معنای کارآمد بودن آنها نیست. در بسیاری از موارد، شبکه‌های سطحی برای مدل‌سازی توابع پیچیده نیاز به تعداد بسیار زیادی نورون دارند که این موضوع از نظر محاسباتی و زمانی غیرعملی است. در مقابل، شبکه‌های عمیق می‌توانند همان توابع را با تعداد کمتری نورون و لایه به طور کارآمدتری مدل‌سازی کنند.

و)

ناپدید شدن گرادیان‌ها زمانی رخ می‌دهد که گرادیان‌های محاسبه شده در فرآیند پس‌انتشار به تدریج کاهش یافته و به صفر نزدیک می‌شوند. این اتفاق باعث می‌شود که وزن‌های لایه‌های اولیه شبکه به‌طور بسیار کند به‌روزرسانی شوند و یا اصلاً به‌روزرسانی نشوند و در نتیجه یادگیری شبکه متوقف شود.

توابع فعال‌سازی مانند ReLU که دارای مشتق با اندازه بزرگ‌تر یا مساوی یک در ناحیه مثبت هستند، از کاهش گرادیان‌ها در لایه‌های عمیق جلوگیری می‌کنند (چرا که دیگر گرادیان‌ها در حین فرآیند پس‌انتشار، مکرراً در مشتق‌های جزئی با اندازه کوچک‌تر از یک ضرب نمی‌شوند و اندازه‌شان بسیار کوچک نمی‌شود). این ویژگی اجازه می‌دهد که گرادیان‌های قوی‌تر به‌طور مؤثرتری از لایه‌های بالاتر به پایین‌تر منتقل شوند. در نتیجه، شبکه‌های عصبی با استفاده از ReLU سریع‌تر و کارآمدتر آموزش می‌بینند و مشکل ناپدید شدن گرادیان‌ها کاهش می‌یابد.

## ز) نادرست.

هنگام استفاده از SGD، گام‌هایی که برداشته می‌شوند بسیار سریع هستند (فرکانس به‌روزرسانی زیاد است) اما به دلیل استفاده از تنها یک نمونه آموزشی برای به‌روزرسانی وزن‌ها، گام‌های برداشته شده دقت بالایی ندارند و حتی ممکن است پارامترهای شبکه را از پارامترهای بهینه دورتر کنند.

در مقابل، Mini-Batch GD دارای فرکانس به‌روزرسانی کمتری است (گام‌ها کندتر برداشته می‌شوند) اما چون از تعداد بیشتری نمونه آموزشی برای به‌روزرسانی وزن‌ها استفاده می‌کند، تابع هزینه حاصل، تقریب دقیق‌تری از تابع هزینه برای کل نمونه‌ها می‌باشد و در نتیجه گام برداشته شده دقیق‌تر خواهد بود.

بنابراین می‌توان نتیجه گرفت که معمولاً SGD نسبت به Mini-Batch GD، نیاز به تعداد به‌روزرسانی‌های بیشتری برای همگرایی دارد و در نتیجه حتی ممکن است دیرتر همگرا شود.

ی) مراحل به‌روزرسانی در ادامه به ترتیب آورده شده‌اند.

۱. محاسبه گرادیان تابع: گرادیان تابع نسبت به  $x$  محاسبه می‌شود:

$$\frac{dy}{dx} = 1/2x^3 - 0/3x^2 - 4x - 0/8$$

۲. الگوریتم SGD با ممان: فرمول به‌روزرسانی مقدار  $x$  با استفاده از ممان به صورت زیر است:

$$m_i = \mu m_{i-1} + (1 - \mu)g_{i-1} \quad , \quad x_i = x_{i-1} - \gamma m_i$$

۳. محاسبات مرحله به مرحله:

مرحله ۱: محاسبه  $y_0$  و  $g_0$  در  $x_0$

$$y_0 = 0.3(-2/8)^4 - 0.1(-2/8)^3 - 2(-2/8)^2 - 0.8(-2/8) = 7/1949$$

$$g_0 = 1/2(-2/8)^3 - 0.3(-2/8)^2 - 4(-2/8) - 0.8 = -18/2944$$

مرحله ۲: محاسبه  $m_1$  و  $x_1$

فرض می‌کنیم  $m_0 = 0$ :

$$m_1 = \mu m_0 + (1 - \mu)g_0 = 0.7(0) + 0.3(-18/2944) = -5/4883$$

$$x_1 = x_0 - \gamma m_1 = -2/8 - 0.05(-5/4883) = -2/5256$$

مرحله ۳: محاسبه  $y_1$  و  $g_1$  در  $x_1$

$$y_1 = 0.3(-2/5256)^4 - 0.1(-2/5256)^3 - 2(-2/5256)^2 - 0.8(-2/5256) = 3/0803$$

$$g_1 = 1/2(-2/5256)^3 - 0.3(-2/5256)^2 - 4(-2/5256) - 0.8 = -11/9428$$

مرحله ۴: محاسبه  $m_2$  و  $x_2$

$$m_2 = \mu m_1 + (1 - \mu)g_1 = 0.7(-5/4883) + 0.3(-11/9428) = -7/4246$$

$$x_2 = x_1 - \gamma m_2 = -2/5256 - 0.05(-7/4246) = -2/1543$$

مرحله ۵: محاسبه  $y_2$  در  $x_2$

$$y_2 = 0.3(-2/1543)^4 - 0.1(-2/1543)^3 - 2(-2/1543)^2 - 0.8(-2/1543) = -0.0971$$

