

رسالة محمد



# یادگیری عمیق

مدرس: محمدرضا محمدی

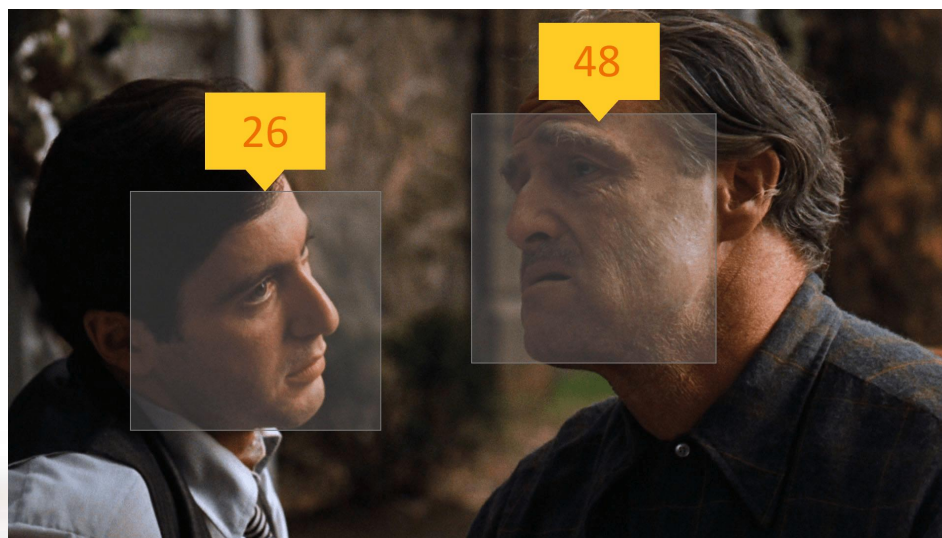
زمستان ۱۴۰۱

# شبکه‌های عصبی خطی

## Linear Neural Networks

# رگرسیون

- رگرسیون به مجموعه‌ای از روش‌ها برای مدل‌سازی رابطه بین یک یا چند متغیر مستقل و یک متغیر وابسته اشاره دارد
- تخمین یک مقدار عددی مورد نظر است



# رگرسیون خطی

- فرض می‌شود رابطه میان متغیرهای مستقل  $\mathbf{X}$  و متغیر وابسته  $y$  خطی است
  - جمع وزن‌دار
- نویز موجود رفتار مناسبی دارد (مانند توزیع نرمال)
- فرض کنید می‌خواهیم قیمت خانه را برحسب مساحت و سال ساخت تخمین بزنیم

features		target
area	age	price
63	13	1.4
78	2	1.95
95	9	2.12
$\vdots$	$\vdots$	$\vdots$

$$\text{price} = w_{\text{area}} \cdot \text{area} + w_{\text{age}} \cdot \text{age} + b$$

sample  $\mathbf{x}^{(i)} = [x_1^{(i)}, x_2^{(i)}]^T, \quad y^{(i)}$

# رگرسیون خطی

- رگرسیون خطی برای یک بردار ویژگی  $d$  بعدی

$$\hat{y} = w_1 x_1 + \dots + w_d x_d + b$$

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

$$\mathbf{x} \in \mathbb{R}^d, \mathbf{w} \in \mathbb{R}^d$$

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

- برای یک مجموعه داده شامل  $n$  نمونه

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w} + b$$

- هدف این است که بتوانیم پارامترهای  $\mathbf{w}$  و  $b$  را به گونه‌ای تخمین بزنیم که برای یک نمونه داده جدید که از همان توزیع  $\mathbf{X}$  نمونه‌برداری شده باشد، پیش‌بینی با حداقل خطا را داشته باشد

