

رسالة محمد



یادگیری عمیق

مدرس: محمدرضا محمدی

زمستان ۱۴۰۱

پرسپترون چندلایه

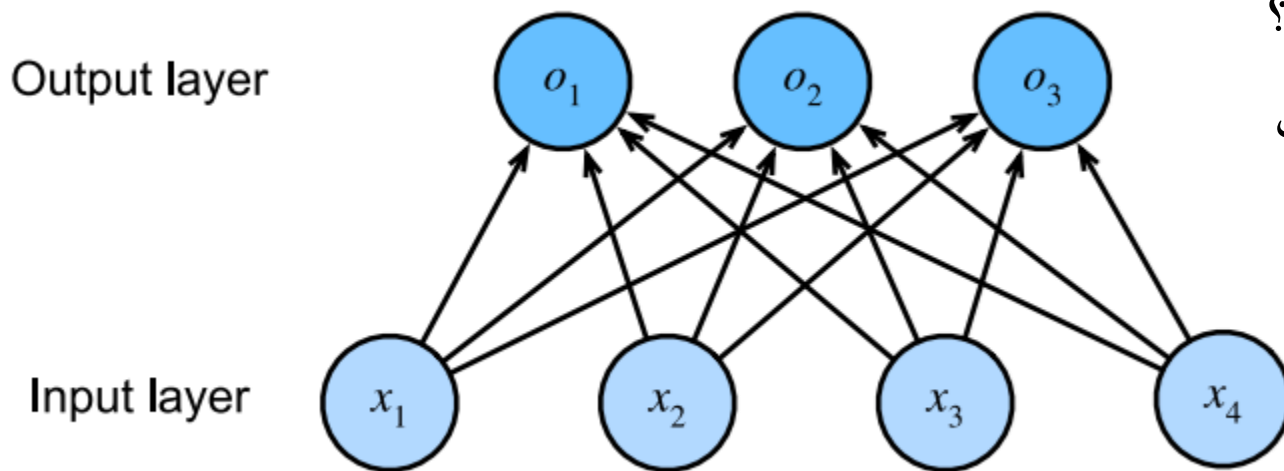
Multilayer Perceptron

مدل‌های خطی

- بر فرض یکنواختی دلالت دارد

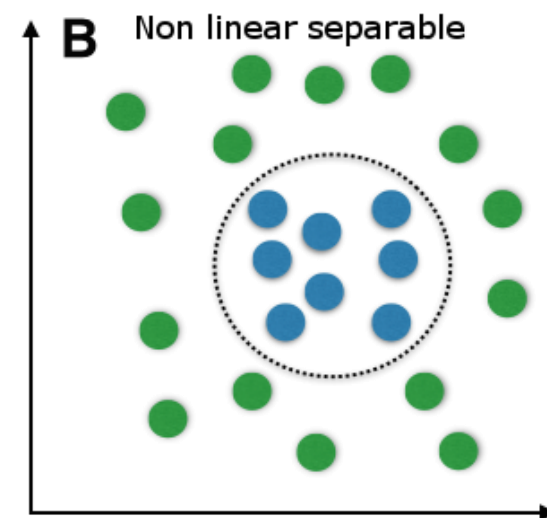
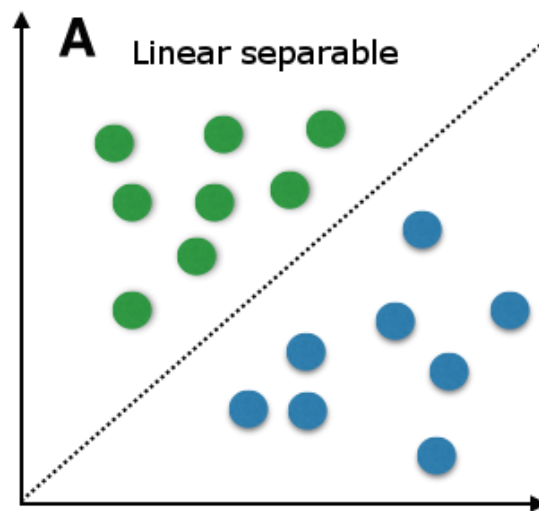
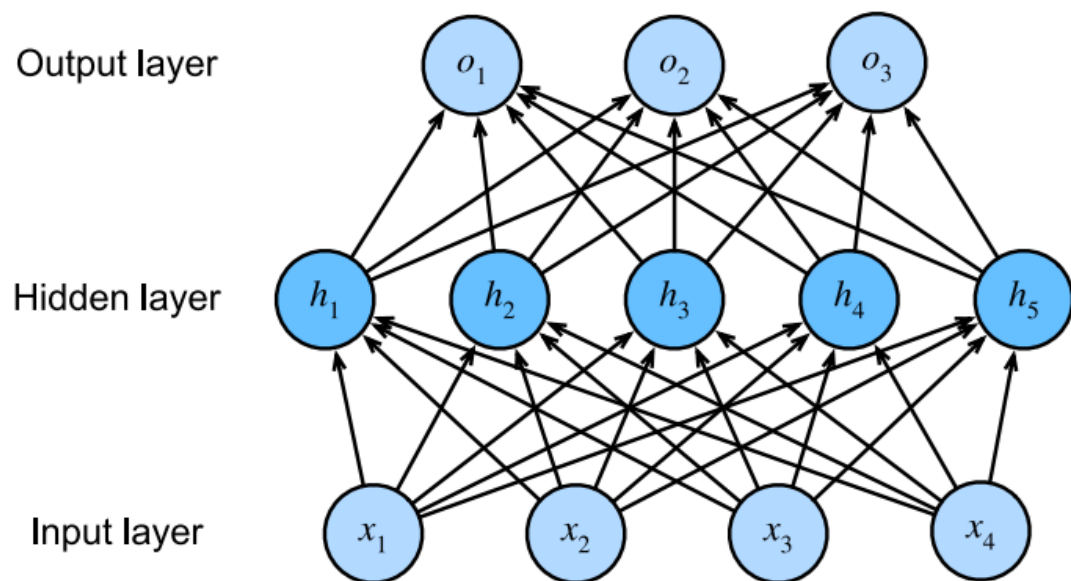
- هر افزایشی در ویژگی همواره منجر به افزایش (کاهش) خروجی می‌شود اگر وزن مربوطه + (-) باشد
- مثال: احتمال بازپرداخت اقساط بر حسب ویژگی‌هایی شامل درآمد
 - اثر افزایش درآمد از ۰ به ۵ میلیون متناظر با افزایش از ۱۰۰ به ۱۰۵ میلیون است
- مثال: احتمال مرگ بر حسب ویژگی‌هایی شامل دمای بدن
 - هر چه دما بالاتر باشد ریسک بیشتری دارد یا کمتر؟

