

شبکه‌های عصبی

یادگیری عمیق

گردآورنده: سینا جعفری

۱ مقدمه

شبکه‌های عصبی یکی از شاخه‌های مهم و پرکاربرد حوزه یادگیری نظارت‌شده و روش‌های محاسباتی نوین تلقی می‌شوند. این شاخه، با بهره‌گیری از ساختار عصبی مغز انسان‌ها طراحی شده و می‌تواند تقریب و راهکارهای بسیار مناسبی برای حل مسائل مختلف نظارت‌شده در یادگیری ماشین ارائه دهد. شباهت این روش به ساختار عصبی انسان را می‌توان از دو جهت زیر بررسی کرد.

۱. ساختار عصبی انسان همانند ساختار مدل‌های شبکه عصبی، متناسب با یک مجموعه اطلاعات آموزش می‌بیند.

۲. هر دو ساختار، براساس یک سیستم وزن‌دهی از پارامترها کار می‌کنند.

یکی از ساده‌ترین انواع شبکه‌های عصبی مصنوعی^۱، پرسپترون‌ها^۲ هستند که توسط وارن مک‌کالک^۳ و والتر پیتس^۴ معرفی شدند [۱]. اما این الگوریتم اولین بار در آزمایشگاه کرنل آروننیکال توسط فرانک روزنبلات^۵ پیاده‌سازی شد. این الگوریتم در دو نوع تک لایه و چندلایه کار می‌کند که در ادامه فصل به‌طور کامل به هر دوی آن خواهیم پرداخت.

۱.۱ اجزای شبکه عصبی

شبکه‌های عصبی به نوع گسترده‌ای از مدل‌های غیرخطی اشاره می‌کنند که شامل ترکیبی از ضرب‌های ماتریسی و سایر عملیات غیرخطی بر روی ورودی هستند. بطور کلی شبکه عصبی از چهار بخش اصلی ورودی، وزن‌ها (پارامترهای قابل یادگیری مدل^۶)، لایه مخفی^۷ و خروجی تشکیل می‌شود که در ادامه هرکدام را شرح خواهیم داد.

ورودی. ورودی شبکه‌های عصبی باید ماتریسی از اعداد باشند. اگر مجموعه داده S را در اختیار داشته باشیم که:

$$S = ((x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)) \quad s.t. \quad x_i \in \mathbb{R}^d, \quad y_i \in \mathbb{R}$$

که در این جا m و d به ترتیب تعداد نمونه‌ها و اندازه ابعاد فضای ویژگی هستند، می‌توانیم یک فضای ورودی از نوع ماتریس تعریف کنیم که فقط شامل x_i ها می‌باشد. بنابراین داریم:

$$X = (x_1, x_2, \dots, x_m) \quad s.t. \quad x_i \in \mathbb{R}^d$$

¹Artificial Neural Network

²Perceptron

³Warren

⁴Pitts

⁵Frank Rosenblatt

⁶Trainable Parameters

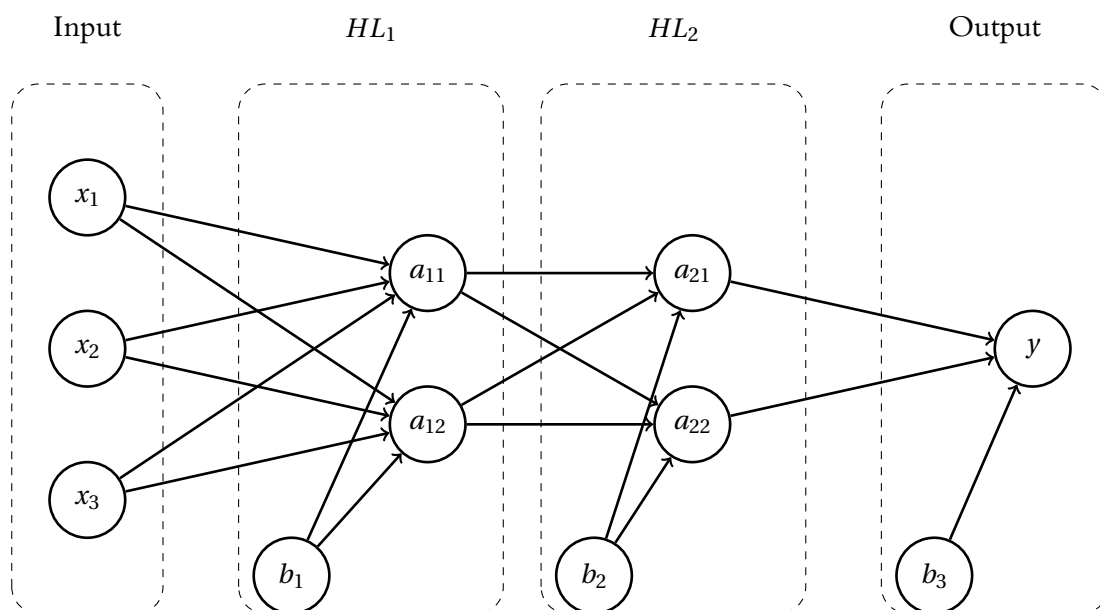
⁷Hidden Layer

توجه کنید در مسائلی که مجموعه داده‌ها همگی تصویر هستند، هر نمونه (تصویر) یک ماتریس با ابعاد $w \times h \times c$ می‌باشد. یعنی در این حالت فضای ویژگی برابرست با:

$$X = (x_1, x_2, \dots, x_m) \quad s.t. \quad x_i \in \mathbb{R}^{w \times h \times c}$$

که در این جا w ، h و c به ترتیب برابرند با طول، ارتفاع و تعداد کانال‌های هر تصویر.

لایه مخفی. یکی از مهم‌ترین اجزای سازنده هر شبکه عصبی، لایه‌های مخفی هستند. تعداد این لایه‌ها و نوع معماری که در آن‌ها قرار دارند، تأثیر بسیار مستقیمی در عملکرد یک شبکه عصبی می‌توانند داشته باشند. هر لایه مخفی در شبکه عصبی، مانند یک پشته^۸ دارای تعدادی نورون^۹ است که مقدار هر کدام از آن‌ها یک ابرپارامتر محسوب می‌شود و نمی‌توان به سادگی در مورد تعداد آن برای یک مسئله خاص نظر داد.



شکل ۱: یک شبکه با چند لایه مخفی.

به عنوان مثال، در شکل ۱ یک شبکه عصبی با چند لایه مخفی (در اصطلاح به این نوع شبکه‌ها، شبکه عصبی چندلایه یا پرسپترون چند لایه که به اختصار *MLP*^{۱۰} است، گوئیم) داریم که دارای دو لایه مخفی است. همان‌طور که مشاهده می‌کنید، هر لایه مخفی دارای سه نورون است که هر نورون برابرست با حاصل ضرب نورون‌های لایه قبل در وزن‌های مربوطه.

وزن‌ها. هر لایه در شبکه عصبی، توسط وزن‌ها به یکدیگر متصل می‌شوند. این وزن‌ها دقیقاً همان پارامترهای قابل یادگیری مدل شبکه عصبی هستند که در طول فرآیند یادگیری، باید بروزرسانی شوند تا به مقدار بهینه خود نزدیک شوند. معمولاً مقداردهی اولیه این وزن‌ها، به صورت تصادفی و از توزیع نرمال انتخاب می‌شود، اما روش‌های بهینه‌تر و کاراتری نیز برای مقداردهی اولیه وزن‌ها وجود دارد. یکی از این روش‌های مؤثر که اخیراً نیز در خیلی از شبکه‌های عصبی استفاده شده است، استفاده از مدل‌های از پیش آموزش دیده^{۱۱} است، که در این روش از وزن‌های یک شبکه عصبی آموزش دیده به عنوان مقادیر اولیه وزن‌های شبکه عصبی موردنظر استفاده می‌شود. بنابراین در شکل ۱، بین لایه ورودی و لایه مخفی اول ماتریس وزن اول، بین لایه مخفی اول و دوم ماتریس وزن دوم و بین لایه مخفی دوم و خروجی ماتریس وزن سوم قرار دارند که به ترتیب با نمادهای $W^{[1]}$ ، $W^{[2]}$ و $W^{[3]}$ آن‌ها را نمایش می‌دهیم. ماتریس وزن کلی یک شبکه عصبی را W نمایش داده و به صورت زیر تعریف می‌کنیم:

$$W = [W^{[1]}, W^{[2]}, W^{[3]}]$$

⁸Stack

⁹Neuron

¹⁰Multi-Layer Perceptron

¹¹Pre-trained Models

خروجی. لایه خروجی هر شبکه عصبی، معمولاً دو حالت مختلف دارد که متناسب است با نوع مسئله. اگر مسئله رگرسیون باشد، در لایه آخر یا خروجی تنها یک نورون وجود دارد که مقدار آن برابرست با مقدار پیش‌بینی شده توسط مدل. اما اگر مسئله از نوع دسته‌بندی باشد، در لایه خروجی به تعداد دسته‌های مسئله نورون خواهیم داشت که هر کدام یک مقدار خاص دارند. اما در این حالت چگونه می‌توان دسته پیش‌بینی شده توسط مدل را بدست آورد؟ اعدادی که هر نورون در لایه آخر دارند ممکن است از نظر عددی با هم اختلاف داشته باشند و نمی‌توان آن‌ها را در این مقیاس مقایسه کرد. یکی از راه‌هایی که برای حل این مشکل ارائه شده است، استفاده از لایه softmax می‌باشد. تابع softmax یک ورودی بردار دریافت کرده و همه اعداد درون بردار را به مقیاس ۰ تا ۱ می‌برد، به‌طوری که مجموع همه اعضای بردار برابر شود با ۱. به عبارت دیگر، مقدار تمامی اعضای بردار خروجی، تابعی از احتمال می‌شوند. این تابع به صورت زیر تعریف می‌شود:

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \quad (1)$$

که در این جا، K برابر با طول بردار می‌باشد. به عنوان مثال، اگر بردار \vec{z} برابر باشد با:

$$\vec{z} = [12, 15, 17, 10]$$

مقدار تابع softmax بر روی این بردار برابرست با:

$$\sigma(\vec{z}) = [0.0059, 0.1184, 0.8749, 0.0008]$$

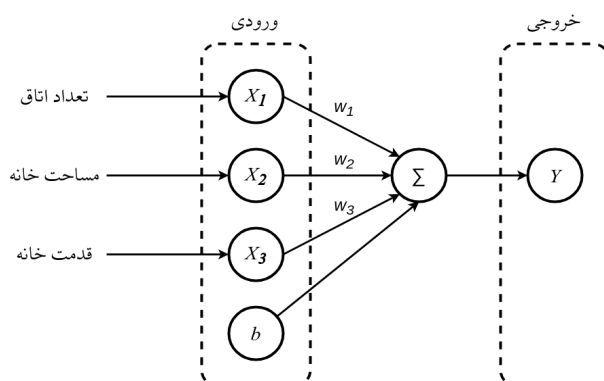
که مجموع اعضای بردار بدست آمده برابرست با ۱.

۲ پرسپترون تک لایه

در این جا برای سادگی، یک شبکه عصبی با ورودی x را با نماد $\bar{h}_\theta(x)$ نمایش می‌دهیم. بنابراین برای تابع \bar{h}_θ داریم:

$$\bar{h}_\theta : X \rightarrow Y \quad (2)$$

یعنی تابع \bar{h}_θ یک نگاشت از فضای ویژگی به فضای برچسب است. حال فرض کنید مجموعه داده‌ای با ویژگی‌های قد، سن و وزن افراد در اختیار داریم و قرار است با این سه ویژگی، دیابت داشتن یا نداشتن یک شخص را پیش‌بینی کنیم، یعنی $y \in \{0, 1\}$. روال کار یک پرسپترون تک لایه بسیار ساده می‌باشد و مقدار هر ویژگی در یک پارامتر قابل یادگیری ضرب و سپس همه‌ی مقادیر با هم جمع می‌شوند.



شکل ۲: یک شبکه پرسپترون تک لایه.

سپس بعد از جمع کردن تمام مقادیر، یک تابع روی خروجی آن اعمال می‌شود که به آن تابع فعالساز^{۱۲} گوئیم. این توابع اغلب غیرخطی هستند و درواقع این توابع غیرخطی در شبکه عصبی باعث غیرخطی شدن عملکرد آن‌ها می‌شود. برای مثال یکی از معروفترین توابع

¹²Activation Function

فعالساز، تابع $ReLU^{13}$ است که به صورت زیر تعریف می‌شود:

$$ReLU(x) = x^+ = \max(0, x) = \frac{x + |x|}{2} = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{o.w.} \end{cases} \quad (3)$$

