

به نام خدا

تمرین سری سوم
درس یادگیری ماشین
دکتر علی شریفی زارچی

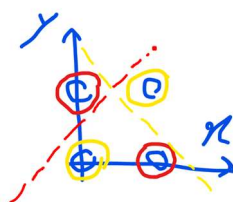
فرزان رحمانی
۴۰۳۲۱۰۷۲۵

سوال اول

الف) از ۱۶ تابع بولین مجزا در دو ورودی، ۱۴ تابع را می توان با یک **perceptron** نشان داد.

استدلال: یک پرسپترون می تواند هر تابع بولین خطی قابل تفکیک (linearly separable boolean function) را نشان دهد. در مورد دو ورودی، چهار ترکیب ورودی ممکن وجود دارد (۰،۰)، (۰،۱)، (۱،۰)، (۱،۱)، که منجر به $2^4 = 16$ تابع بولین ممکن می شود. با این حال، تنها توابعی که می توانند با یک خط مستقیم در فضای ورودی از هم جدا شوند، به صورت خطی قابل تفکیک هستند.

دو تابعی که به صورت خطی قابل تفکیک نیستند، توابع XOR (exclusive OR) و XNOR (exclusive NOR) هستند. این توابع را نمی توان با پرسپترون نشان داد زیرا هیچ مرز خطی وجود ندارد که بتواند کلاس های خروجی آنها را از هم جدا کند. بنابراین، یک پرسپترون می تواند بقیه توابع یعنی ۱۴ تابع از ۱۶ تابع بولین را برای دو ورودی نشان دهد.



توضیحات بیشتر:

- ما می توانیم فضای ورودی را به صورت هندسی به صورت ۴ نقطه تجسم کنیم: (۰،۰)، (۰،۱)، (۱،۰)، و (۱،۱).
- یک تابع تفکیک پذیر خطی به این معنی است که می توانیم یک خط مستقیم رسم کنیم که خروجی های درست را از خروجی های نادرست جدا می کند.

بنابراین توابع که به صورت خطی قابل تفکیک هستند عبارتند از:

- توابع Single input (۴ تابع - بستگی به x_1 یا x_2 یا نقیض آنها دارند)
- توابع ثابت یا Constant (۲ تابع - همیشه درست یا همیشه نادرست)
- توابعی که شامل ۳ خروجی درست و ۱ خروجی نادرست یا برعکس ۳ خروجی نادرست و ۱ خروجی درست هستند. (۸ تابع - شامل AND, OR, NAND, NOR, ...)
- $۸ + ۲ + ۴ = ۱۴$

بنابراین، از ۱۶ تابع بولی ممکن، تنها ۱۴ تابع را می توان با یک پرسپترون نشان داد. ۲ تابع باقیمانده، یعنی XOR و XNOR، به صورت خطی قابل تفکیک نیستند و نمی توان آنها را با یک پرسپترون نشان داد.

ب) غلط. هر دو tanh و sigmoid می توانند از مشکلات vanishing gradient رنج ببرند. در حالی که tanh بر خلاف سیگموئید صفر محور است، همچنان دارای گرادیان هایی است که برای ورودی های مثبت یا منفی بزرگ بسیار کوچک می شوند.

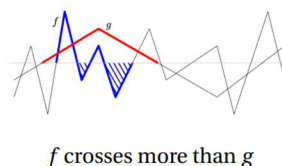
ج) صحیح. هر تابع منطقی را می توان با یک شبکه عصبی دو لایه با یک لایه پنهان با اندازه کافی نشان داد. توابع فعال سازی ReLU می توانند مرزهای تصمیم گیری غیرخطی را بازنمایی کنند و میتوانند هر تابع منطقی را تقریب بزنند.
 Universal Approximation Theorem: Any continuous function can be approximated by a two-layer neural network with appropriate activation function.

د) غلط. افزایش عمق و عرض می تواند منجر به overfitting، افزایش زمان آموزش، و vanishing gradients (برای افزایش عمق) شود و همیشه عملکرد را بهبود نمی بخشد. همچنین ممکن است آموزش را دشوارتر کند. بنابراین، همیشه عملکرد بهتر را تضمین نمی کند.

ه) غلط. در حالی که یک شبکه کم عمق و عریض (hallow and wide network) از نظر تئوری می تواند هر تابعی را تقریب بزند (the Universal Approximation Theorem)، شبکه های عمیق به دلیل ویژگی ها و ساختار سلسله مراتبی خود اغلب در نمایش توابع پیچیده با نوروں های کمتر کارآمدتر هستند.

Depth Separation Theorem

- This theorem can be proved by showing that:
 - Functions with few oscillations poorly approximate functions with many oscillations.
 - Functions computed by networks with few layers must have few oscillations.
 - Functions computed by networks with many layers can have many oscillations.
- The key idea is that just a few function compositions (layers) suffice to construct a highly oscillatory function, whereas function addition (adding nodes but keeping depth fixed) gives a function with few oscillations. Thereafter, an elementary counting argument suffices to show that low-oscillation functions can not approximate high-oscillation functions.



و) مشکل ناپدید شدن گرادیان: مشکل ناپدید شدن گرادیان در طول آموزش شبکه های عصبی عمیق زمانی رخ می دهد که گرادیان های مورد استفاده برای به روزرسانی وزن ها بسیار کوچک می شوند، زیرا به لایه های قبلی منتشر می شوند. این امر به روزرسانی وزن ها را در این لایه ها کاهش می دهد و به طور موثر شبکه را از یادگیری وابستگی های دوربرد جلوگیری می کند. به عبارت دیگر، این اتفاق می افتد زیرا:

- توابع فعال سازی سنتی مانند سیگموئید و \tanh دارای گرادیان هایی هستند که برای ورودی های مثبت یا منفی بزرگ بسیار کوچک هستند.
- هنگامی که این شیب های کوچک در حین انتشار پس از انتشار چند بار ضرب می شوند، حتی کوچکتر می شوند.
- که منجر به یادگیری بسیار کند در لایه های اولیه می شود.

چگونه ReLU کمک می کند؟ توابع فعال سازی مانند Rectified Linear Unit (ReLU) مشکل گرادیان ناپدید شدن را کاهش می دهند زیرا مشتق آن به جای کاهش یافتن، ثابت است (۱ برای ورودی های مثبت و ۰ برای ورودی های منفی). برخلاف توابع سیگموئید یا \tanh ، که مقادیر ورودی را در محدوده کوچکی قرار می دهند و مشتقات کمتر از ۱ دارند، ReLU به گرادیان ها اجازه می دهد بدون کوچک کردن آنها عبور کنند و تضمین می کند که لایه های نزدیک تر به ورودی در طول آموزش به روزرسانی های قابل توجهی دریافت می کنند.