- استخراج ویژگی قبل از لایه خطی بسیار مهم است



پرسپترون چندلایه

- ساده‌ترین شبکه عمیق است که از چندین لایه کاملاً متصل تشکیل می‌شود
- وظیفه لایه‌های میانی استخراج بازنمایی مناسبی است که مسئله مورد نظر در فضای جدید به صورت خطی قابل حل باشد



پرسپترون چندلایه

$$\mathbf{X} \in \mathbb{R}^{n \times d}$$

$$\mathbf{H} = \mathbf{XW}^{(1)} + \mathbf{b}^{(1)}$$

$$\mathbf{H} \in \mathbb{R}^{n \times h}$$

$$\mathbf{W}^{(1)} \in \mathbb{R}^{d \times h}$$

$$\mathbf{b}^{(1)} \in \mathbb{R}^{1 \times h}$$

$$\mathbf{O} = \mathbf{HW}^{(2)} + \mathbf{b}^{(2)}$$

$$\mathbf{O} \in \mathbb{R}^{n \times q}$$

$$\mathbf{W}^{(2)} \in \mathbb{R}^{h \times q}$$

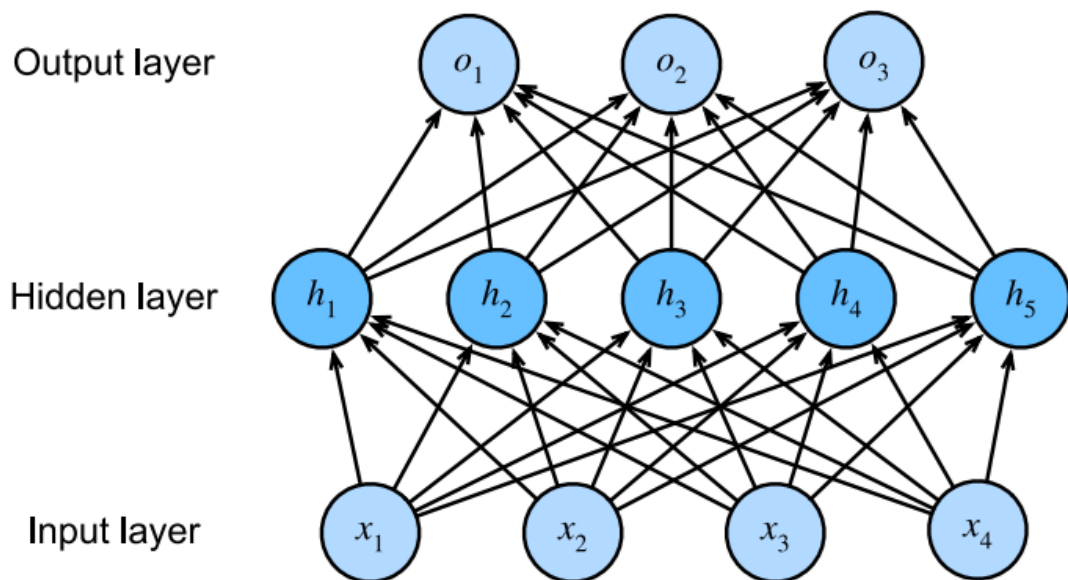
$$\mathbf{b}^{(2)} \in \mathbb{R}^{1 \times q}$$

$$\mathbf{O} = (\mathbf{XW}^{(1)} + \mathbf{b}^{(1)})\mathbf{W}^{(2)} + \mathbf{b}^{(2)}$$

$$= \mathbf{XW}^{(1)}\mathbf{W}^{(2)} + \mathbf{b}^{(1)}\mathbf{W}^{(2)} + \mathbf{b}^{(2)}$$

$$= \mathbf{XW} + \mathbf{b}$$

• دو لایه خطی متوالی هیچ تفاوتی با یک لایه ندارد



مثال

$$L = \frac{1}{4} \sum_{i=1}^4 (y^{(i)} - \hat{y}^{(i)})^2$$

• تابع ضرر MSE:

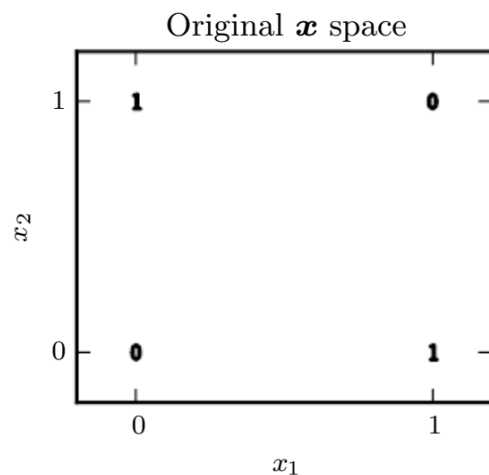
$$\hat{y}^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)} + b$$

• مدل خطی:

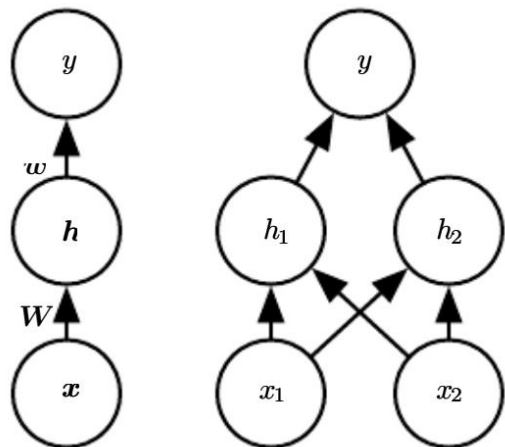
• با استفاده از بهینه‌سازی تحلیلی به پاسخ زیر می‌رسیم

$$\mathbf{w} = \mathbf{0}, b = 0.5$$

- $L = 0.25$ خواهد شد



مثال



$$L = \frac{1}{4} \sum_{i=1}^4 (y^{(i)} - \hat{y}^{(i)})^2$$

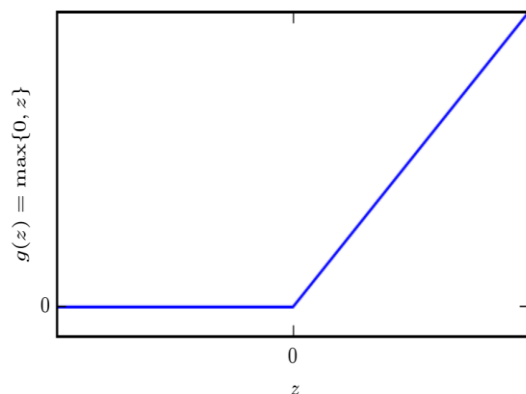
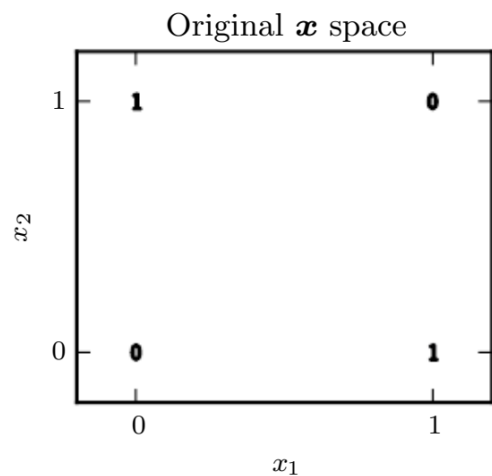
• تابع ضرر MSE:

$$\hat{y}^{(i)} = f^{(2)}(f^{(1)}(\mathbf{x}))$$

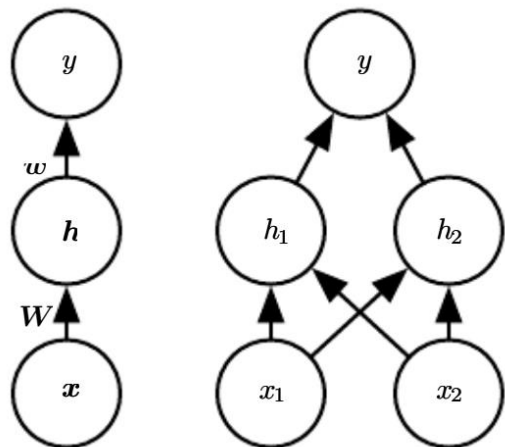
• مدل ۲ لایه:

$$\hat{y}^{(i)} = \mathbf{w}^T g(\mathbf{W}^T \mathbf{x}^{(i)} + \mathbf{c}) + b$$

$$\hat{y}^{(i)} = \mathbf{w}^T \max(0, \mathbf{W}^T \mathbf{x}^{(i)} + \mathbf{c}) + b$$



مثال

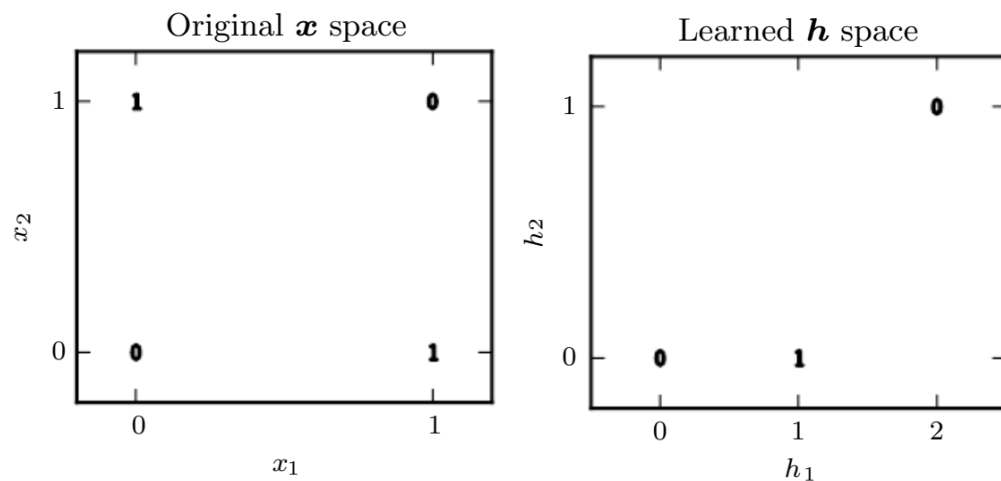


$$L = \frac{1}{4} \sum_{i=1}^4 (y^{(i)} - \hat{y}^{(i)})^2$$

• تابع ضرر MSE:

• مدل ۲ لایه: $\hat{y}^{(i)} = \mathbf{w}^T \max(0, \mathbf{W}^T \mathbf{x}^{(i)} + \mathbf{c}) + b$

• پاسخ زیر با $L = 0$ قابل محاسبه است



$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

$$\mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix} \quad b = 0$$

پرسپترون چندلایه

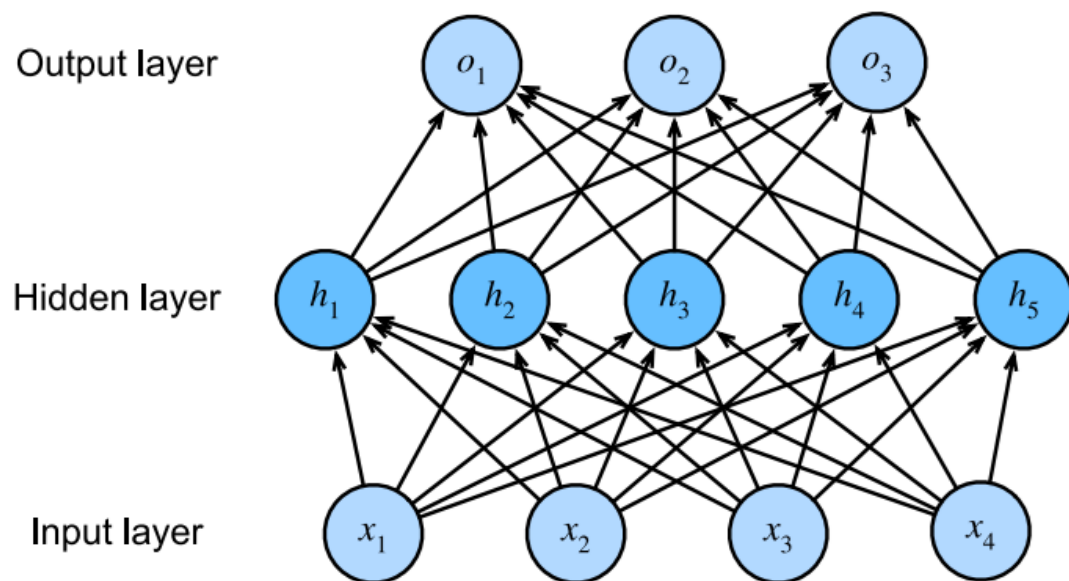
- یک شبکه MLP با تابع فعال‌سازی مناسب در لایه‌های میانی و تعداد نورون کافی می‌تواند هر تابعی را پیاده‌سازی کند

- اصطلاحاً Universal Approximator است

- حتی فقط با یک لایه میانی با تعداد نورون به اندازه کافی زیاد

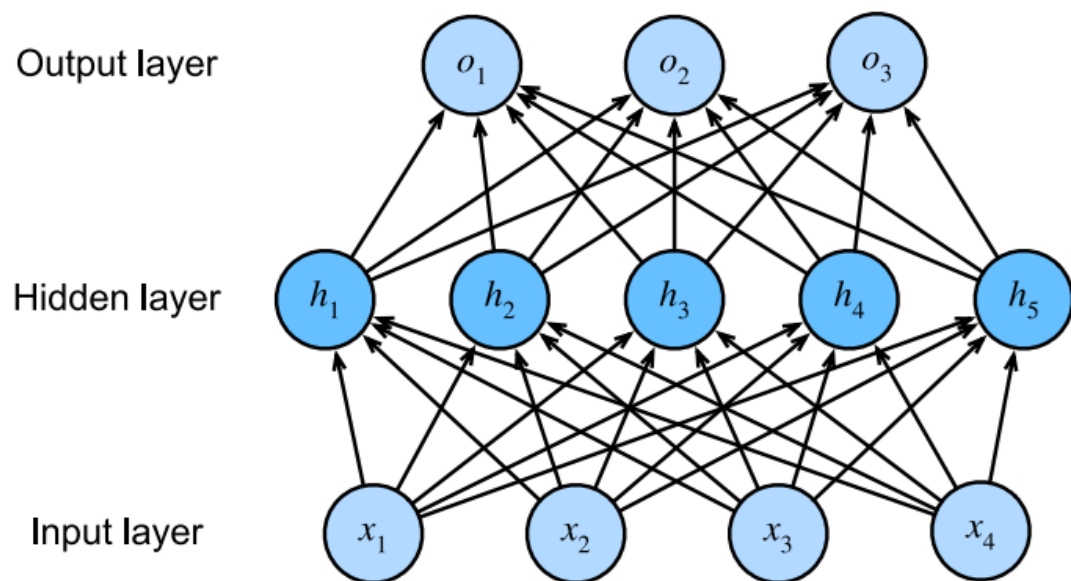
- قسمت سخت کار آموزش وزن‌های مدل است

- برای تقریب بسیاری از توابع می‌توانیم از شبکه‌های عمیق‌تر (بجای عریض‌تر) استفاده کنیم که بسیار فشرده‌تر و کاراتر باشد