نتایج نهایی

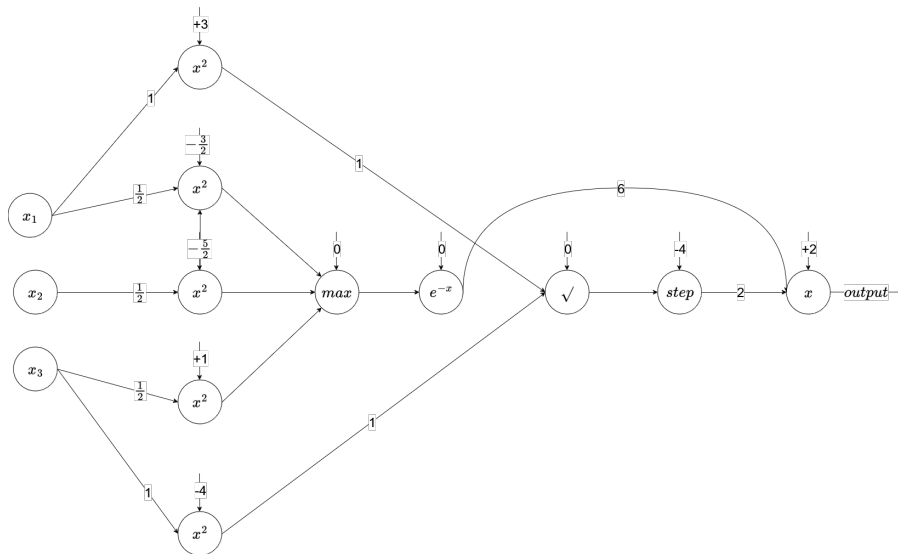
$$x_0 = -2/8, \quad m_0 = 0, \quad g_0 = -18/2944, \quad y_0 = 7/1949$$

$$x_1 = -2/5256, \quad m_1 = -5/4883, \quad g_1 = -11/9428, \quad y_1 = 3/0803$$

$$x_2 = -2/1543, \quad m_2 = -7/4246, \quad y_2 = -0.0971$$

## پاسخ مسئله‌ی ۲.

الف) در شکل زیر یک شبکه عصبی مطلوب را مشاهده می کنید.



شکل ۳: یک شبکه عصبی مطلوب

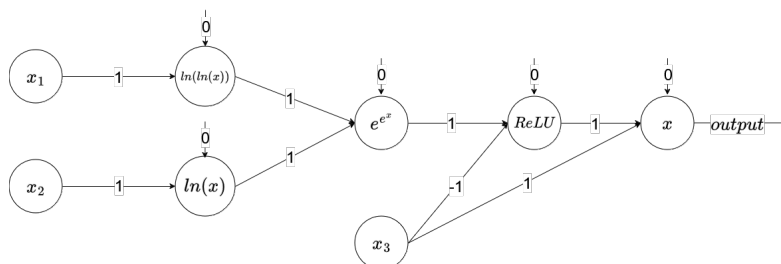
دقت کنید که تابع فعال سازی هر نورون داخل آن و بایاس هر نورون در بالای آن نوشته شده است. همچنین در نورون هایی که یال های ورودی وزن دار دارند، تابع ورودی جمع وزن دار ورودی ها است؛ اما در نورون هایی که یال های ورودیشان بدون وزن است، تابع ورودی همانی است (به عبارت دیگر همه ورودی ها مستقیماً به تابع فعال سازی وارد می شوند). همچنین توجه نمایید که می توان بجای نورون با تابع فعال سازی  $\max$ ، از دو نورون با تابع فعال سازی ReLU استفاده کرد که  $\max(a, b, c) = \text{ReLU}(\text{ReLU}(a - b) + b - c) + c$  را خروجی می دهند.

ب) دقت کنید که می توان نوشت: ( فرض می کنیم که ورودی های  $x_1$  و  $x_2$  مثبت اند، چرا که در غیر این صورت تابع  $x_1^{x_2}$  ممکن است تعریف نشود. به علاوه، بدون این فرض ساختن آن با پیچیدگی بسیار زیادی همراه خواهد بود)

$$x_1^{x_2} = e^{x_2 \ln x_1} = e^{\ln x_2 + \ln \ln x_1}$$

$$\max(a, b) = \text{ReLU}(a - b) + b$$

با توجه به تساوی های بالا، یک شبکه عصبی مطلوب به شکل زیر می باشد. مشابه قسمت قبل، تابع فعال سازی و بایاس هر نورون به ترتیب درون و بالای آن نوشته شده اند.



شکل ۴: یک شبکه عصبی مطلوب

### پاسخ مسئله‌ی ۳.

۱. ابعاد پارامترهای ذکر شده به شکل زیر است: (دقت کنید که چون  $z_2$  به عنوان ورودی به تابع  $\sigma(x)$  داده شده است، باید یک عدد اسکالر باشد)

$$\begin{aligned}W_1 &\rightarrow D_{a_1} \times D_x \\b_1 &\rightarrow D_{a_1} \times 1 \\W_2 &\rightarrow 1 \times D_{a_1} \\b_2 &\rightarrow 1 \times 1\end{aligned}$$

حال فرض کنید بخواهیم شبکه را بر روی  $k$  نمونه vectorize کنیم. در این صورت ابعاد پارامترهای بالا به شکل زیر تغییر می کنند:

$$\begin{aligned}W_1 &\rightarrow D_{a_1} \times D_x \\b_1 &\rightarrow D_{a_1} \times k \\W_2 &\rightarrow k \times D_{a_1} \\b_2 &\rightarrow k \times 1\end{aligned}$$

همچنین در این حالت ابعاد  $X$  و  $Y$  به شکل زیر خواهند بود:

$$\begin{aligned}X &\rightarrow D_x \times k \\Y &\rightarrow k \times 1\end{aligned}$$

۲.

$$\begin{aligned}\frac{\partial J}{\partial \hat{y}^{(i)}} &= \frac{\partial J}{\partial L^{(i)}} \times \frac{\partial L^{(i)}}{\partial \hat{y}^{(i)}} \\&= \frac{\partial \left( -\frac{1}{m} \sum_{j=1}^m L^{(j)} \right)}{\partial L^{(i)}} \times \frac{\partial \left( y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) (1 - \log \hat{y}^{(i)}) \right)}{\partial \hat{y}^{(i)}} \\&= \left( -\frac{1}{m} \right) \times \left( \frac{y^{(i)}}{\hat{y}^{(i)}} - \frac{1 - y^{(i)}}{1 - \hat{y}^{(i)}} \right) \\&= \frac{\hat{y}^{(i)} - y^{(i)}}{m \hat{y}^{(i)} (1 - \hat{y}^{(i)})} \\&\Rightarrow \frac{\partial J}{\partial \hat{y}} = \frac{\hat{y} - y}{m \hat{y} \odot (1 - \hat{y})}\end{aligned}$$

به طوری که در عبارت بالا  $\odot$  نشان دهنده ضرب elementwise بوده و تقسیم نیز به صورت elementwise انجام گرفته است.