به عبارت دیگر، ReLU به حل این مشکل به شیوه زیر کمک می کند:

- داشتن گرادیان ثابت ۱ برای همه ورودی های مثبت
- محاسبه مشتق ساده (۱ برای $x > 0$ ، ۰ برای $x < 0$)
- عدم اشباع (saturation) در ناحیه مثبت
- فعال سازی Sparse (مقادیر منفی را صفر می کند)

ز) غلط. SGD با اندازه دسته‌ای کوچک معمولاً دارای واریانس بالاتری در به‌روزرسانی‌ها است و در حالی که به‌روزرسانی‌های مکرر را انجام می‌دهد، همگرایی سریع‌تری را نسبت به mini-batch gradient descent تضمین نمی‌کند. به عبارت دیگر، در حالی که اندازه‌های دسته‌ای کوچکتر در SGD منجر به به‌روزرسانی‌های مکرر می‌شوند، واریانس بالاتری را در برآورد گرادیان ایجاد می‌کنند که می‌تواند منجر به نوسانات و همگرایی ناپایدار شود.

ی) از فرمول‌های زیر استفاده می‌کنیم:

First Moment (Momentum)

- **Definition:** The first moment, m_t , represents a moving average of past gradients. It builds "velocity" that propels learning in a consistent direction.

- **Update Rule:**

$$m_{t+1} = \beta_1 m_t + (1 - \beta_1) \nabla_w J(w_t)$$

$$w_{t+1} = w_t - \eta m_{t+1}$$

where:

- β_1 : Decay rate, usually 0.9 or 0.99, which controls the weight of past gradients.
- η : Learning rate.
- **Why Use First Momentum?**
 - Inspired by the idea of rolling momentum, it smooths and accelerates learning by sustaining direction from prior gradients.
 - This type of momentum is ideal for traversing narrow valleys or regions where standard gradient descent would oscillate.

داریم:

$$(y = 0.3x^4 - 0.1x^3 - 2x^2 - 0.8x) \quad -$$

$$\text{Initial point: } (x_0 = -2.8) \quad -$$

$$\text{Learning rate: } (\gamma = 0.05) \quad -$$

$$\text{Momentum coefficient: } (\mu = 0.7) \quad -$$

$$\text{Number of iterations: } (i = 2) \quad -$$

$$\text{Initial momentum: } (m_0 = 0) \quad -$$

ابتدا مشتق y را نسبت به x پیدا می‌کنیم.

$$\frac{dy}{dx} = y' = 1.2x^3 - 0.3x^2 - 4x - 0.8$$

محاسبات در هر مرحله به شکل زیر است:

$$g_i = 1.2x_i^3 - 0.3x_i^2 - 4x_i - 0.8$$

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

$$x_i = x_{i-1} - \gamma m_i$$

$$y_i = 0.3x_i^4 - 0.1x_i^3 - 2x_i^2 - 0.8x_i$$

توجه: اصلاح بایاس در هر مرحله انجام نشده است.

تکرار ۰ (وضعیت ابتدایی):

$$m_0 = 0, \quad x_0 = -2.8, \quad y_0 = 7.1949$$

تکرار ۱:

1. محاسبه مشتق (g_1) at ($x_0 = -2.8$):

$$g_1 = 1.2(-2.8)^3 - 0.3(-2.8)^2 - 4(-2.8) - 0.8 = -18.2944$$

2. به روز رسانی ممان (m_1):

$$m_1 = \mu m_0 + (1 - \mu)g_1 = 0.7 \times 0 + 0.3 \times (-18.2944) = -5.4883$$

3. به روز رسانی (x_1):

$$x_1 = x_0 - \gamma m_1 = -2.8 - 0.05 \times (-5.4883) = -2.5256$$

4. محاسبه (y_1):

$$y_1 = 0.3(-2.5256)^4 - 0.1(-2.5256)^3 - 2(-2.5256)^2 - 0.8(-2.5256) \approx 3.0801$$

تکرار ۲:

1. محاسبه مشتق (g_2) at ($x_1 = -2.5256$):

$$g_2 = 1.2(-2.5256)^3 - 0.3(-2.5256)^2 - 4(-2.5256) - 0.8 \approx -11.9428$$

2. به روز رسانی ممان (m_2):