اما در پرسپترون تک لایه این تابع، خطی و نام آن تابع پله^{۱۴} است، که به صورت زیر تعریف می‌شود:

$$f(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{o.w.} \end{cases} \quad (4)$$

بنابراین، حاصل یک شبکه پرسپترون تک لایه (شکل ۲) برابرست با:

$$y = f(w_1x_1 + w_2x_2 + w_3x_3 + b) = f\left(\sum_{i=1}^3 w_i x_i + b\right) \quad (5)$$

که در این رابطه، w_i وزن پارامتر i ام، x_i مقدار ویژگی i ام و در آخر b نیز مقدار بایاس^{۱۵} است. در شبکه‌های عصبی، بایاس به مقدار ثابتی گفته می‌شود که درون هر لایه از یک شبکه قرار دارد و با اضافه شدن به مجموع حاصل ضرب پارامترها در ویژگی‌ها، می‌تواند نقش اساسی‌ای در تعیین مقدار برچسب داشته باشد. البته در بسیاری از معماری‌ها و لایه‌های مختلف یک شبکه عصبی مقدار بایاس می‌تواند گذاشته نشود، اما در مدل‌های پایه شبکه عصبی این مقدار ثابت گذاشته می‌شود. در ادامه، همانند رگرسیون که می‌توانستیم ویژگی‌ها و پارامترهای مدل را درون دو بردار مجزا قرار دهیم، در شبکه‌های عصبی نیز می‌توانیم این کار را انجام دهیم. درواقع با انجام این کار عملیات ضرب و جمع بین دو عدد تبدیل به ضرب ماتریسی شده و برای حل حالات پیچیده‌تری که در ادامه ذکر می‌شود، محاسبات را ساده‌تر می‌کند. بنابراین اگر بردارهای X و W به ترتیب وزن‌های پارامتر مدل و مقادیر ویژگی‌ها باشند، داریم:

$$y = f(W^T X + b) \quad \text{s.t.} \quad X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \quad W = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}. \quad (6)$$

¹³Rectified Linear Unit

¹⁴Step Function or Binary Step Function

¹⁵Bias

۳ توابع فعالساز

می‌توان از توابع غیرخطی معروف دیگری به جای تابع $ReLU$ که در بخش قبل آورده شد، استفاده کرد. همه این توابع باید یک نگاشت از فضای اعداد حقیقی به اعداد حقیقی باشند، یعنی:

$$\sigma: \mathbb{R} \rightarrow \mathbb{R} \quad (۷)$$

برخی از این توابع غیرخطی عبارتند از:

$$\sigma(z) = \max\{\alpha z, z\} = \begin{cases} z & \text{if } z > 0 \\ \alpha z & \text{o.w.} \end{cases} \quad \text{s.t.} \quad \alpha \in (0, 1) \quad (\text{leaky-ReLU}) \quad (۸)$$

$$\sigma(z) = \begin{cases} z & \text{if } z > 0 \\ \alpha(e^z - 1) & \text{o.w.} \end{cases} \quad \text{s.t.} \quad \alpha > 0 \quad (\text{ELU}) [2] \quad (۹)$$

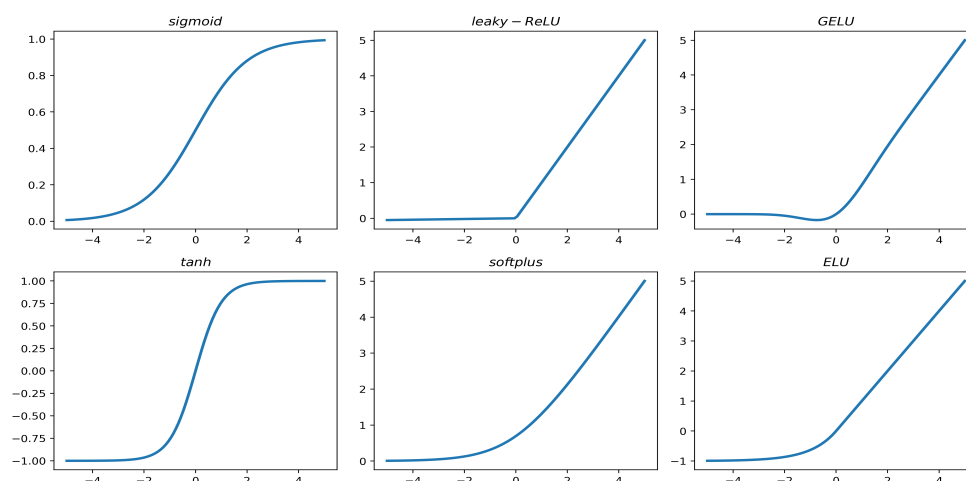
$$\sigma(z) = \frac{z}{2} \left[1 + \operatorname{erf}\left(\frac{z}{\sqrt{2}}\right) \right] \quad (\text{GELU}) [3] \quad (۱۰)$$