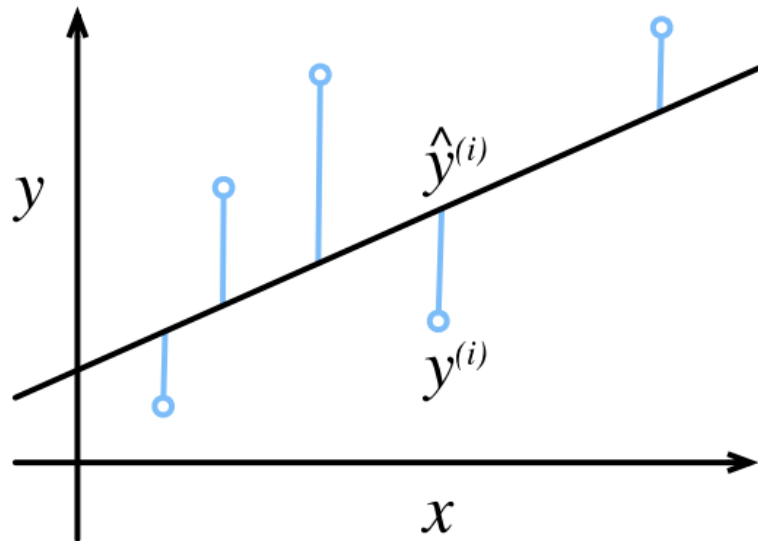
# رگرسیون خطی

$$\hat{y} = \mathbf{X}\mathbf{w} + b$$

- حتی اگر بهترین مدل برای پیش‌بینی  $y$  بر اساس  $\mathbf{x}$  خطی باشد، نمی‌توان توقع داشت در مجموعه داده واقعی تمام مقادیر  $\mathbf{w}^T \mathbf{x}^{(i)} + b$  دقیقاً برابر با  $y^{(i)}$  باشند  
- ممکن است مقداری خطای اندازه‌گیری داشته باشیم (نویز)
- برای یافتن مقادیر مناسب پارامترهای مدل نیاز است تا یک معیار مناسب برای ارزیابی عملکرد مدل تعریف کنیم

# تابع ضرر

- تابع ضرر فاصله میان مقادیر واقعی و مقادیر پیش‌بینی شده را اندازه‌گیری می‌کند
- به طور معمول مقدار تابع ضرر یک عدد نامنفی است که مقادیر کوچکتر بهتر هستند و در حالت کاملاً دقیق به ۰ می‌رسد
- مربع خطای یکی از توابع پرکاربرد در مسئله رگرسیون است



$$l^{(i)}(\mathbf{w}, b) = (\hat{y}^{(i)} - y^{(i)})^2$$

$$L(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n l^{(i)}(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}^{(i)} + b - y^{(i)})^2$$

$$\mathbf{w}^*, b^* = \arg \min_{\mathbf{w}, b} L(\mathbf{w}, b)$$



# بهینه‌سازی

$$\mathbf{w}^*, b^* = \arg \min_{\mathbf{w}, b} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (\mathbf{w}^T \mathbf{x}^{(i)} + b - y^{(i)})^2$$

- برای سادگی می‌توانیم  $b$  را به  $\mathbf{w}$  و یک مقدار ۱ را به هر  $\mathbf{x}^{(i)}$  الحاق کنیم

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

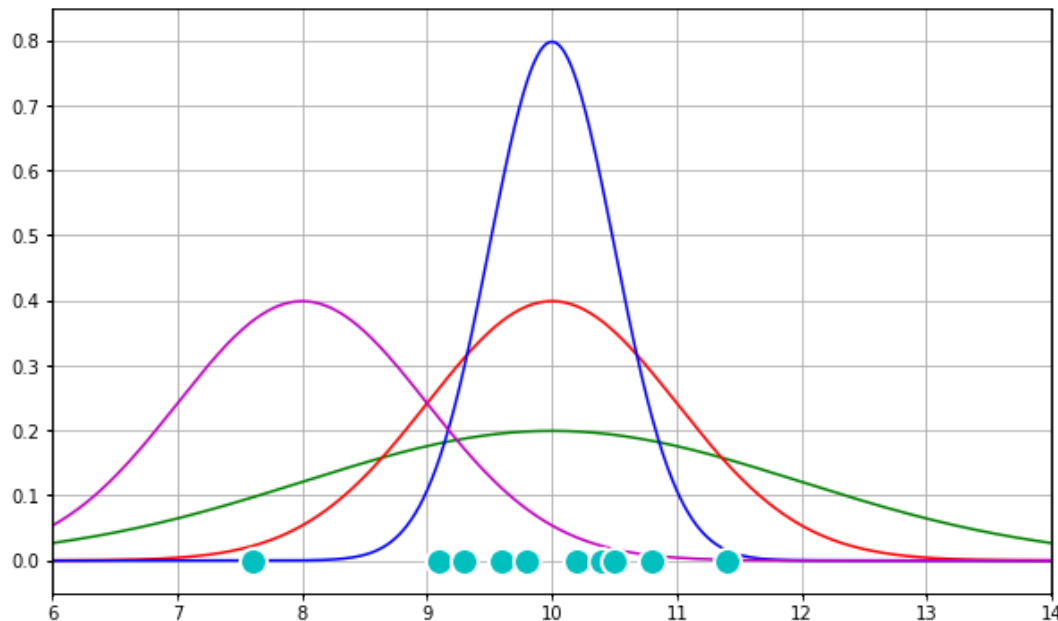
- این مسئله یک مینیمم محلی دارد که مشتق در آن صفر است و به سادگی قابل محاسبه است