توابع فعال سازی

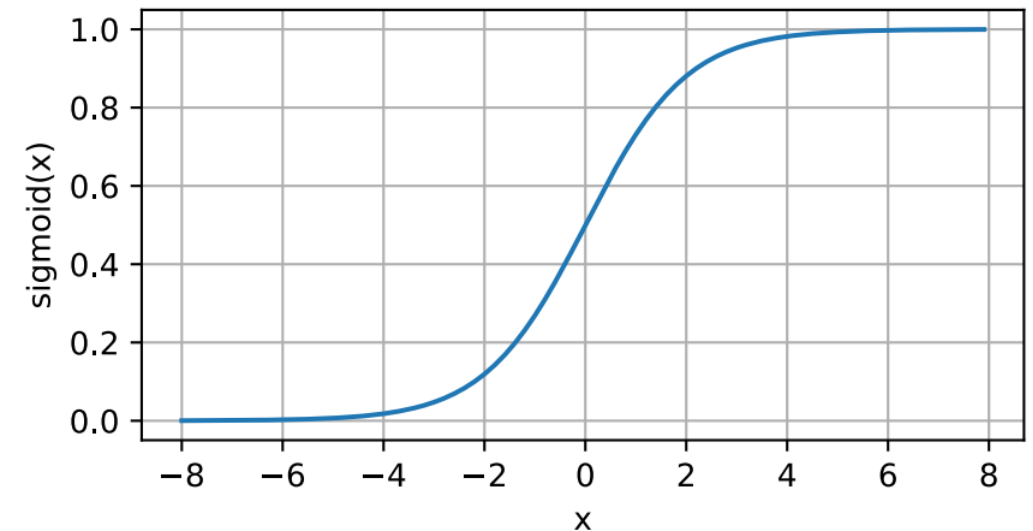
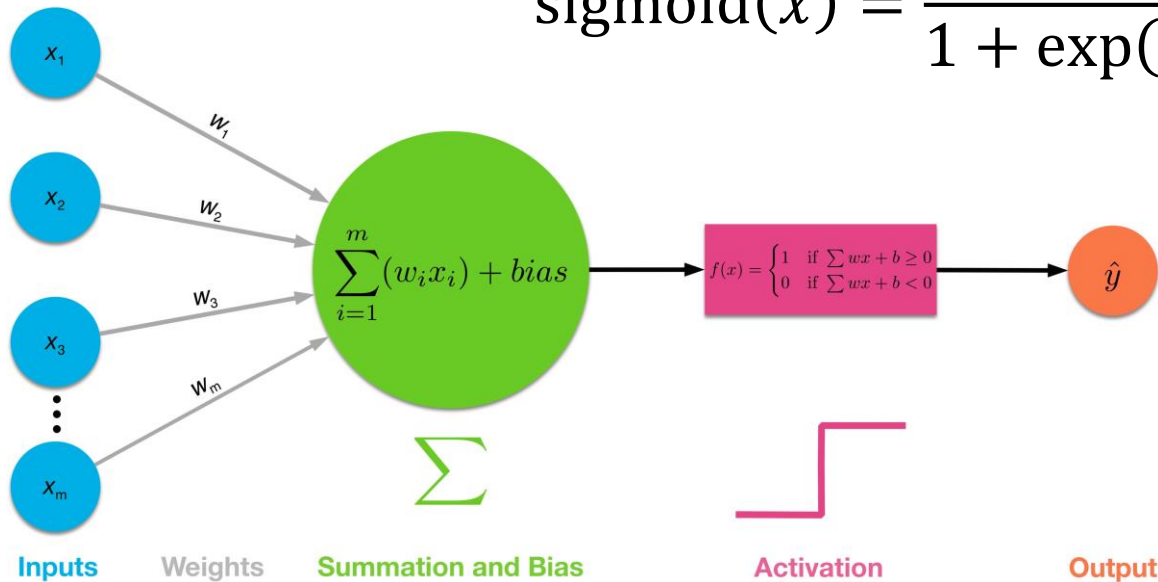
- توابع فعال سازی برای هر نورون مشخص می کنند که خروجی بخشی خطی در چه شرایطی منجر به فعال شدن نورون شود
- به خصوص در شبکه های عمیق بسیار با اهمیت هستند



Sigmoid

- در اولین شبکه‌های عصبی، دانشمندان علاقه‌مند به مدل‌سازی نورون‌های بیولوژیکی بودند که یا فعال می‌شوند یا نمی‌شوند
- با توسعه روش‌های یادگیری مبتنی بر گرادیان، نیاز بود تا تقریب مشتق‌پذیر استفاده شود