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad (\text{sigmoid}) \quad (۱۱)$$

$$\sigma(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad (\text{tanh}) \quad (۱۲)$$

$$\sigma(z) = \frac{1}{\beta} \log(1 + e^{\beta z}) \quad \text{s.t.} \quad \beta > 0 \quad (\text{softplus}) [4] \quad (۱۳)$$

در ادامه نمودار هرکدام از این توابع در شکل ۳ آورده شده است.



شکل ۳: نمودار توابع فعالساز.

یکی از ویژگی‌های مهمی که توابع فعالساز غیرخطی باید دارا باشند، مشتق پذیر بودن است. همان‌طور که می‌دانید در روش گرادیان کاهشی، مدام باید از عمل مشتق‌گیری برای یافتن جواب بهینه استفاده کرد. بنابراین مشتق پذیر بودن این توابع از اهمیت زیادی برخوردار هستند. اگر تابع فعالسازی که انتخاب می‌کنیم، در بعضی نقاط مشتق پذیر نباشد، آنگاه موجب صفر شدن بعضی ضرایب می‌شود که این اتفاق باعث توقف یادگیری می‌شود. در ادامه مشتق برخی توابع فعالساز ذکر شده را یادآوری می‌کنیم.

$$\frac{d}{dz}\sigma(z) = \begin{cases} 1 & \text{if } z > 0 \\ \alpha & \text{o.w.} \end{cases} \quad \text{s.t.} \quad \alpha \in (0, 1) \quad (\text{derivative of leaky-ReLU}) \quad (14)$$

$$\frac{d}{dz}\sigma(z) = \begin{cases} 1 & \text{if } z > 0 \\ \alpha e^z & \text{o.w.} \end{cases} \quad \text{s.t.} \quad \alpha > 0 \quad (\text{derivative of ELU}) \quad (15)$$

$$\frac{d}{dz}\sigma(z) = \sigma(z)(1 - \sigma(z)) \quad (\text{derivative of sigmoid}) \quad (16)$$

$$\frac{d}{dz}\sigma(z) = 1 - \sigma(z)^2 \quad (\text{derivative of tanh}) \quad (17)$$

$$\frac{d}{dz}\sigma(z) = \frac{1}{1 + e^{-z}} \quad (\text{derivative of softplus}) \quad (18)$$

برای بررسی دقیق‌تر مشتق تابع $GELU$ که در رابطه ۱۰ قرار دارد، می‌توانید مقاله [۳] را بررسی کنید.

۴ الگوریتم پیش‌خور

فرآیند یادگیری در شبکه‌های عصبی طی دو مرحله رفت و برگشت صورت می‌گیرد که به ترتیب توسط الگوریتم‌های پیش‌خور^{۱۶} و پس‌انتشار^{۱۷} صورت می‌گیرد. در الگوریتم پیش‌خور تنها کافی است ورودی را به ترتیب در لایه‌های شبکه عصبی ضرب کرده و نتیجه حاصل را بدست آورد. در این مرحله، هیچ‌گونه بروزرسانی در پارامترهای شبکه صورت نمی‌گیرد و تنها خروجی و سپس مقدار تابع هزینه را می‌توان بدست آورد. در مرحله بعد، با استفاده از الگوریتم پس‌انتشار مشتق تابع هزینه را نسبت به تمامی پارامترهای شبکه محاسبه کرده و مقدار حاصل را جایگزین می‌کنیم. این روش در بخش‌های بعدی به‌طور کامل شرح داده خواهد شد.