$$\nabla_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 = 2\mathbf{X}^T(\mathbf{y} - \mathbf{X}\mathbf{w})$$

$$\mathbf{X}^T \mathbf{X} \mathbf{w}^* = \mathbf{X}^T \mathbf{y} \Rightarrow \mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{X}^\dagger \mathbf{y}$$

# برآوردگر Maximum Likelihood

$$\max \prod_{i=1}^n \mathcal{N}(x^{(i)} | \mu, \sigma^2)$$



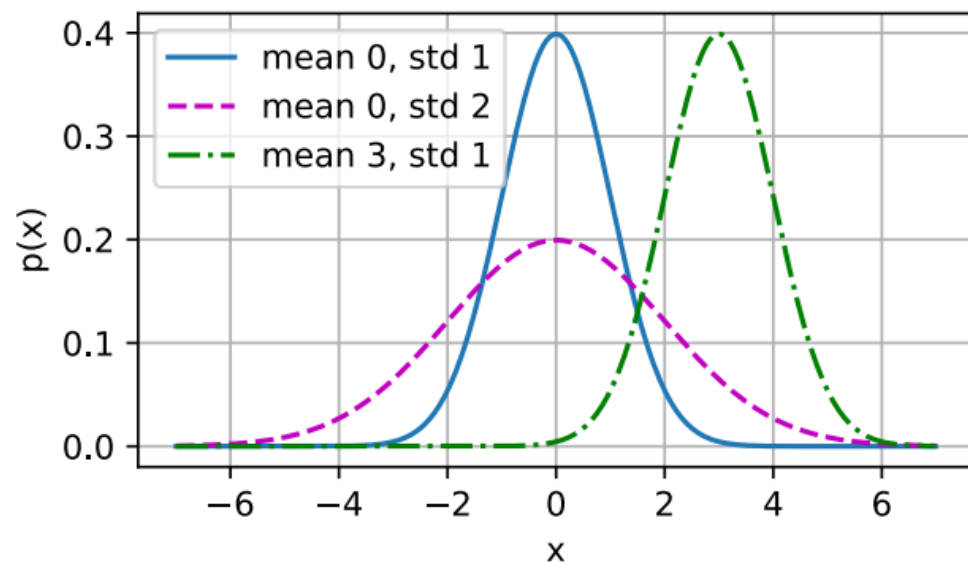
- فرض کنید ۱۰ داده از یک فرآیند مشاهده شده است
  - به عنوان نمونه، زمان پاسخ دادن به یک سوال بر حسب ثانیه
- چگونه این فرآیند را مدل کنیم؟
  - یک راه این است که یک مدل برای توزیع آماری فرآیند فرض کنیم
  - سپس، پارامترهای مدل را به گونه‌ای تعیین کنیم که تابع درست‌نمایی مربوط به مشاهده این نمونه‌ها از آن توزیع را بیشینه کند

# رگرسیون خطی

- فرض کنید مشاهدات از یک مدل خطی همراه با نویز نرمال (گوسی) جمع‌آوری شده‌اند

$$\hat{y} = \mathbf{w}^T \mathbf{x} + b + \epsilon$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2)$$