۳.

$$\begin{aligned}\frac{\partial \hat{y}^{(i)}}{\partial z_2} &= \frac{\partial \sigma(z_2)}{\partial z_2} \\&= \sigma(z_2) (1 - \sigma(z_2)) \\&= \hat{y}^{(i)} (1 - \hat{y}^{(i)})\end{aligned}$$

$$\begin{aligned}\frac{\partial z_{\mathfrak{r}}}{\partial a_{\mathfrak{I}}} &= \frac{\partial \left(W_{\mathfrak{r}}a_{\mathfrak{I}} + b_{\mathfrak{r}}\right)}{\partial a_{\mathfrak{I}}} \\ &= W_{\mathfrak{r}}\end{aligned}\tag{.4}$$

$$\begin{aligned}\left(\frac{\partial a_{\mathfrak{I}}}{\partial z_{\mathfrak{I}}}\right)_i &= \left(\frac{\partial \text{ReLU}(z_{\mathfrak{I}})}{\partial z_{\mathfrak{I}}}\right)_i \\ &= \begin{cases} \cdot & z_{\mathfrak{I}_i} \leqslant \cdot \\ \mathfrak{I} & z_{\mathfrak{I}_i} > \cdot \end{cases}\end{aligned}\tag{.5}$$

$$\begin{aligned}\frac{\partial z_{\mathfrak{I}}}{\partial W_{\mathfrak{I}}} &= \frac{\partial \left(W_{\mathfrak{I}}x^{(i)} + b_{\mathfrak{I}}\right)}{\partial W_{\mathfrak{I}}} \\ &= x^{(i)T}\end{aligned}\tag{.6}$$

$$\begin{aligned}\frac{\partial J}{\partial W_{\mathfrak{I}}} &= -\frac{\mathfrak{I}}{m}\sum_{i=\mathfrak{I}}^m\frac{\partial J}{\partial \hat{y}^{(i)}}\cdot\frac{\partial \hat{y}^{(i)}}{\partial z_{\mathfrak{r}}}\cdot\left(\frac{\partial z_{\mathfrak{r}}}{\partial a_{\mathfrak{I}}}\odot\frac{\partial a_{\mathfrak{I}}}{\partial z_{\mathfrak{I}}}\right)\cdot\frac{\partial z_{\mathfrak{I}}}{\partial W_{\mathfrak{I}}} \\ &= -\frac{\mathfrak{I}}{m}\sum_{i=\mathfrak{I}}^m\delta_{\mathfrak{I}}^{(i)}\cdot\delta_{\mathfrak{r}}^{(i)}\cdot\left(\delta_{\mathfrak{r}}^{(i)}\odot\delta_{\mathfrak{r}}^{(i)}\right)\cdot\delta_{\mathfrak{I}}^{(i)}\end{aligned}\tag{.7}$$

## پاسخ مسئله‌ی ۴.

الف) مقادیر خروجی نورون‌ها در ادامه آورده شده‌اند.

لایه مخفی اول:

$$h_1^{(1)} = \text{Sigmoid}(a) = \frac{1}{1 + e^{-a}}$$

$$h_2^{(1)} = \text{Sigmoid}(-a) = \frac{1}{1 + e^a}$$

لایه مخفی دوم:

$$h_1^{(2)} = \text{ReLU}(ah_1^{(1)} - bh_2^{(1)}) = \text{ReLU}\left(\frac{a}{1 + e^{-a}} - \frac{b}{1 + e^a}\right)$$

$$h_2^{(2)} = \text{ReLU}(bh_1^{(1)} - ah_2^{(1)}) = \text{ReLU}\left(\frac{b}{1 + e^{-a}} - \frac{a}{1 + e^a}\right)$$

لایه مخفی سوم:

$$\begin{aligned} h_1^{(3)} &= \text{Softmax}(ah_1^{(2)} - ah_2^{(2)}, bh_1^{(2)} - bh_2^{(2)})_1 = \frac{e^{ah_1^{(2)} - ah_2^{(2)}}}{e^{ah_1^{(2)} - ah_2^{(2)}} + e^{bh_1^{(2)} - bh_2^{(2)}}} \\ \Rightarrow h_1^{(3)} &= \frac{e^{a \text{ReLU}\left(\frac{a}{1+e^{-a}} - \frac{b}{1+e^a}\right) - a \text{ReLU}\left(\frac{b}{1+e^{-a}} - \frac{a}{1+e^a}\right)}}{e^{a \text{ReLU}\left(\frac{a}{1+e^{-a}} - \frac{b}{1+e^a}\right) - a \text{ReLU}\left(\frac{b}{1+e^{-a}} - \frac{a}{1+e^a}\right)} + e^{b \text{ReLU}\left(\frac{b}{1+e^{-a}} - \frac{a}{1+e^a}\right) - b \text{ReLU}\left(\frac{a}{1+e^{-a}} - \frac{b}{1+e^a}\right)}} \\ h_2^{(3)} &= \text{Softmax}(ah_1^{(2)} - ah_2^{(2)}, bh_1^{(2)} - bh_2^{(2)})_2 = \frac{e^{bh_1^{(2)} - bh_2^{(2)}}}{e^{ah_1^{(2)} - ah_2^{(2)}} + e^{bh_1^{(2)} - bh_2^{(2)}}} \\ \Rightarrow h_2^{(3)} &= \frac{e^{b \text{ReLU}\left(\frac{b}{1+e^{-a}} - \frac{a}{1+e^a}\right) - b \text{ReLU}\left(\frac{a}{1+e^{-a}} - \frac{b}{1+e^a}\right)}}{e^{a \text{ReLU}\left(\frac{a}{1+e^{-a}} - \frac{b}{1+e^a}\right) - a \text{ReLU}\left(\frac{b}{1+e^{-a}} - \frac{a}{1+e^a}\right)} + e^{b \text{ReLU}\left(\frac{b}{1+e^{-a}} - \frac{a}{1+e^a}\right) - b \text{ReLU}\left(\frac{a}{1+e^{-a}} - \frac{b}{1+e^a}\right)}} \end{aligned}$$