۵ پرسپترون چند لایه

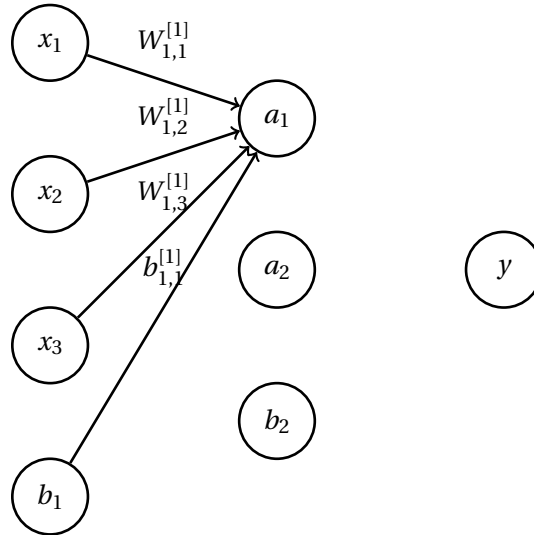
اگر در مجموعه داده موردنظر، الگوهای خطی وجود نداشته باشد، آنگاه دیگر مدل‌های خطی نظیر پرسپترون تک لایه و رگرسیون خطی دیگر جوابگو نیستند. اضافه کردن لایه‌های مخفی متعدد به همراه توابع فعالساز غیرخطی، موجب ساخته شدن مدل‌های غیرخطی می‌شود که توانایی پیدا کردن الگوهای غیرخطی را دارند. می‌دانیم که هر لایه مخفی تعدادی وزن در ورودی دریافت و تعدادی وزن نیز در خروجی برمی‌گرداند. حال یک مثال از یک شبکه پرسپترون تک لایه به همراه ۳ نورون را بررسی می‌کنیم. در شکل ۴، مقدار نورون

¹⁶Feed-Forward

¹⁷Backpropagation

اول در لایه اول یعنی a_1 برابرست با:

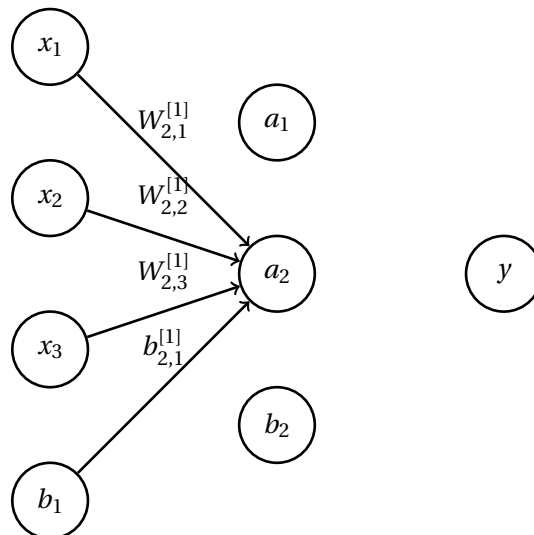
$$\begin{aligned} a_1 &= f(W_{1,1}^{[1]}x_1 + W_{1,2}^{[1]}x_2 + W_{1,3}^{[1]}x_3 + b_{1,1}^{[1]}) \\ &= f\left(\sum_{i=1}^3 W_{1,i}^{[1]}x_i + b_{1,1}^{[1]}\right) \end{aligned} \quad (19)$$



شکل ۴: یک شبکه پرسپترون با یک لایه مخفی به همراه سه نرون.

همچنین در شکل ۵ مقدار a_2 همانند حالت قبل برابرست با:

$$a_2 = f\left(\sum_{i=1}^3 W_{2,i}^{[1]}x_i + b_{2,1}^{[1]}\right) \quad (20)$$



شکل ۵: یک شبکه پرسپترون با یک لایه مخفی به همراه سه نرون.

که در این جا، $W^{[1]}$ نماد وزن لایه اول و نماد $W_{1,i}^{[1]}$ به معنی i -امین وزن از لایه اول که به نرون اول لایه بعد متصل می شود است. همچنین f نیز تابع فعالساز شبکه است که در این جا برای سادگی فرض کرده ایم تمام لایه ها دارای یک تابع فعالساز هستند. حال می توانیم مقادیر وزن ها را به صورت مجزا درون یک بردار قرار دهیم. پس بردارهای زیر را این گونه تعریف می کنیم:

$$W_1^{[1]} = [W_{1,1}^{[1]}, \quad W_{1,2}^{[1]}, \quad W_{1,3}^{[1]}] \quad (21)$$