$$\text{sigmoid}(x) = \frac{1}{1 + \exp(-x)}$$

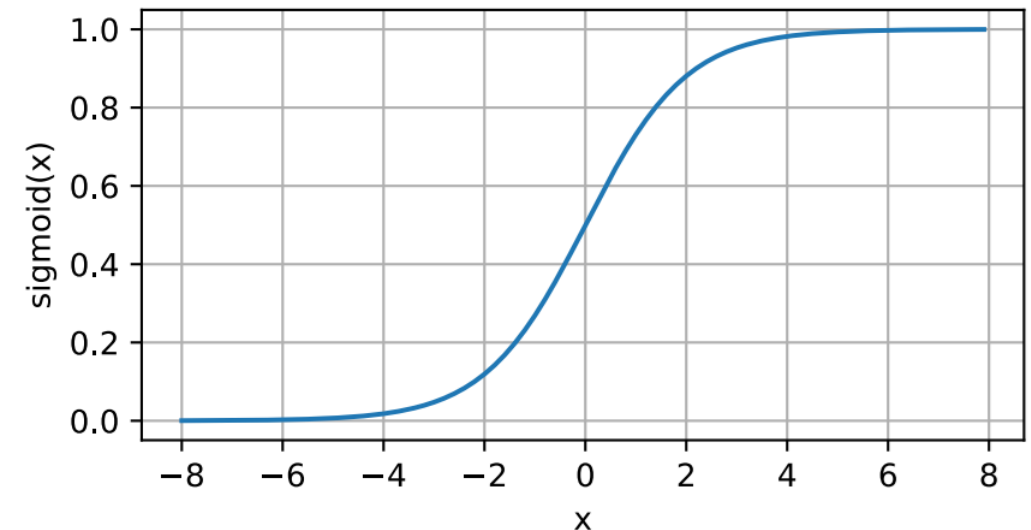
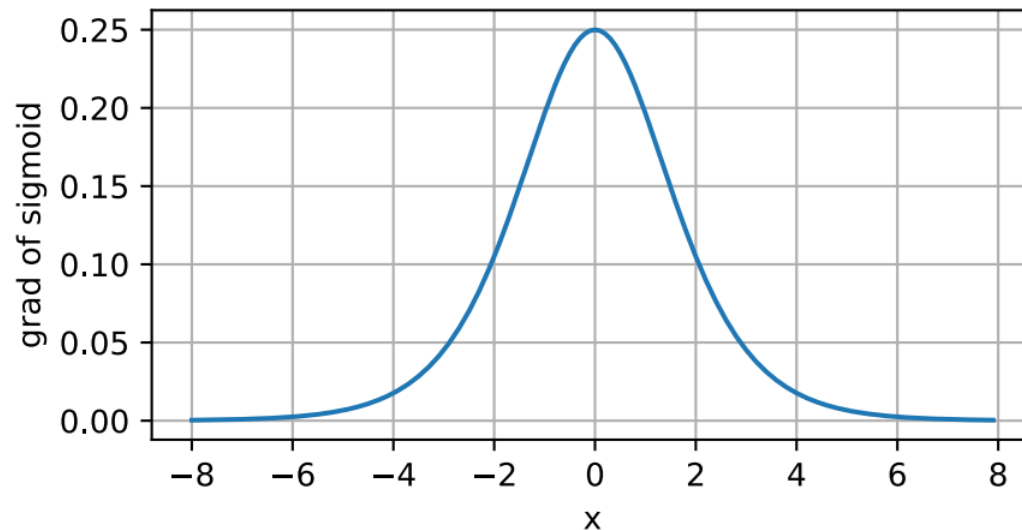


Sigmoid

- در بسیاری از مقادیر، گرادیان نزدیک به صفر است که بهینه‌سازی شبکه‌های عمیق را پیچیده می‌کند
- یکی از ایرادات sigmoid این است که خروجی آن همواره مثبت است

$$\frac{d}{dx} \text{sigmoid}(x) = \text{sigmoid}(x)(1 - \text{sigmoid}(x))$$

$$\text{sigmoid}(x) = \frac{1}{1 + \exp(-x)}$$



Sigmoid

- در بسیاری از مقادیر، گرادیان نزدیک به صفر است که بهینه‌سازی شبکه‌های عمیق را پیچیده می‌کند
- یکی از ایرادات sigmoid این است که خروجی آن همواره مثبت است

- فرض کنید تمام ورودی‌های یک لایه مثبت باشند

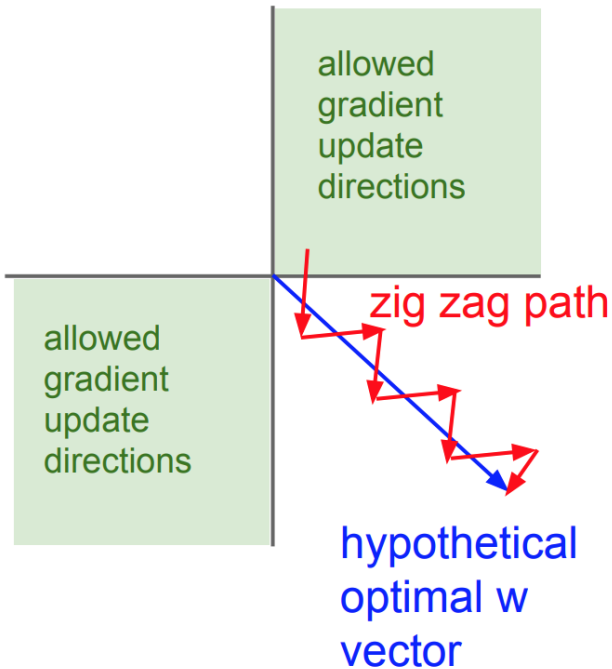
$$z = \sum_k w_k h_k + b$$

- راجع به گرادیان نسبت به \mathbf{w} چه می‌توانیم بگوئیم؟

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial w_i} = \frac{\partial L}{\partial z} h_i$$

▪ تمام مقادیر مثبت یا تمام مقادیر منفی خواهند بود!