لایه خروجی:

$$\begin{aligned} a_1^{(3)} &= ah_1^{(3)} + \omega ah_2^{(3)} \\ \Rightarrow a_1^{(3)} &= \frac{ae^{a \text{ReLU}\left(\frac{a}{1+e^{-a}} - \frac{b}{1+e^a}\right) - a \text{ReLU}\left(\frac{b}{1+e^{-a}} - \frac{a}{1+e^a}\right)} + \frac{a}{\gamma} e^{b \text{ReLU}\left(\frac{b}{1+e^{-a}} - \frac{a}{1+e^a}\right) - b \text{ReLU}\left(\frac{a}{1+e^{-a}} - \frac{b}{1+e^a}\right)}}{e^{a \text{ReLU}\left(\frac{a}{1+e^{-a}} - \frac{b}{1+e^a}\right) - a \text{ReLU}\left(\frac{b}{1+e^{-a}} - \frac{a}{1+e^a}\right)} + e^{b \text{ReLU}\left(\frac{b}{1+e^{-a}} - \frac{a}{1+e^a}\right) - b \text{ReLU}\left(\frac{a}{1+e^{-a}} - \frac{b}{1+e^a}\right)}} \\ a_2^{(3)} &= ah_1^{(3)} - bh_2^{(3)} \\ \Rightarrow a_2^{(3)} &= \frac{ae^{b \text{ReLU}\left(\frac{b}{1+e^{-a}} - \frac{a}{1+e^a}\right) - b \text{ReLU}\left(\frac{a}{1+e^{-a}} - \frac{b}{1+e^a}\right)} - be^{a \text{ReLU}\left(\frac{a}{1+e^{-a}} - \frac{b}{1+e^a}\right) - a \text{ReLU}\left(\frac{b}{1+e^{-a}} - \frac{a}{1+e^a}\right)}}{e^{a \text{ReLU}\left(\frac{a}{1+e^{-a}} - \frac{b}{1+e^a}\right) - a \text{ReLU}\left(\frac{b}{1+e^{-a}} - \frac{a}{1+e^a}\right)} + e^{b \text{ReLU}\left(\frac{b}{1+e^{-a}} - \frac{a}{1+e^a}\right) - b \text{ReLU}\left(\frac{a}{1+e^{-a}} - \frac{b}{1+e^a}\right)}} \end{aligned}$$



(ب)

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

با جایگذاری مقادیر  $a_1^{(r)}$  و  $a_2^{(r)}$  (به دست آمده در قسمت الف) در عبارت بالا، مقدار دقیق  $E$  به دست می آید.

ج) گرادیان  $E$  را نسبت به تمامی پارامترها، از لایه آخر به لایه اول به ترتیب محاسبه می کنیم.

لایه خروجی:

$$\frac{\partial E}{\partial a_1^{(r)}} = a_1^{(r)}$$

$$\frac{\partial E}{\partial a_2^{(r)}} = a_2^{(r)} - 1$$

$$\frac{\partial E}{\partial W_{11}^{(r)}} = \frac{\partial E}{\partial a_1^{(r)}} \cdot \frac{\partial a_1^{(r)}}{\partial W_{11}^{(r)}} = a_1^{(r)} h_1^{(r)}$$

$$\frac{\partial E}{\partial W_{21}^{(r)}} = \frac{\partial E}{\partial a_1^{(r)}} \cdot \frac{\partial a_1^{(r)}}{\partial W_{21}^{(r)}} = a_1^{(r)} h_2^{(r)}$$

$$\frac{\partial E}{\partial W_{12}^{(r)}} = \frac{\partial E}{\partial a_2^{(r)}} \cdot \frac{\partial a_2^{(r)}}{\partial W_{12}^{(r)}} = \left( a_2^{(r)} - 1 \right) h_1^{(r)}$$

$$\frac{\partial E}{\partial W_{22}^{(r)}} = \frac{\partial E}{\partial a_2^{(r)}} \cdot \frac{\partial a_2^{(r)}}{\partial W_{22}^{(r)}} = \left( a_2^{(r)} - 1 \right) h_2^{(r)}$$

$$\frac{\partial E}{\partial b_1^{(r)}} = \frac{\partial E}{\partial a_1^{(r)}} \cdot \frac{\partial a_1^{(r)}}{\partial b_1^{(r)}} = a_1^{(r)}$$

$$\frac{\partial E}{\partial b_2^{(r)}} = \frac{\partial E}{\partial a_2^{(r)}} \cdot \frac{\partial a_2^{(r)}}{\partial b_2^{(r)}} = a_2^{(r)} - 1$$

لایه مخفی سوم:

$$\frac{\partial E}{\partial h_1^{(r)}} = \frac{\partial E}{\partial a_1^{(r)}} \cdot \frac{\partial a_1^{(r)}}{\partial h_1^{(r)}} + \frac{\partial E}{\partial a_2^{(r)}} \cdot \frac{\partial a_2^{(r)}}{\partial h_1^{(r)}} = a_1^{(r)} a - \left( a_2^{(r)} - 1 \right) b$$

$$\frac{\partial E}{\partial h_2^{(r)}} = \frac{\partial E}{\partial a_1^{(r)}} \cdot \frac{\partial a_1^{(r)}}{\partial h_2^{(r)}} + \frac{\partial E}{\partial a_2^{(r)}} \cdot \frac{\partial a_2^{(r)}}{\partial h_2^{(r)}} = \frac{1}{2} a_1^{(r)} a + \left( a_2^{(r)} - 1 \right) a$$

$$\begin{aligned} \frac{\partial E}{\partial W_{11}^{(r)}} &= \frac{\partial E}{\partial h_1^{(r)}} \cdot \frac{\partial h_1^{(r)}}{\partial W_{11}^{(r)}} + \frac{\partial E}{\partial h_2^{(r)}} \cdot \frac{\partial h_2^{(r)}}{\partial W_{11}^{(r)}} = \left( a_1^{(r)} a - \left( a_2^{(r)} - 1 \right) b \right) h_1^{(r)} (1 - h_1^{(r)}) h_1^{(r)} \\ &\quad - \left( \frac{1}{2} a_1^{(r)} a + \left( a_2^{(r)} - 1 \right) a \right) h_1^{(r)} h_2^{(r)} h_1^{(r)} \end{aligned}$$



$$\frac{\partial E}{\partial W_{\Psi}^{(\mathbf{r})}} = \frac{\partial E}{\partial h_{\Psi}^{(\mathbf{r})}} \cdot \frac{\partial h_{\Psi}^{(\mathbf{r})}}{\partial W_{\Psi}^{(\mathbf{r})}} = \frac{\partial E}{\partial h_{\Psi}^{(\mathbf{r})}} \cdot \text{ReLU}'(bh_{\Psi}^{(\mathbf{r})} - ah_{\Psi}^{(\mathbf{r})})h_{\Psi}^{(\mathbf{r})}$$