$$W_2^{[1]} = [W_{2,1}^{[1]}, \quad W_{2,2}^{[1]}, \quad W_{2,3}^{[1]}]$$

پس می‌توانیم مقادیر a_1 و a_2 را به صورت ضرب دو بردار بنویسیم. داریم:

$$a_1 = f(W_1^{[1]T} X + b_{1,1}^{[1]}) \quad (22)$$

$$a_2 = f(W_2^{[1]T} X + b_{2,1}^{[1]})$$

حال می‌توانیم مقادیر a_1 و a_2 را درون برداری به نام h_1 قرار دهیم که به معنای مقدار لایه مخفی اول است. بنابراین داریم:

$$h_1 = [a_1, \quad a_2] \quad (23)$$

در ادامه می‌توانیم دو برداری که در رابطه ۲۱ تعریف شده‌اند را درون یک بردار دیگر به نام $W^{[1]}$ و همچنین می‌توان دو مقدار بایاس لایه اول یعنی $b_{1,1}^{[1]}$ و $b_{2,1}^{[1]}$ را درون یک بردار به نام $b^{[1]}$ قرار داد. پس داریم:

$$W^{[1]} = [W_1^{[1]}, \quad W_2^{[1]}] \quad (24)$$

$$b^{[1]} = [b_{1,1}^{[1]}, \quad b_{2,1}^{[1]}]$$

حال رابطه ۲۳ را به صورت زیر بازنویسی می‌کنیم:

$$h_1 = W^{[1]T} X + b^{[1]} \quad (25)$$

اگر یک لایه مخفی دیگر به شکل ۵ اضافه کنیم، رابطه این لایه با استنتاج از حالت قبل، برابرست با:

$$h_2 = W^{[2]T} h_1 + b^{[2]} \quad (26)$$

بنابراین می‌توان نتیجه گرفت در لایه i -ام از یک شبکه پرسپترون چند لایه، مقدار این لایه برابرست با:

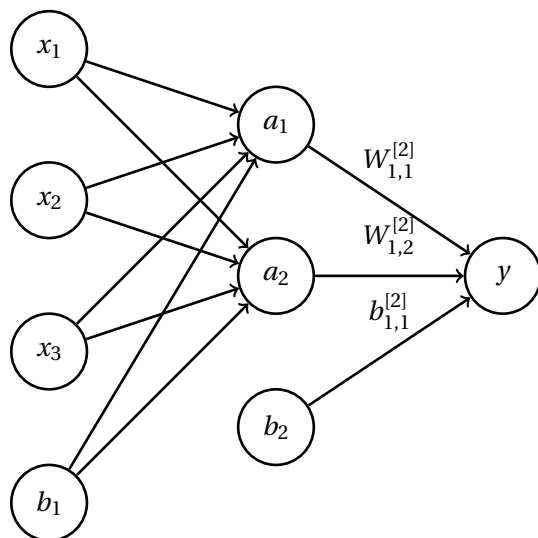
$$h_i = W^{[i]T} h_{i-1} + b^{[i]} \quad (27)$$

پس خروجی یک شبکه پرسپترون با r لایه، که با نماد \bar{h}_θ نمایش می‌دهیم، برابرست با:

$$\bar{h}_\theta(X) = W^{[r]T} h_{r-1} + b^{[r]} \quad (28)$$

بنابراین مقدار y در شکل ۶ برابرست با:

$$\begin{aligned} y &= a_1 W_{1,1}^{[2]} + a_2 W_{1,2}^{[2]} + b_{1,1}^{[2]} \\ &= W^{[2]T} h_1 + b^{[2]} \end{aligned} \quad (29)$$



شکل ۶: یک شبکه پرسپترون با یک لایه مخفی به همراه سه نورون.

مراجع

- [1] W. Mcculloch and W. Pitts, "A logical calculus of ideas immanent in nervous activity," *Bulletin of Mathematical Biophysics*, vol. 5, pp. 127–147, 1943.
- [2] D.-A. Clevert, T. Unterthiner, and S. Hochreiter, "Fast and accurate deep network learning by exponential linear units (elus)," *arXiv preprint arXiv:1511.07289*, 2015.
- [3] D. H. and Kevin Gimpel, "Bridging nonlinearities and stochastic regularizers with gaussian error linear units," *CoRR*, vol. abs/1606.08415, 2016.
- [4] X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks," in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pp. 315–323, JMLR Workshop and Conference Proceedings, 2011.