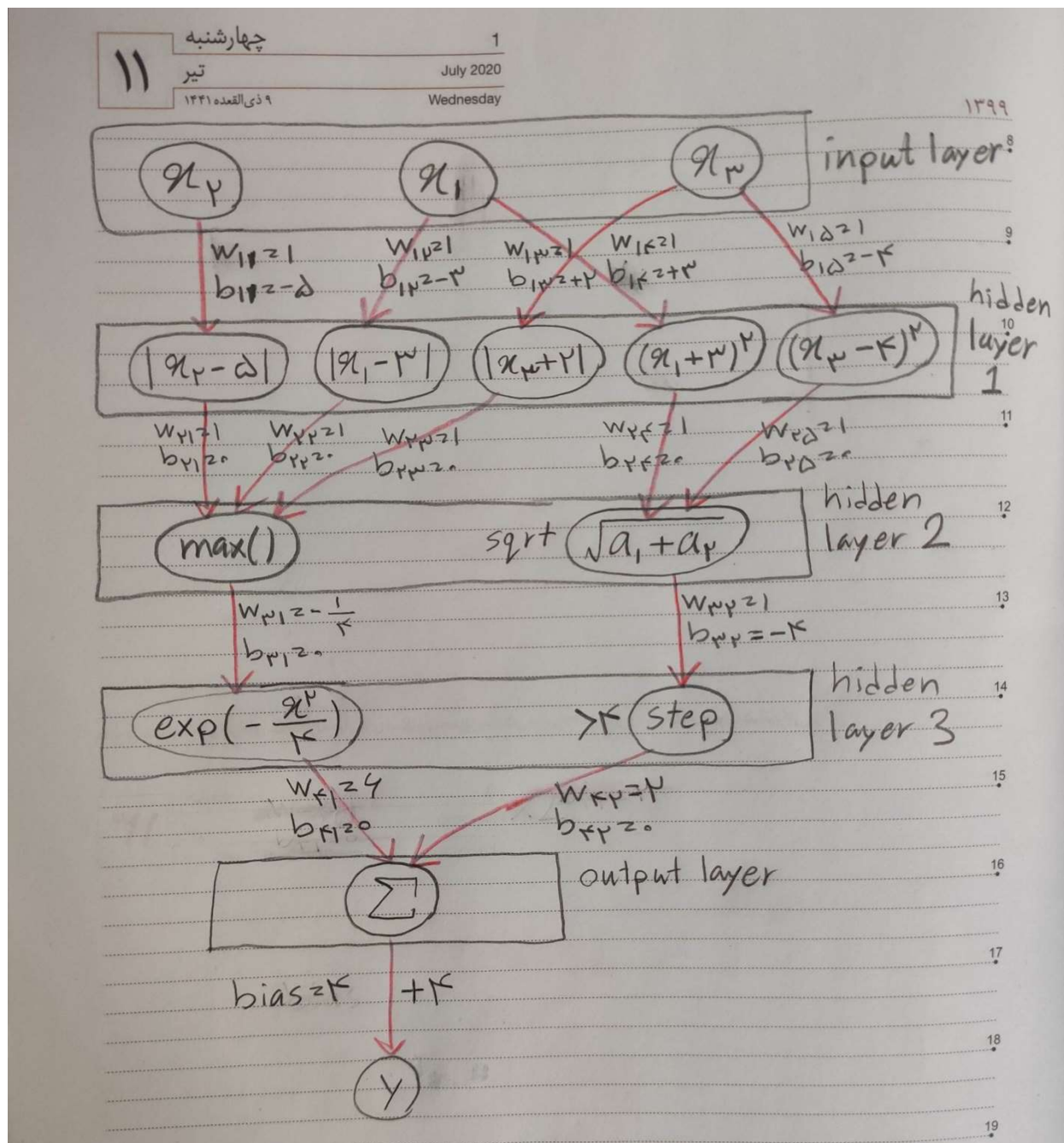
$$m_2 = \mu m_1 + (1 - \mu)g_2 = 0.7 \times (-5.4883) + 0.3 \times (-11.9428) = -7.4247$$

3. به روز رسانی (x_2):

$$x_2 = x_1 - \gamma m_2 = -2.5256 - 0.05 \times (-7.4247) = -2.1544$$

4. محاسبه (y_2):

$$y_2 = 0.3(-2.1544)^4 - 0.1(-2.1544)^3 - 2(-2.1544)^2 - 0.8(-2.1544) \approx -0.0968$$



همچنین جزئیات بیشتر هر نورون در ادامه آمده است:

چون بیشتر اسم ها به انگلیسی است ادامه سوال را به انگلیسی نگارش میکنیم تا مفهوم تر باشد.

input Layer:

شامل ورودی ها (x_1, x_2, x_3) است.

1st hidden layer:

- Absolute value neurons (abs1, abs2, abs3: قدر مطلق ها):
 - $f_{\text{net}}: w * x + b$
 - $f_{\text{act}}: |x|$
 - $f_{\text{out}}: \text{identity}$
 - Parameters: $w_1=1, b_1=-3$ for abs1; $w_1=1, b_1=-5$ for abs2; $w_1=1, b_1=2$ for abs3
- Square neurons (توان دوها):
 - $f_{\text{net}}: w * x + b$
 - $f_{\text{act}}: x^2$
 - $f_{\text{out}}: \text{identity}$
 - Parameters: $w_1=1, b_1=3, w_2=1, b_2=-4$

2nd hidden Layer:

- Maximum neuron (ماکزیمم میگیرد):
 - $f_{\text{net}}: \text{takes three inputs } (w_1*a_1+b_1, w_2*a_2 + b_2, w_3*a_3 + b_3)$
 - $f_{\text{act}}: \text{max function } (\max(z_1, z_2, z_3) + \text{bias})$
 - $f_{\text{out}}: \text{identity}$
 - Parameters: $w_1=w_2=w_3=1, b_1=b_2=b_3=0, \text{bias}=0$
- Square root neuron (جذر میگیرد):
 - $f_{\text{net}}: \text{takes two inputs } (w_1*a_1 + w_2*a_2 + \text{bias})$
 - $f_{\text{act}}: \text{sqrt function}$
 - $f_{\text{out}}: \text{identity}$
 - Parameters: $w_1=w_2=1, \text{bias}=0$

3rd hidden Layer:

- Exponential neuron (توان نمایی):
 - $f_{\text{net}}: -x^2/4$
 - $f_{\text{act}}: \exp(x)$
 - $f_{\text{out}}: \text{identity}$
 - Parameters: $w=-1/4, \text{bias}=0$
- Step function neuron:

- $f_{\text{net}}: w \cdot x + \text{bias}$
- f_{act} : step function
- f_{out} : identity
- Parameters: $w=1$, $\text{bias}=-4$