- بهینه‌سازی وزن‌ها در برخی راستاها زیگ‌زاگی خواهد بود

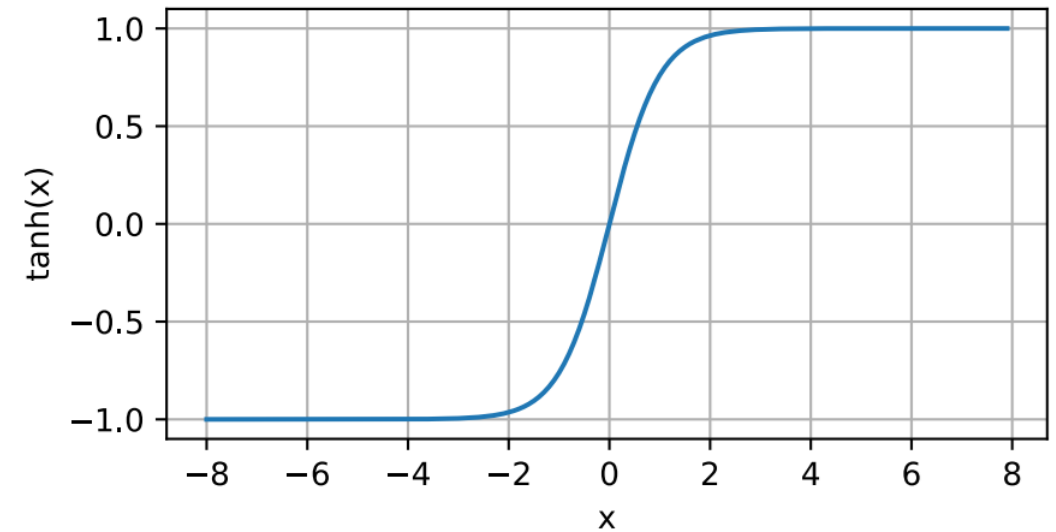
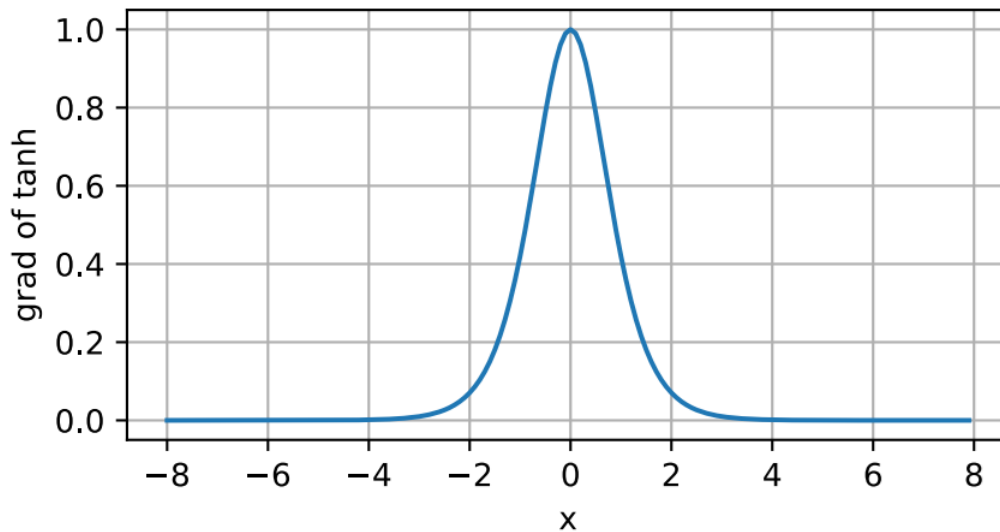


Tanh

- مشابه با تابع Sigmoid خروجی آن محدود است اما به بازه -1 تا $+1$
- اشباع شدن تابع \tanh ، بهینه‌سازی را دشوار می‌کند و از لحاظ محاسباتی هم به دلیل \exp پرهزینه است

$$\frac{d}{dx} \tanh(x) = 1 - \tanh^2(x)$$

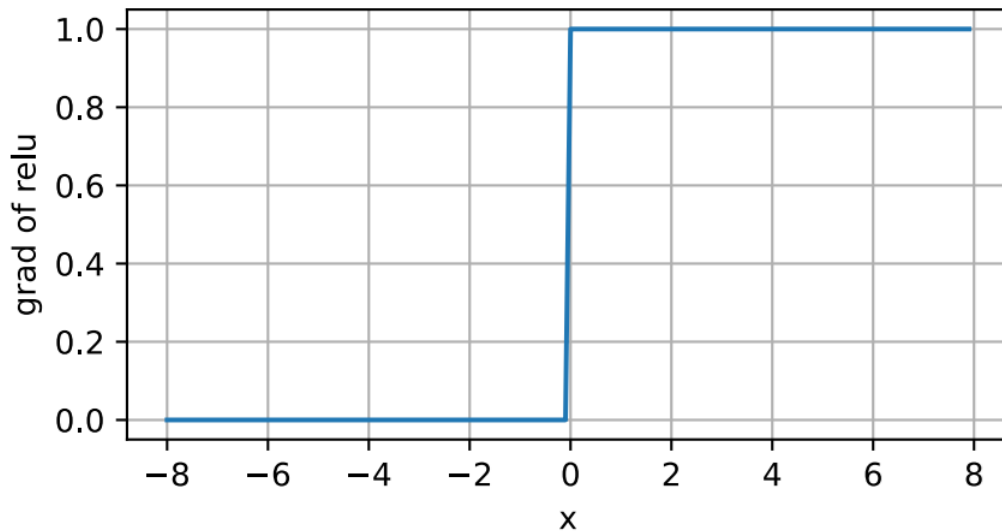
$$\tanh(x) = \frac{1 - \exp(-2x)}{1 + \exp(-2x)} = 2\sigma(2x) - 1$$