- توزیع نرمال

$$p(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (x - \mu)^2\right)$$

# رگرسیون خطی

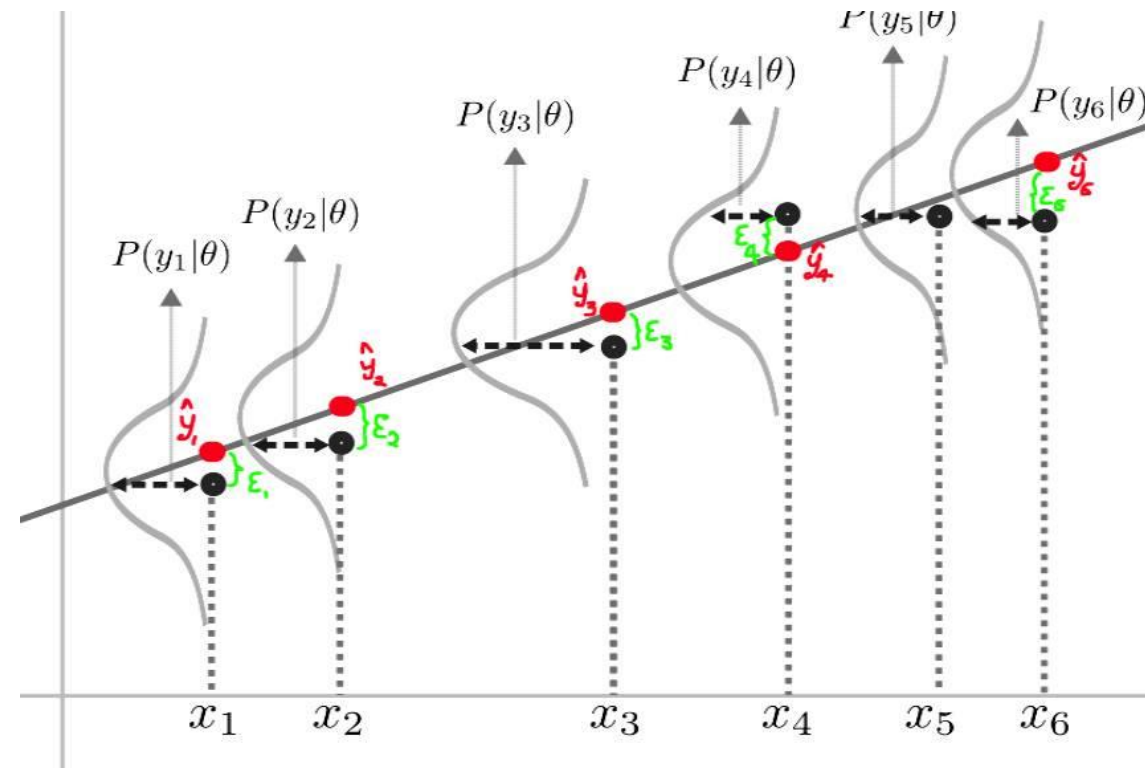
- فرض کنید مشاهدات از یک مدل خطی همراه با نویز نرمال (گوسی) جمع‌آوری شده‌اند

$$\hat{y} = \mathbf{w}^T \mathbf{x} + b + \epsilon$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$P(y | \mathbf{x}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (y - \mathbf{w}^T \mathbf{x} - b)^2\right)$$

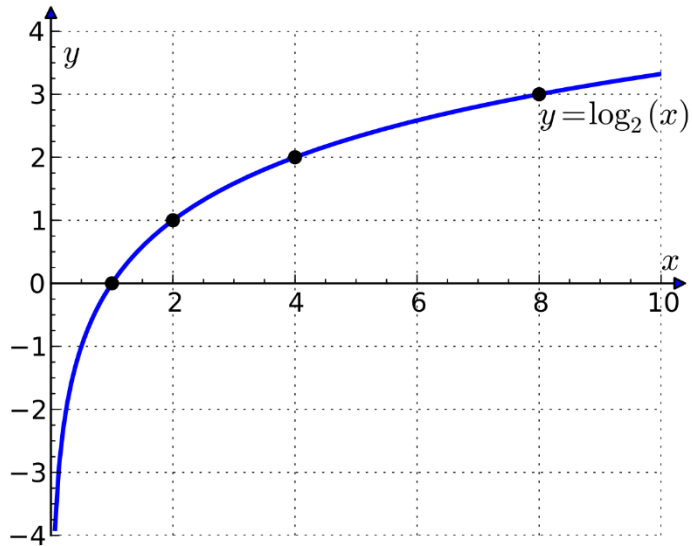
$$P(\mathbf{y} | \mathbf{X}) = \prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)})$$



# Maximum Likelihood

$$P(y | \mathbf{x}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (y - \mathbf{w}^T \mathbf{x} - b)^2\right)$$

$$P(\mathbf{y} | \mathbf{X}) = \prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)})$$



- بجای بیشینه‌سازی  $P(\mathbf{y} | \mathbf{X})$  می‌توان لگاریتم آن را بیشینه کرد - ضرب به جمع تبدیل می‌شود

$$\begin{aligned} \log P(\mathbf{y} | \mathbf{X}) &= \sum_{i=1}^n \log P(y^{(i)} | \mathbf{x}^{(i)}) \\ &= \sum_{i=1}^n -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (y - \mathbf{w}^T \mathbf{x} - b)^2 \end{aligned}$$