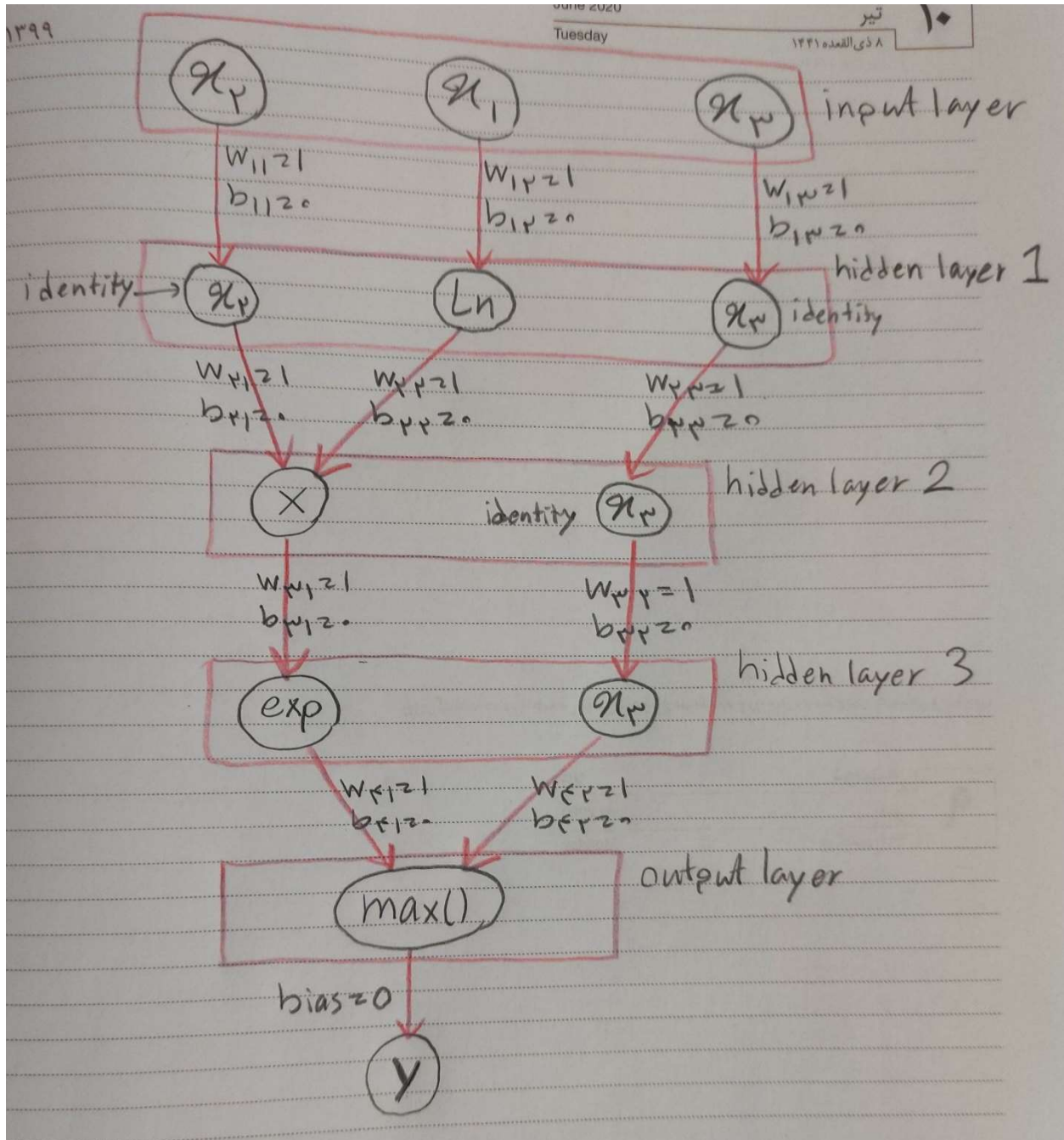
Output Layer:

- Sum neuron (جمع میکنیم):
 - $f_{\text{net}}: w_1x_1 + w_2x_2 + b$
 - f_{act} : identity
 - f_{out} : identity
 - Parameters: $w_1=6$, $w_2=2$, $\text{bias}=4$

ب) نکته: برای محاسبه $x_1^{x_2}$ از خاصیت زیر استفاده میکنیم:

$$x_1^{x_2} = e^{\ln(x_1^{x_2})} = e^{x_2 \times \ln(x_1)}$$

پس ابتدا لگاریتم میگیریم، سپس ضرب میکنیم و در نهایت از exp استفاده میکنیم.



همچنین جزییات بیشتر هر نورون در ادامه آمده است:

چون بیشتر اسم ها به انگلیسی است ادامه سوال را به انگلیسی نگارش میکنیم تا مفهوم تر باشد.

input Layer:

شامل ورودی ها (x_1, x_2, x_3) است.

1st hidden layer:

- Natural Log neuron (x_1 برای \ln):
 - $f_{\text{net}}: w \cdot x + b$
 - $f_{\text{act}}: \ln(x)$
 - $f_{\text{out}}: \text{identity}$
 - Parameters: $w=1, b=0$
- Identity neurons (x_2, x_3 برای تابع همانی):
 - $f_{\text{net}}: w \cdot x + b$
 - $f_{\text{act}}: x$ (identity)
 - $f_{\text{out}}: \text{identity}$
 - Parameters: $w=1, b=0$

2nd hidden layer:

- Multiplication neuron (ضرب):
 - $f_{\text{net}}: a_1=w_1x_1+b_1, a_2=w_2x_2+b_2$
 - $f_{\text{act}}: \text{multiplication}(a_1 * a_2)$
 - $f_{\text{out}}: \text{identity}$
 - Parameters: $w_1=1, w_2=1, b_1=0, b_2=0$
- Identity neurons (x_3 برای تابع همانی):
 - $f_{\text{net}}: w \cdot x + b$
 - $f_{\text{act}}: x$ (تابع همانی)
 - $f_{\text{out}}: \text{identity}$
 - Parameters: $w=1, b=0$

3rd hidden layer

- Exponential neuron (exp):
 - $f_{\text{net}}: w \cdot x + b$
 - $f_{\text{act}}: \exp(x)$
 - $f_{\text{out}}: \text{identity}$

- Parameters: $w=1, b=0$
- Identity neurons (تابع همانی برای x_3):
 - $f_{\text{net}}: w \cdot x + b$
 - $f_{\text{act}}: x$ (تابع همانی)
 - $f_{\text{out}}: \text{identity}$
 - Parameters: $w=1, b=0$

Output layer:

- Maximum neuron (max):
 - f_{net} : takes two inputs ($w_1 \cdot a_1 + b_1, w_2 \cdot a_2 + b_2$)
 - f_{act} : max function ($\max(a_1, a_2) + \text{bias}$)
 - $f_{\text{out}}: \text{identity}$
 - Parameters: $w_1=w_2=1, b_1=b_2=0, \text{bias}=0$

$$\delta_1^{(i)} = \frac{-1}{m} \left[\frac{y^{(i)} - \hat{y}^{(i)}}{\hat{y}^{(i)}(1 - \hat{y}^{(i)})} \right] \xrightarrow{\text{vectorized form}} \frac{\partial J}{\partial \hat{Y}} = \frac{-1}{m} \left[\frac{Y - \hat{Y}}{\hat{Y}(1 - \hat{Y})} \right]$$

$$\frac{\partial \hat{y}^{(i)}}{\partial z_r} = \frac{\partial}{\partial z_r} (\sigma(z_r)) = \sigma(z_r)(1 - \sigma(z_r)) \quad (3)$$

$$\delta_r^{(i)} = \sigma(z_r)(1 - \sigma(z_r))$$

$$\frac{\partial z_r}{\partial a_1} = \frac{\partial}{\partial a_1} (W_r a_1 + b_r) = W_r \quad (4)$$