$$\frac{\partial E}{\partial b_{\mathbf{y}}^{(\mathbf{y})}} = \frac{\partial E}{\partial h_{\mathbf{y}}^{(\mathbf{y})}} \cdot \frac{\partial h_{\mathbf{y}}^{(\mathbf{y})}}{\partial b_{\mathbf{y}}^{(\mathbf{y})}} = \frac{\partial E}{\partial h_{\mathbf{y}}^{(\mathbf{y})}} \cdot \text{ReLU}'(bh_{\mathbf{y}}^{(\mathbf{y})} - ah_{\mathbf{y}}^{(\mathbf{y})})$$

$$\begin{aligned} \frac{\partial E}{\partial h_{\mathfrak{l}}^{(1)}} &= \frac{\partial E}{\partial h_{\mathfrak{l}}^{(\mathfrak{r})}} \cdot \frac{\partial h_{\mathfrak{l}}^{(\mathfrak{r})}}{\partial h_{\mathfrak{l}}^{(1)}} + \frac{\partial E}{\partial h_{\mathfrak{r}}^{(\mathfrak{r})}} \cdot \frac{\partial h_{\mathfrak{r}}^{(\mathfrak{r})}}{\partial h_{\mathfrak{l}}^{(1)}} = \frac{\partial E}{\partial h_{\mathfrak{l}}^{(\mathfrak{r})}} \cdot \text{ReLU}'(ah_{\mathfrak{l}}^{(1)} - bh_{\mathfrak{r}}^{(1)})a \\ &\quad + \frac{\partial E}{\partial h_{\mathfrak{r}}^{(\mathfrak{r})}} \cdot \text{ReLU}'(bh_{\mathfrak{l}}^{(1)} - ah_{\mathfrak{r}}^{(1)})b \end{aligned}$$

$$\frac{\partial E}{W_{\lambda}^{(1)}} = \frac{\partial E}{\partial h_{\lambda}^{(1)}} \cdot \frac{\partial h_{\lambda}^{(1)}}{\partial W_{\lambda}^{(1)}} = \frac{\partial E}{\partial h_{\lambda}^{(1)}} \cdot h_{\lambda}^{(1)}(1 - h_{\lambda}^{(1)})$$

$$\frac{\partial E}{b_{\lambda}^{(\lambda)}} = \frac{\partial E}{\partial h_{\lambda}^{(\lambda)}} \cdot \frac{\partial h_{\lambda}^{(\lambda)}}{\partial b_{\lambda}^{(\lambda)}} = \frac{\partial E}{\partial h_{\lambda}^{(\lambda)}} \cdot h_{\lambda}^{(\lambda)}(1 - h_{\lambda}^{(\lambda)})$$

$$\frac{\partial E}{W_{\gamma}^{(1)}} = \frac{\partial E}{\partial h_{\gamma}^{(1)}} \cdot \frac{\partial h_{\gamma}^{(1)}}{\partial W_{\gamma}^{(1)}} = ,$$

$$\frac{\partial E}{W_{\downarrow\downarrow}^{(\downarrow)}} = \frac{\partial E}{\partial h_{\downarrow}^{(\downarrow)}} \cdot \frac{\partial h_{\downarrow}^{(\downarrow)}}{\partial W_{\downarrow\downarrow}^{(\downarrow)}} = \frac{\partial E}{\partial h_{\downarrow}^{(\downarrow)}} \cdot h_{\downarrow}^{(\downarrow)}(1 - h_{\downarrow}^{(\downarrow)})$$

$$\frac{\partial E}{W_{\Upsilon\Upsilon}^{(1)}} = \frac{\partial E}{\partial h_{\Upsilon}^{(1)}} \cdot \frac{\partial h_{\Upsilon}^{(1)}}{\partial W_{\Upsilon\Upsilon}^{(1)}} = ,$$

$$\frac{\partial E}{b_{\Psi}^{(\Psi)}} = \frac{\partial E}{\partial h_{\Psi}^{(\Psi)}} \cdot \frac{\partial h_{\Psi}^{(\Psi)}}{\partial b_{\Psi}^{(\Psi)}} = \frac{\partial E}{\partial h_{\Psi}^{(\Psi)}} \cdot h_{\Psi}^{(\Psi)}(1 - h_{\Psi}^{(\Psi)})$$

$$\text{ReLU}'(x) = \begin{cases} 1 & x \geq 0 \\ 0 & x < 0 \end{cases}$$

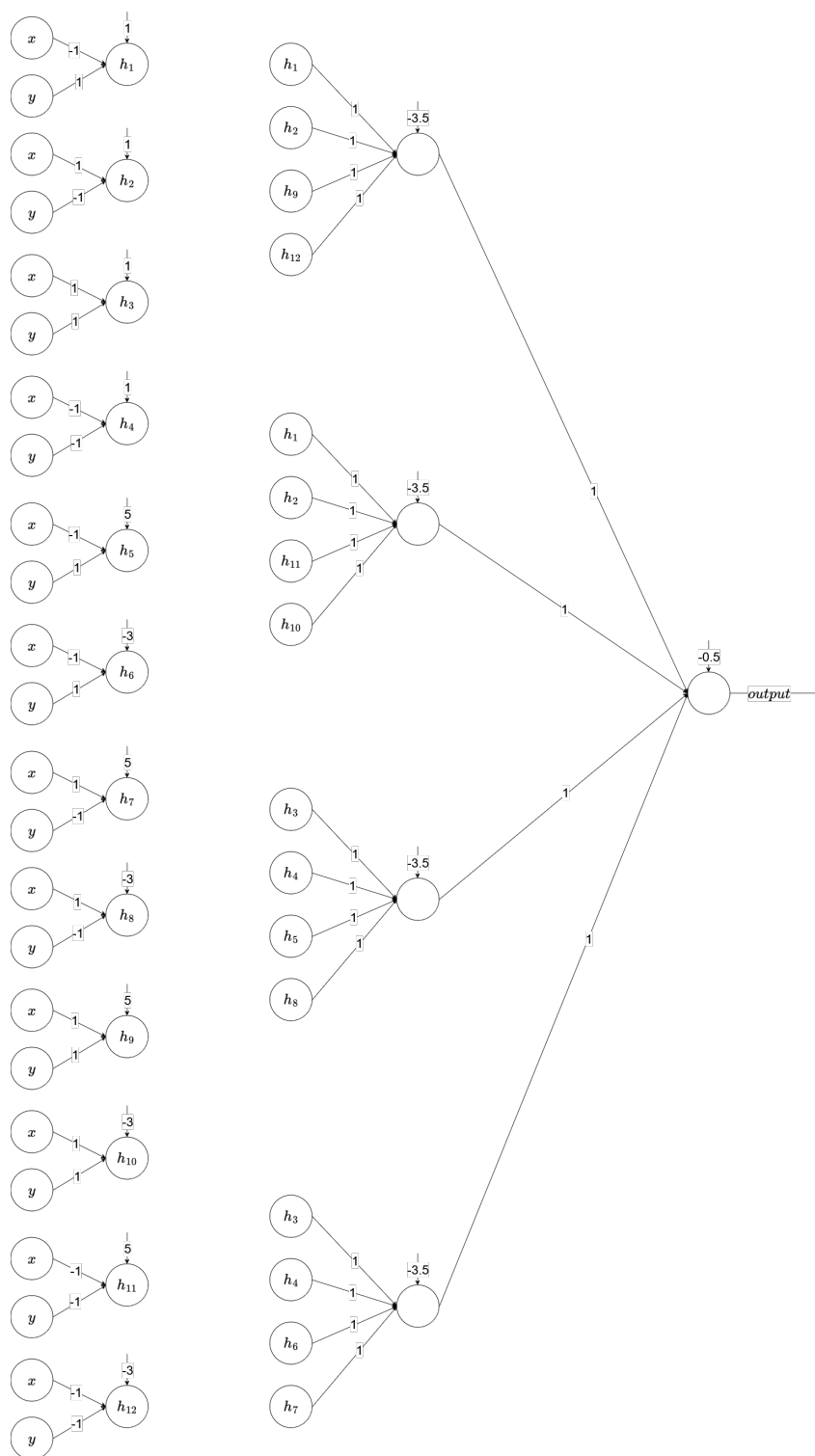
د) برای هر کدام از  $W_{ij}^{(l)}$  ها و نیز هر کدام از  $b_i^{(l)}$  ها ( $i, j \in \{1, 2\} \wedge l \in \{1, 2, 3, 4\}$ )، مقدار پارامتر را به شکل زیر بروزرسانی می کنیم:

$$param_{\text{new}} = param - \eta \frac{\partial E}{\partial param}$$

با استفاده از گرادیان های به دست آمده در قسمت قبل و همچنین مقادیر اولیه پارامتر ها (داده شده در صورت سوال)، می توان پارامتر های موجود در شبکه عصبی را بروزرسانی نمود.

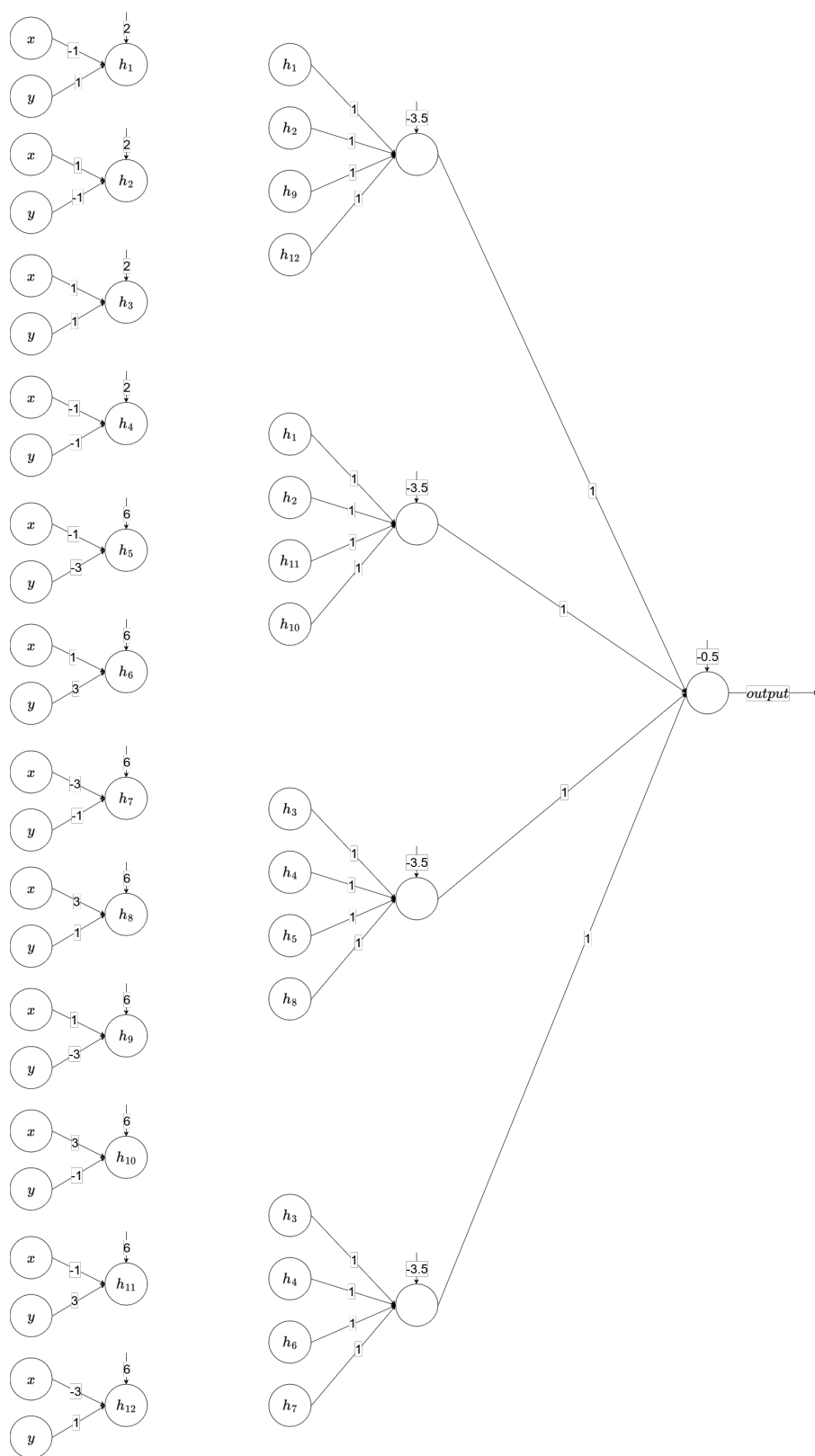
## پاسخ مسئله‌ی ۵.