# Maximum Likelihood

- به طور معمول منفی لگاریتم تابع درست‌نمایی را کمینه می‌کنیم

$$-\log P(\mathbf{y} | \mathbf{X}) = \sum_{i=1}^n \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} (y - \mathbf{w}^T \mathbf{x} - b)^2$$

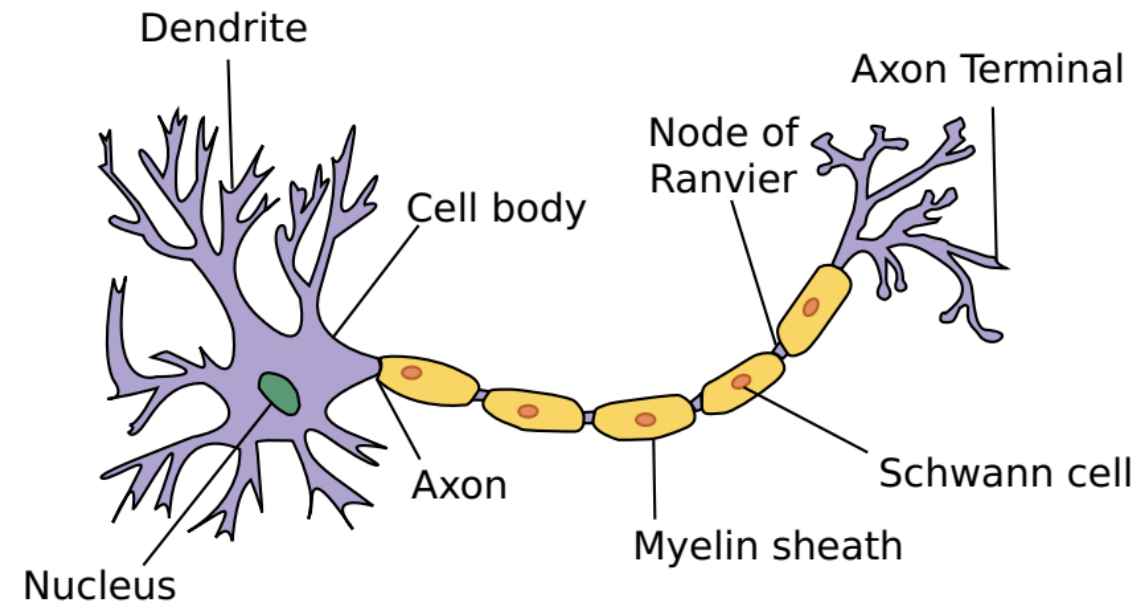
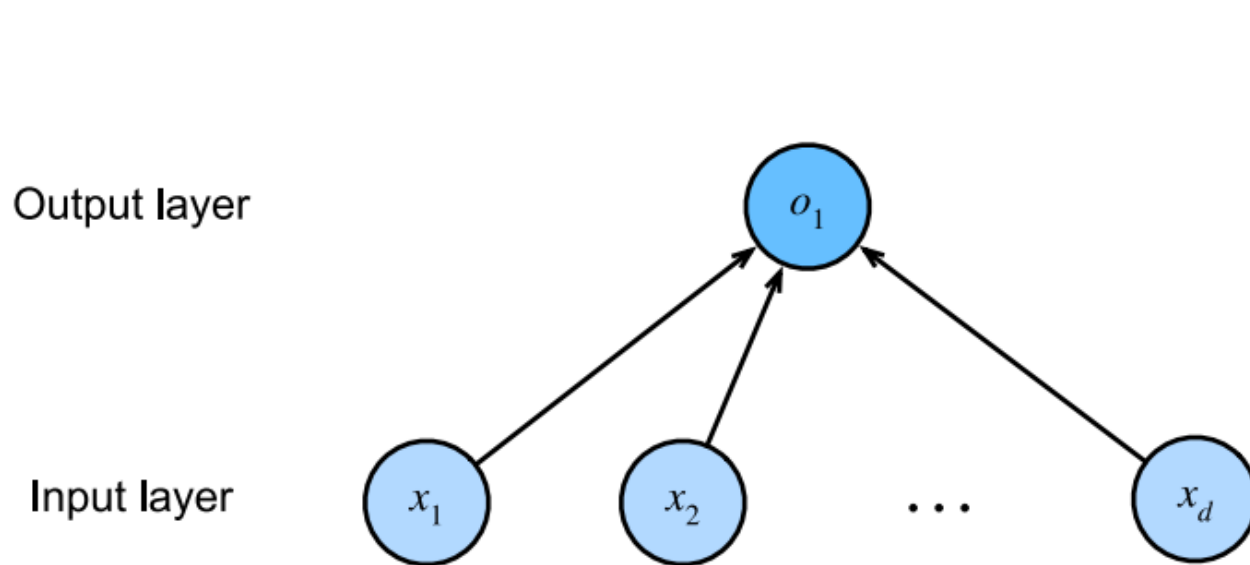
- اگر  $\sigma^2$  را ثابت در نظر بگیریم