$$\delta_r^{(i)} = W_r$$

$$\frac{\partial a_1}{\partial z_1} = \frac{\partial}{\partial z_1} (\text{ReLU}(z_1)) = I(z_1 > 0) = \begin{cases} 1 & \text{if } z_1 > 0 \\ 0 & \text{if } z_1 \leq 0 \end{cases} \quad (5)$$

$$\delta_r^{(i)} = I(z_1 > 0)$$

Indicator function

$$\frac{\partial z_1}{\partial w_1} = \frac{\partial}{\partial w_1} (W_1 x^{(i)} + b_1) = x^{(i)T} \Rightarrow \delta_\omega^{(i)} = x^{(i)T} \quad (6)$$

$$\frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial \hat{y}^{(i)}} \times \frac{\partial y^{(i)}}{\partial z_r} \times \frac{\partial z_r}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \times \frac{\partial z_1}{\partial w_1} \quad (7)$$

$$\frac{\partial J}{\partial w_1} = \delta_1^{(i)} \times \delta_r^{(i)} \times \delta_\mu^{(i)} \times \delta_\tau^{(i)} \times \delta_\omega^{(i)}$$

$$\frac{\partial J}{\partial w_1} = -\frac{1}{m} \left[\frac{y^{(i)} - \hat{y}^{(i)}}{\hat{y}^{(i)}(1 - \hat{y}^{(i)})} \right] \times \sigma(z_1)(1 - \sigma(z_1)) \times W_1 \times I(z_1 > 0) \times \mathcal{X}^{(i)T}$$

برای پیاده سازی vectorized با m نمونه این مقدار برابر با جمع همه نمونه ها می شود.

$$\frac{\partial J}{\partial w_1} = \sum_{i=1}^m \delta_1^{(i)} \times \delta_2^{(i)} \times \delta_3^{(i)} \times \delta_4^{(i)} \times \delta_5^{(i)}$$

$$\frac{\partial J}{\partial w_1} = -\frac{1}{m} \sum_{i=1}^m \frac{y^{(i)} - \hat{y}^{(i)}}{\hat{y}^{(i)}(1 - \hat{y}^{(i)})} \cdot \sigma(z_1^{(i)})(1 - \sigma(z_1^{(i)})) \times W_1 \times I(z_1^{(i)} > 0) \times \mathcal{X}^{(i)T}$$

1399

- 8. 1st hidden layer -

سوال ۴ - الف) forward pass

$$Z^{(1)} = \begin{bmatrix} a & -a \\ b & a \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix}$$

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

12 • 2nd hidden layer:

13 (2)
$$\begin{bmatrix} a & b \\ -a & -b \end{bmatrix} \begin{bmatrix} \sigma(a) \\ \sigma(b) \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} \frac{a}{1+e^{-a}} + \frac{b}{1+e^{-b}} \\ -\frac{b}{1+e^{-a}} - \frac{a}{1+e^{-b}} \end{bmatrix}$$

$$h^{(2)} = \text{ReLU}(z^{(2)}) = \begin{bmatrix} \max(0, \frac{a}{1+e^{-a}} + \frac{b}{1+e^{-b}}) \\ \max(0, \frac{-b}{1+e^{-a}} + \frac{-a}{1+e^{-b}}) \end{bmatrix} = \begin{bmatrix} h_1^{(2)} \\ h_2^{(2)} \end{bmatrix}$$

16. 3rd hidden layer:

$$Z^{(3)} = \begin{bmatrix} a & -b \\ -a & b \end{bmatrix} \begin{bmatrix} h_1^{(2)} \\ h_2^{(2)} \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} ah_1^{(2)} - bh_2^{(2)} \\ -ah_1^{(2)} + bh_2^{(2)} \end{bmatrix} = \begin{bmatrix} z_1^{(3)} \\ z_2^{(3)} \end{bmatrix}$$

$$h^{(3)} = \text{softmax}(Z^{(3)}) = \begin{bmatrix} \frac{e^{z_1^{(3)}}}{e^{z_1^{(3)}} + e^{z_2^{(3)}}} \\ \frac{e^{z_2^{(3)}}}{e^{z_1^{(3)}} + e^{z_2^{(3)}}} \end{bmatrix} = \begin{bmatrix} h_1^{(3)} \\ h_2^{(3)} \end{bmatrix}$$

output layer:

$$z^{(4)} = \begin{bmatrix} a & -b \\ 0.5a & a \end{bmatrix} \begin{bmatrix} \frac{e^{z_1^{(3)}}}{e^{z_1^{(3)}} + e^{z_2^{(3)}}} \\ \frac{e^{z_2^{(3)}}}{e^{z_1^{(3)}} + e^{z_2^{(3)}}} \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} ah_1^{(3)} - bh_2^{(3)} \\ 0.5ah_1^{(3)} + ah_2^{(3)} \end{bmatrix}$$

$$a^{(4)} = h^{(4)} = \text{identity}(z^{(4)}) = \begin{bmatrix} z_1^{(4)} \\ z_2^{(4)} \end{bmatrix} = \begin{bmatrix} ah_1^{(3)} - bh_2^{(3)} \\ 0.5ah_1^{(3)} + ah_2^{(3)} \end{bmatrix}$$

(ما به جای a از h استفاده کردیم) همان خروجی a خواسته شده است

$$a^{(4)} = h^{(4)}$$

(ب) محاسبه خطای E

$$E = \frac{1}{2} ((a_1^{(4)} - 0)^2 + (a_2^{(4)} - 1)^2)$$

$$E = \frac{1}{2} ((h_1^{(4)} - 0)^2 + (h_2^{(4)} - 1)^2)$$

backward pass (2)

$$\frac{\partial E}{\partial w^{(4)}} = \frac{\partial E}{\partial h^{(4)}} \times \frac{\partial h^{(4)}}{\partial z^{(4)}} \times \frac{\partial z^{(4)}}{\partial w^{(4)}}$$

$$\frac{\partial E}{\partial w^{(4)}} = \begin{bmatrix} h_1^{(4)} \\ h_2^{(4)} - 1 \end{bmatrix} \times 1 \times h^{(3)T} = \begin{bmatrix} h_1^{(4)} \\ h_2^{(4)} - 1 \end{bmatrix} \begin{bmatrix} h_1^{(3)} & h_2^{(3)} \end{bmatrix}$$