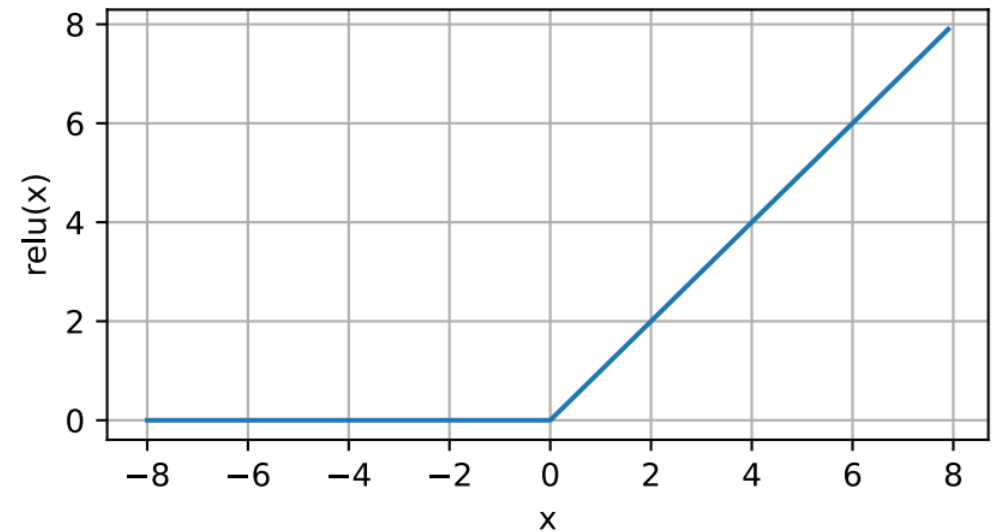
Rectified Linear Unit (ReLU)

- یک تابع غیرخطی بسیار ساده و پرکاربرد است
- بهینه‌سازی ReLU بسیار آسان است زیرا بسیار شبیه به واحدهای خطی است
- مشتق ReLU برای مقادیری که فعال است همواره بزرگ است

$$\frac{d}{dx} \text{ReLU}(x) = x > 0$$



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



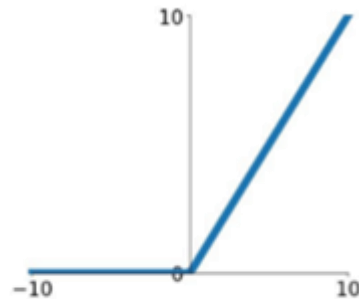
ReLU با شیب مخالف صفر

0	0
1	1
2	2
3	3
4	4
5	5
6	6
7	7
8	8
9	9

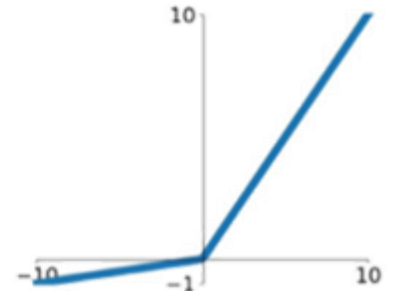
- تابع قدر مطلق ($g(z) = |z|$) با استفاده از $\alpha_i = -1$
- تابع Leaky ReLU از یک مقدار ثابت کوچک مانند $\alpha_i = 0.01$ استفاده می کند
- نسخه Parametric ReLU (PReLU) α_i را یک متغیر قابل آموزش در نظر می گیرد

$$h_i = g(\mathbf{z}, \boldsymbol{\alpha})_i = \max(0, z_i) + \alpha_i \min(0, z_i)$$

ReLU
 $\max(0, x)$

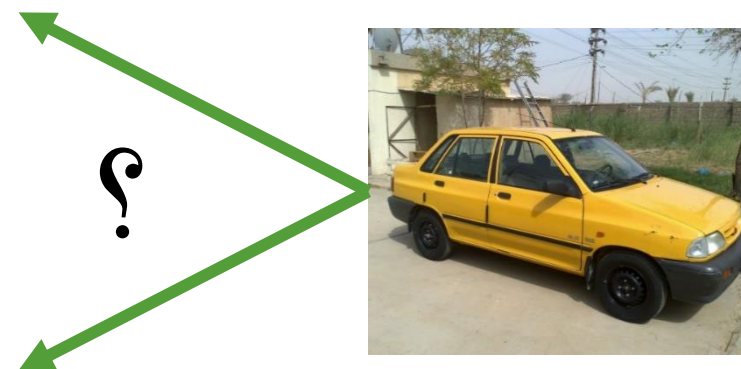


Leaky ReLU
 $\max(0.1x, x)$



انتخاب مدل

- در یادگیری ماشین، هدف ما کشف الگوها است
- اما چگونه می‌توانیم مطمئن باشیم که واقعاً یک الگوی عمومی را کشف کرده‌ایم و صرفاً داده‌های خود را حفظ نکرده‌ایم؟



انتخاب مدل

- چگونه الگوهایی را کشف کنیم که تعمیم‌پذیر باشند؟
- چالش بزرگ این است که در زمان آموزش مدل فقط به مجموعه کوچکی از داده‌ها دسترسی داریم