الف) شکل زیر یک شبکه مطلوب با ۱۷ عدد TLU را نمایش می‌دهد. برای جلوگیری از پیچیدگی زیاد شکل، ورودی‌ها و نیز برخی از نورون‌ها چند بار رسم شده‌اند.



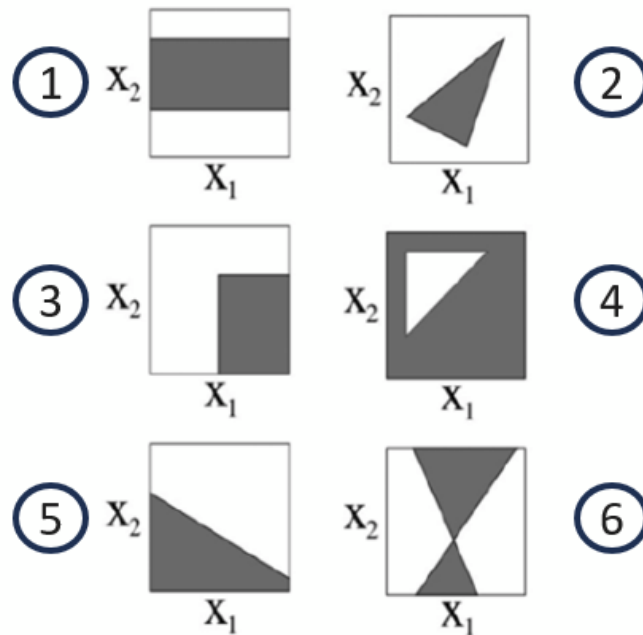
شکل ۵: معماری شبکه مطلوب

ب) بله. شبکه عصبی مطلوب را در شکل زیر مشاهده می کنید.



شکل ۶: معماری شبکه مطلوب

ج) شکل ها را به صورت زیر شماره گذاری می کنیم:



شکل ۷: شماره گذاری اشکال

شبکه  $A$  با شکل شماره ۵ مطابقت دارد، زیرا در این شبکه تنها می توان یک خط رسم کرد که صفحه را به وسیله آن به دو بخش مجزا تقسیم می کنیم.

شبکه  $B$  منطبق بر شکل شماره ۱ است؛ دلیل آن این است که در این شکل دو خط وجود دارد. در این شبکه نیز دو خط ایجاد می کنیم و سپس در لایه دوم، خروجی ها را با استفاده از عملگر منطقی AND ترکیب می کنیم تا ناحیه مشترک تعیین شود.

شبکه  $C$  با شکل شماره ۶ مرتبط است. در این حالت، ابتدا در لایه اول دو خط ترسیم می کنیم، اما برای ترکیب خروجی ها نیاز به عملگر XOR داریم، که این امر مستلزم افزودن دو لایه دیگر به شبکه است.

شبکه  $D$  به شکل شماره ۳ اختصاص دارد؛ زیرا می توانیم در لایه اول یک خط افقی و یک خط عمودی رسم کنیم. سپس، در لایه دوم، با بهره گیری از عملگر منطقی AND، ناحیه بین این دو خط را مشخص می کنیم.

شبکه  $E$  متناظر با شکل شماره ۲ است. در این شبکه، نیاز است که در لایه اول سه خط ترسیم کنیم و سپس در لایه بعدی، خروجی ها را با هم توسط عملگر AND ترکیب کنیم تا ناحیه مورد نظر حاصل شود.

شبکه  $F$  مرتبط با شکل شماره ۴ است؛ زیرا این شکل از یک خط افقی، یک خط عمودی و یک خط مورب تشکیل شده است. در این شبکه، یکی از نوروها فقط ورودی  $x_1$  را دریافت می کند، دیگری تنها ورودی  $x_2$ ، و دیگری هر دو ورودی را دریافت می کند. سپس در لایه بعدی، خروجی ها را با عملگر منطقی OR ترکیب می کنیم تا ناحیه مطلوب تعیین شود.

## پاسخ مسئله‌ی ۶.

الف) در ابتدا دقت کنید که لایه  $l$  ام شبکه عصبی را با  $\mathcal{L}_l$  نمایش می‌دهیم، به طوری که  $\mathcal{L}$  لایه ورودی و  $\mathcal{L}_{L+1}$  لایه خروجی بوده و برای هر  $l \in \{1, \dots, L\}$  همان  $l$  امین لایه مخفی شبکه عصبی می‌باشد. برای  $i$  امین نورون مخفی در  $\mathcal{L}_l$ ،  $\sigma(z_i^{(l)})$  به شکل زیر است:

$$\sigma(z_i^{(l)}) = \begin{cases} z_i^{(l)} & z_i^{(l)} \geq 0 \\ 0 & z_i^{(l)} < 0 \end{cases}$$

و در نتیجه  $\sigma(z_i^{(l)})$  از دو بخش خطی تشکیل شده است. حال توجه نمایید که برای یک شبکه عصبی داده شده، یک نمونه  $x \in \mathcal{X}$  مقدار  $z_i^{(l)}$  و در نتیجه مقدار  $\sigma(z_i^{(l)})$  را به طور یکتا تعیین می‌کند. به عبارت دیگر همانطور که در صورت سوال هم اشاره شده است، می‌توان وضعیت نورون مخفی  $i$  ام در  $\mathcal{L}_l$  را با یک متغیر دودویی مانند  $c_i^{(l)} \in \{0, 1\}$  تعیین کرد، به طوری که  $z_i^{(l)} \geq 0$  اگر و تنها اگر  $c_i^{(l)} = 1$  و  $z_i^{(l)} < 0$  اگر و تنها اگر  $c_i^{(l)} = 0$ .

فرض کنید  $c^{(l)} = [c_1^{(l)}, \dots, c_{n_l}^{(l)}]$  ماتریس وضعیت تمامی نورون‌های مخفی در  $\mathcal{L}_l$  بوده و  $C = [c^{(1)}, \dots, c^{(L)}]$  نشان دهنده وضعیت تمامی نورون‌های مخفی در شبکه عصبی باشد. برای یک شبکه عصبی داده شده، ماتریس  $C$  با داشتن نمونه  $x$  به طور یکتا تعیین می‌شود.