۳۰	۲۹	۲۸	۲۷	۲۶	۲۵	۲۴	۲۳	۲۲	۲۱	۲۰	۱۹	۱۸	۱۷	۱۶	۱۵	۱۴	۱۳	۱۲	۱۱	۱۰	۹	۸	۷	۶	۵	۴	۳	۲	۱
----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	---	---	---	---	---	---	---	---	---

(2)

$$\frac{\partial E}{\partial W^{(2)}} = W^{(3)T} \times 1 \times dz^{(3)} \times h^{(1)}$$

$$\frac{\partial E}{\partial w^{(2)}} = W^{(3)T} W^{(4)T} \begin{bmatrix} h_1^{(3)}(1-h_1^{(3)}) & -h_1^{(3)}h_2^{(3)} \\ -h_2^{(3)}h_1^{(3)} & h_2^{(3)}(1-h_2^{(3)}) \end{bmatrix} \begin{bmatrix} h_1^{(4)} \\ h_2^{(4)} - 1 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \end{bmatrix}^T$$

$$\frac{\partial E}{\partial b^{(2)}} = dz^{(2)} \times \frac{\partial z^{(2)}}{\partial b^{(2)}} = W^{(3)T} W^{(4)T} \begin{bmatrix} h_1^{(3)}(1-h_1^{(3)}) & -h_1^{(3)}h_2^{(3)} \\ -h_2^{(3)}h_1^{(3)} & h_2^{(3)}(1-h_2^{(3)}) \end{bmatrix} \begin{bmatrix} h_1^{(4)} \\ h_2^{(4)} - 1 \end{bmatrix}$$

$$\frac{\partial E}{\partial w^{(1)}} = \frac{\partial E}{\partial z^{(2)}} \times \frac{\partial z^{(2)}}{\partial h^{(1)}} \times \frac{\partial h^{(1)}}{\partial z^{(1)}} \times \frac{\partial z^{(1)}}{\partial w^{(1)}}$$

جمعة ١٤
 أذر December 2020 4
 ١٨ ربيع الثاني ١٤٤٢ Friday

$$\frac{\partial E}{\partial w^{(1)}} = 2W^{(2)T} \left[W^{(3)T} W^{(4)T} \begin{bmatrix} h_1^{(3)}(1-h_1^{(3)}) & -h_1^{(3)}h_2^{(3)} \\ -h_2^{(3)}h_1^{(3)} & h_2^{(3)}(1-h_2^{(3)}) \end{bmatrix} \begin{bmatrix} h_1^{(4)} \\ h_2^{(4)} - 1 \end{bmatrix} \right] \underbrace{h^{(1)}(1-h^{(1)})}_{\sigma(z^{(1)})} \sigma'(z^{(1)})^T$$

$$\frac{\partial E}{\partial b^{(1)}} = d z^{(2)} \times \frac{\partial z^{(2)}}{\partial b^{(1)}}$$

$$\frac{\partial E}{\partial b^{(1)}} = z^{(2)T} W^{(2)T} W^{(4)T} \begin{bmatrix} h_1^{(3)}(1-h_1^{(3)}) & -h_1^{(3)}h_1^{(3)} \\ -h_2^{(3)}h_1^{(3)} & h_2^{(3)}(1-h_2^{(3)}) \end{bmatrix} \begin{bmatrix} h_1^{(4)} \\ h_2^{(4)} - 1 \end{bmatrix} \sigma'(z^{(1)})$$

١٩ ربيع الثاني ١٤٤٢

$$W^{(1)'} = W^{(1)} - \eta \frac{\partial E}{\partial W^{(1)}}$$

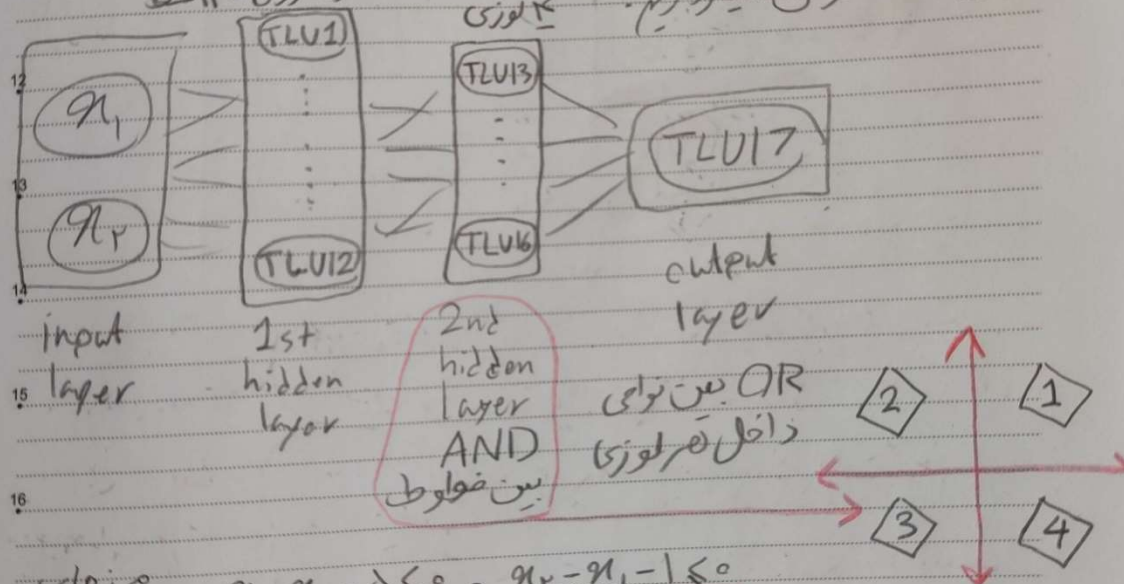
$$\mathcal{H}^T = \begin{bmatrix} 1 & 0 \end{bmatrix}$$

$$W^{(1)'} = \begin{bmatrix} a & -a \\ b & a \end{bmatrix} - \eta \begin{bmatrix} a & -b \\ b & -a \end{bmatrix} \begin{bmatrix} a & -a \\ -b & b \end{bmatrix} \begin{bmatrix} a & 0, \Delta a \\ -b & a \end{bmatrix} A \begin{bmatrix} h_1^{(4)} \\ h_2^{(4)} - 1 \end{bmatrix} h^{(1)}(1-h^{(1)})^8 \mathcal{H}^T$$

$$b^{(1)'} = b^{(1)} - \eta \frac{\partial E}{\partial b^{(1)}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \eta \frac{\partial E}{\partial b^{(1)}}$$

$$b^{(1)'} = -\eta \begin{bmatrix} a & -b \\ b & -a \end{bmatrix} \begin{bmatrix} a & -a \\ -b & b \end{bmatrix} \begin{bmatrix} a & 0, \Delta a \\ -b & a \end{bmatrix} A \begin{bmatrix} h_1^{(4)} \\ h_2^{(4)} - 1 \end{bmatrix} h^{(1)}(1-h^{(1)})$$

۱۳۹۹
۸. الف) برای اینکه ببینیم یک نقطه داخل یک لوزی قرار دارد یا خیر
نیاز به ۴ خط (TLV) داریم. لذا به $4 \times 4 = 16$ TLV در لایه اول نیاز
است ولی با توجه به اینکه ۴ خط مشترک در لوزی ها موجود است بنابراین
به $16 - 4 = 12$ TLV نیاز داریم. فرض کنیم در لایه دوم باید مرزهای
لوزی ها AND بگیریم که نیاز به ۴ TLV داریم در نهایت نیز به یک TLV
برای OR گرفتن نیاز داریم. خطوط مرزی = ۱۲ خط



۱۷. مرزهای شکل ناصفا اول

$$x_2 - x_1 - 1 \leq 0 \text{ و } x_1 + x_2 - 5 \leq 0$$

$$x_2 - x_1 + 1 \geq 0 \text{ و } x_2 + x_1 - 3 \geq 0$$

۱۸. $x_1 > 0$ و $x_2 > 0$

۱۹. مرزهای شکل ناصفا دوم

$$x_2 + x_1 - 1 \leq 0 \text{ و } x_2 + x_1 + 1 \geq 0$$

$$x_2 - x_1 - 3 \geq 0 \text{ و } x_2 - x_1 - 5 \leq 0$$

۲۰. $x_1 < 0$ و $x_2 > 0$

8

9

10

1

12



13

1

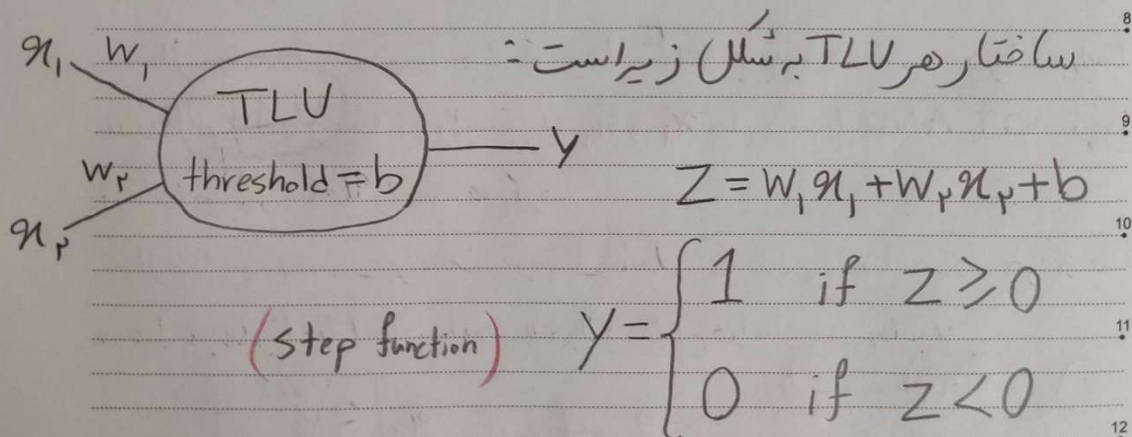
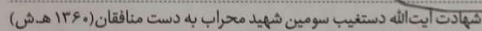
19

20

3

3

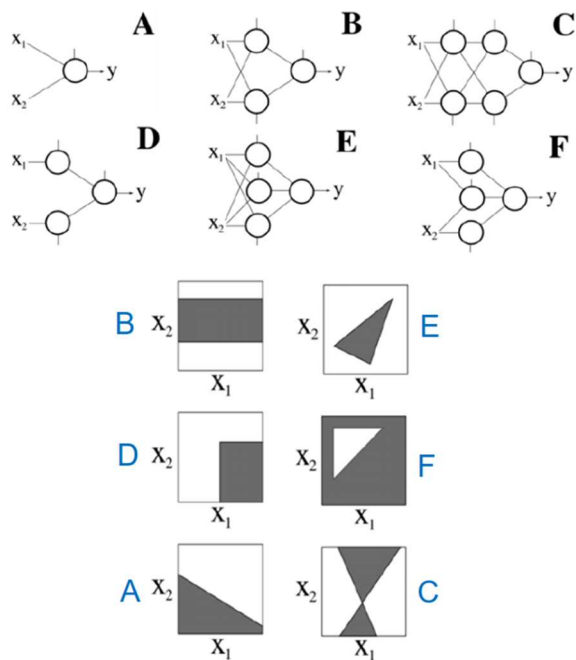
②



۳۰	۲۹	۲۸	۲۷	۲۶	۲۵	۲۴	۲۳	۲۲	۲۱	۲۰	۱۹	۱۸	۱۷	۱۶	۱۵	۱۴	۱۳	۱۲	۱۱	۱۰	۹	۸	۷	۶	۵	۴	۳	۲	۱
----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	---	---	---	---	---	---	---	---	---

④

ج) با توجه به اینکه هر کدام از شبکه ها فقط یک پاسخ را ایجاد کرده، هر ناحیه توسط شبکه ای که در کنارش در شکل زیر نوشته شده است ایجاد شده.



دلایل و علت ها در صفحه بعد آمده است.

سوال ششم

الف) ابتدا بیاید تابع مفاهیم ReLU و ابر صفحه ها و Activation Patterns را بیان کنیم و سپس با استفاده از آنها به اثبات بپردازیم.
تابع فعال سازی ReLU:

تابع فعال سازی واحد خطی اصلاح شده (ReLU) به صورت زیر تعریف می شود:

$$\text{ReLU}(x) = \max(0, x)$$

این تابع تکه ای x را وقتی $x > 0$ و در غیر این صورت 0 را خروجی می دهد. ویژگی های کلیدی ReLU که در اینجا مرتبط هستند عبارتند از:

- خطی در ناحیه مثبت: وقتی $x \geq 0$ باشد، ReLU به عنوان یک تابع identity عمل می کند.
- صفر در ناحیه منفی: وقتی $x < 0$ ، خروجی ReLU صفر (constant) است.