در ادامه، تابعی که نمونه  $x$  را به ماتریس  $C \in \{0, 1\}^N$  (تعداد کل نورون‌های مخفی است) نگاشت می‌کند، با  $\mathcal{C} : \mathcal{X} \rightarrow \{0, 1\}^N$  نمایش می‌دهیم. بنابراین خواهیم داشت:

$$a^{(l)} = \sigma(z^{(l)}) = c^{(l)} \odot z^{(l)}$$

به طوری که در عبارت بالا  $\odot$  نشان دهنده ضرب elementwise می‌باشد. در نتیجه می‌توان نوشت:

$$z^{(l+1)} = W^{(l)} a^{(l)} + b^{(l)} = W^{(l)} (c^{(l)} \odot z^{(l)}) + b^{(l)} = \tilde{W}^{(l)} z^{(l)} + b^{(l)}$$

به طوری که در عبارت بالا،  $\tilde{W}^{(l)} = W^{(l)} \odot c^{(l)}$  می‌باشد. با جایگذاری مکرر تساوی اخیر برای لایه‌های قبل تر، می‌توان  $z^{(l+1)}$  را برای  $l \in \{1, \dots, L\}$  به شکل زیر بازنویسی کرد:

$$z^{(l+1)} = \left( \prod_{k=1}^l \tilde{W}^{(k)} \right) z^{(1)} + \sum_{h=1}^l \left( \prod_{k=h+1}^l \tilde{W}^{(k)} \right) b^{(h)}$$

با جایگذاری  $z^{(1)} = W^{(0)} x + b^{(1)}$  در تساوی اخیر به دست می‌آید:

$$\begin{aligned} z^{(l+1)} &= \left( \prod_{k=1}^l \tilde{W}^{(k)} \right) (W^{(0)} x + b^{(1)}) + \sum_{h=1}^l \left( \prod_{k=h+1}^l \tilde{W}^{(k)} \right) b^{(h)} \\ &= \left( \prod_{k=1}^l \tilde{W}^{(k)} \right) x + b^{(1)} \left( \prod_{k=1}^l \tilde{W}^{(k)} \right) + \sum_{h=1}^l \left( \prod_{k=h+1}^l \tilde{W}^{(k)} \right) b^{(h)} \\ &= \tilde{W}^{(\cdot:l)} x + b^{(\cdot:l)} \end{aligned}$$



به طوری که در تساوی صفحه قبل،  $\tilde{W}^{(\cdot:l)}$  ماتریس ضریب  $x$  و  $b^{(\cdot:l)}$  جمع بایاس های باقی مانده می باشد. با این تعریف، به وضوح داریم:

$$F_A(x, W) = \text{softmax}(z^{(L+1)}) = \text{softmax}(\tilde{W}^{(\cdot:L)}x + b^{(\cdot:L)})$$

برای یک شبکه عصبی و یک نمونه  $x$  مشخص و ثابت،  $\tilde{W}^{(\cdot:l)}$  و  $b^{(\cdot:l)}$  پارامترهای ثابتی هستند که با توجه به ماتریس وضعیت نورون های مخفی یا همان  $C = \mathcal{C}(x)$  به صورت یکتا و منحصر به فرد تعیین می شوند. به علاوه، ماتریس  $C$  منجر به ایجاد یک خانواده از نابرابری ها به شکل زیر می شود: (با توجه به اولین تساوی این بخش)

$$(\mathfrak{z}_C^{(l+1)} - 1) \odot (\tilde{W}^{(\cdot:l)}x + b^{(\cdot:l)}) \geq 0, \quad l \in \{0, \dots, L-1\}$$

که این نمایش دهنده یک convex polytope می باشد و در نتیجه حکم مسئله برقرار است.

برای مشاهده منبع حل این سوال، روی این [لینک](#) کلیک کنید.

ب) در ابتدا نشان می دهیم که  $r(k, m) \leq \mathcal{O}(k^m)$  می باشد. برای این منظور در ابتدا دقت کنید که می توان نوشت:

$$\binom{k}{i} = \frac{k!}{(k-i)!i!} = \frac{k(k-1)\dots(k-i+1)}{i!} \leq \frac{k^i}{i!} \quad (1)$$

از طرف دیگر می توان نوشت:

$$e^i = \sum_{j=0}^{\infty} \frac{i^j}{j!} \geq \frac{i^i}{i!} \implies \frac{1}{i!} \leq \left(\frac{e}{i}\right)^i \quad (2)$$

از نابرابری های (1) و (2) نتیجه می شود:

$$\binom{k}{i} \leq \left(\frac{ek}{i}\right)^i$$

در نتیجه خواهیم داشت:

$$\begin{aligned} r(k, m) &\leq \sum_{i=0}^m \binom{k}{i} \leq \sum_{i=0}^m \left(\frac{ek}{i}\right)^i \\ &\leq 1 + m + \frac{m(m-1)}{2} + \sum_{i=3}^m k^i \\ &= \mathcal{O}(m^2) + \frac{k^{m+1} - k^3}{k-1} = \mathcal{O}(k^m) \end{aligned}$$

دقت کنید که نابرابری های موجود در مراحل بالا برای  $k > 1$  نوشته شده اند، چرا که برای  $k = 1$  به وضوح  $r(k, m) = \mathcal{O}(1)$  می باشد.

حال توجه نمایید که هر یک از نورون های موجود در شبکه عصبی نمایش دهنده یک خط می باشد؛ چرا که خروجی آن به شکل  $\sigma \left( W_i^{(l)} h_i^{(l)} + b_i^{(l)} \right)$  است که مشخص می کند  $h_i^{(l)}$  در بالا یا پایین خطی که با بردار نرمال  $W_i^{(l)}$  و بایاس  $b_i^{(l)}$  تعیین می شود، قرار می گیرد. بنابراین با عبور از هر لایه در یک شبکه عصبی،  $k$  خط جدید در فضای نمونه ها ایجاد می شوند و تعداد نواحی ایجاد شده در این فضا را  $r(k, m) \leq O(k^m)$  برابر می کنند. از آنجا که در ابتدا کل فضا یک ناحیه واحد است، پس از عبور از تمامی  $n$  لایه موجود در شبکه عصبی تعداد کل نواحی موجود در فضای نمونه ها حداکثر  $O((k^m)^n) = O(k^{mn})$  خواهد بود. از طرف دیگر، تناظری یک به یک میان activation pattern ها و تعداد نواحی ایجاد شده در فضای نمونه ها وجود دارد و می توان نتیجه گرفت:

$$\mathcal{A} \left( F_{A(n,k)}(\mathbb{R}^m; W) \right) \leq O(k^{mn})$$

که این همان حکم مسئله است.

---

در حل این تمرین با غسل مسکین همفکری صورت گرفته است.