Activation Patterns و ابر صفحه:

هر نورون در شبکه بر اساس حالت فعال سازی خود، یک ابر صفحه را تعریف می کند:

- نورون فعال: ورودی پیش از فعال سازی نورون (z) غیر منفی است ($z \geq 0$).
- نورون غیر فعال: ورودی قبل از فعال سازی نورون منفی است ($z < 0$).

این شرایط را می توان به صورت نابرابری های خطی نوشت. برای یک نورون i در لایه اول، ورودی پیش از فعال سازی به صورت زیر است:

$$z_i = w_i^T x + b_i$$

- w_i بردار وزن نورون i است.
- x بردار ورودی است.
- b_i بایاس است.

ابر صفحه تعریف شده با $z_i = 0$ فضای ورودی را به دو نیم فضا تقسیم می کند:

- نیم فضای مثبت: $w_i^T x + b_i \geq 0$ (نورون فعال است).
- نیم فاصله منفی: $w_i^T x + b_i < 0$ (نورون غیر فعال است).

حال به اثبات می پردازیم:

برای اثبات اینکه $F_A(\mathbb{R}^m; W)$ فضای ورودی را به چند polytopes محدب (convex) تقسیم می کند که در آن $F_A(\mathbb{R}; W)$ در هر ناحیه خطی است:

۱. ابتدا بیایید ببینیم که چه چیزی یک activation pattern را تعیین می کند:

برای هر نورون ReLU، فعال سازی با مثبت یا منفی بودن ورودی آن تعیین می شود.

ورودی هر نورون تابعی خطی از ورودی های آن است: $wx + b$

بنابراین، هر نورون به طور موثر یک ابر صفحه $wx + b = 0$ ایجاد می کند که فضا را تقسیم می کند.

۲. خواص کلیدی:

هر ابر صفحه فضا را به دو ناحیه (مثبت و منفی) تقسیم می کند.

تابع ReLU خطی است (identity) وقتی $\text{input} > 0$ باشد و

constant (صفر) است وقتی $\text{input} \leq 0$ باشد.

۳. اثبات تحدب (convexity):

بیابید دو نقطه x_1 و x_2 را در یک ناحیه مشترک activation pattern در نظر بگیریم.

برای هر نورون در شبکه، اگر برای x_1 و x_2 فعال باشد، پس:

$$wx_1 + b > 0 \text{ و } wx_2 + b > 0$$

برای هر نقطه $x_t = tx_1 + (1-t)x_2$ که در آن $t \in [0,1]$:

$$wx_t + b = w(tx_1 + (1-t)x_2) + b$$

$$wx_t + b = t(wx_1 + b) + (1-t)(wx_2 + b)$$

$$wx_1 + b > 0 \text{ و } wx_2 + b > 0 \text{ , } t > 0 \text{ , } 1-t > 0 \rightarrow t(wx_1 + b) + (1-t)(wx_2 + b) > 0$$

$$wx_t + b > 0$$

این نشان می دهد که x_t همان activation pattern مشترک را حفظ می کند.

بنابراین، هر ناحیه محدب است.

۴. اثبات خطی بودن در نواحی مختلف:

در یک activation pattern ثابت:

- نورون های فعال از ورودی های خود به صورت خطی عبور می کنند (identity function) که رفتاری خطی است.
- نورون های غیرفعال خروجی صفر دارند ($\text{constant} = 0$) که رفتاری خطی است.

که این ترکیبی از توابع خطی ایجاد می کند.

بنابراین، $F_A(x;W)$ در هر ناحیه، خطی است.

(ب)

برای اثبات حد بالای $O(k^{mn})$ برای $A(F_{A(n,k)}(\mathbb{R}^m; W))$:

بیابید لایه به لایه آنالیز کنیم:

- لایه اول: ورودی را از \mathbb{R}^m می گیرد.
- هر لایه بعدی از k نورون ورودی می گیرد.

۱. تحلیل لایه اول:

- ابعاد ورودی: m
- تعداد نورون ها: k
- هر نورون یک ابر صفحه را در \mathbb{R}^m تعریف می کند.
- با استفاده از حد داده شده در راهنمایی: $r(k, m) \leq \sum_{i=0}^m \binom{k}{i} \leq \sum_{i=0}^m k^i \leq O(k^m)$

۲. برای لایه های بعدی ($i = 2, \dots, n$):

- هر نورون یک تابع خطی و به دنبال آن ReLU را محاسبه می کند
- بعد ورودی لایه i ، k است (عرض لایه قبلی)
- تعداد نورون های لایه i ، k است
- هر نورون می تواند حداکثر ۲ حالت فعال سازی مختلف ایجاد کند (۰ یا ۱)

۳. بینش کلیدی برای حد:

- برای هر لایه i ($i \geq 2$):

- هر نورون می تواند حداکثر ۲ حالت داشته باشد
- با k نورون در هر لایه، حداکثر 2^k الگوی فعال سازی ممکن در هر لایه داریم

۴. ترکیب حد ها:

- لایه اول: تعداد pattern ها $O(k^m)$.
- هر لایه بعدی: تعداد pattern ها $O(2^k)$.
- Total number of patterns $\leq O(k^m) * O(2^k)^{n-1}$
- ۵. از آنجایی که k نورون در هر لایه دلالت میکند $2^k \leq k^k$:

$$O(k^m) * O(2^k)^{n-1} \leq O(k^m) * O(k^k)^{n-1} = O(k^{m+k(n-1)})$$

برای هر عدد ثابت $k \leq 1$ ، داریم:

$$k^{m+k(n-1)} \leq k^{mn}$$

این به این دلیل است که برای شبکه های عصبی، معمولاً $k \leq m$ (عرض کوچک تر از اندازه ابعاد ورودی است)، بنابراین:

$$Total\ number\ of\ regions = A(F_{A(n,k)}(\mathbb{R}^m; W) \leq O(k^m) * O(k^k)^{(n-1)} = O(k^{m+k(n-1)}) = O(k^{mn})$$

کران نهایی $O(k^{mn})$ به این دلیل است که:

- تعداد نواحی نمی تواند از حاصل ضرب حداکثر مناطق ممکن در هر لایه بیشتر شود.
- نواحی هر لایه با مجموع دو جمله ای ارائه شده در راهنمایی محدود می شوند.
- ترکیب این مرزها در سراسر n لایه این نتیجه را به ما می دهد.

پایان