$$\min \sum_{i=1}^n (y - \mathbf{w}^T \mathbf{x} - b)^2$$

- همان کمینه‌سازی مجموع مربعات خطا است!
- با استفاده از رویکرد ML می‌توان توابع ضرر متناسب با مسئله را بدست آورد

# شبکه‌های عصبی

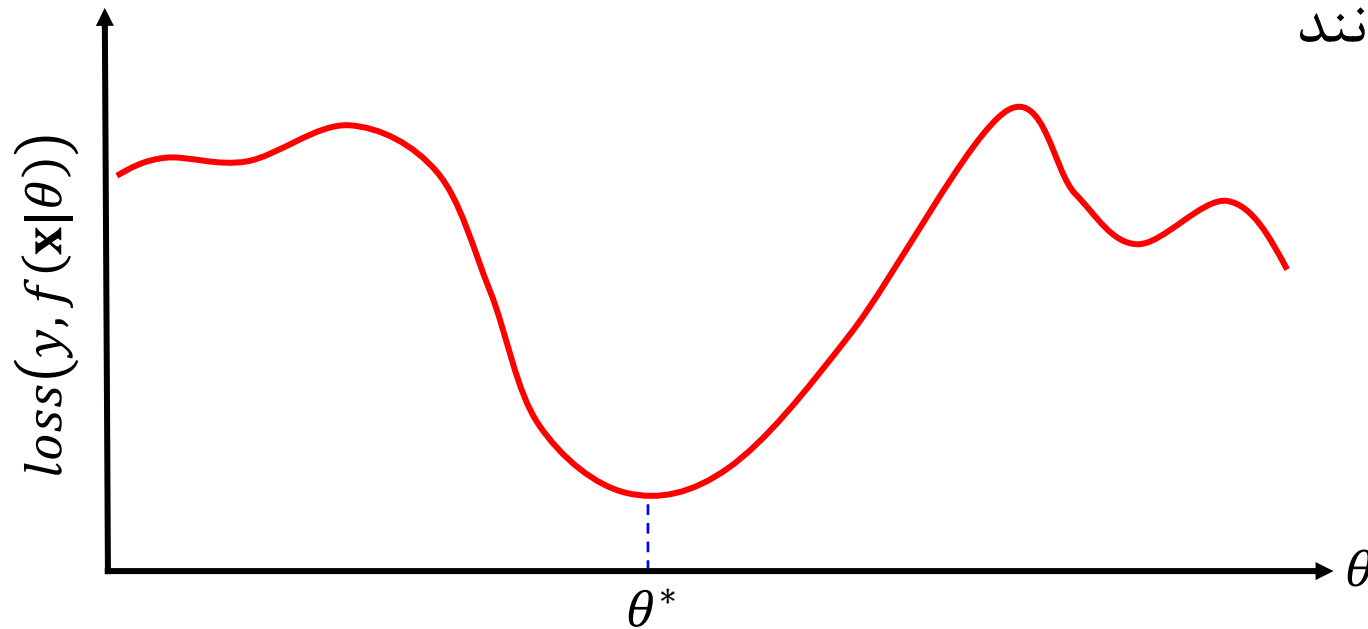
- رگرسیون خطی یک شبکه عصبی یک لایه است
- در این مثال واحد خروجی به تمام واحدهای ورودی متصل است و به آن لایه کاملاً متصل می‌گوئیم
- مدل نورون‌های زیستی از لحاظ ساختاری مشابه با این ساختار است



# بهینه‌سازی

$$\theta^* = \min_{\theta} \text{loss}(y, f(\mathbf{x}|\theta))$$

- در حالت‌هایی که فضای جستجو کوچک باشد می‌توان تمام فضا را جستجو کرد
- اگر فضای جستجو پیوسته اما ساده باشد، می‌توانیم مشتق بگیریم و مساوی با صفر قرار دهیم
- در غیر این صورت، باید از روش‌های تقریبی مانند الگوریتم گرادیان کاهشی استفاده کنیم





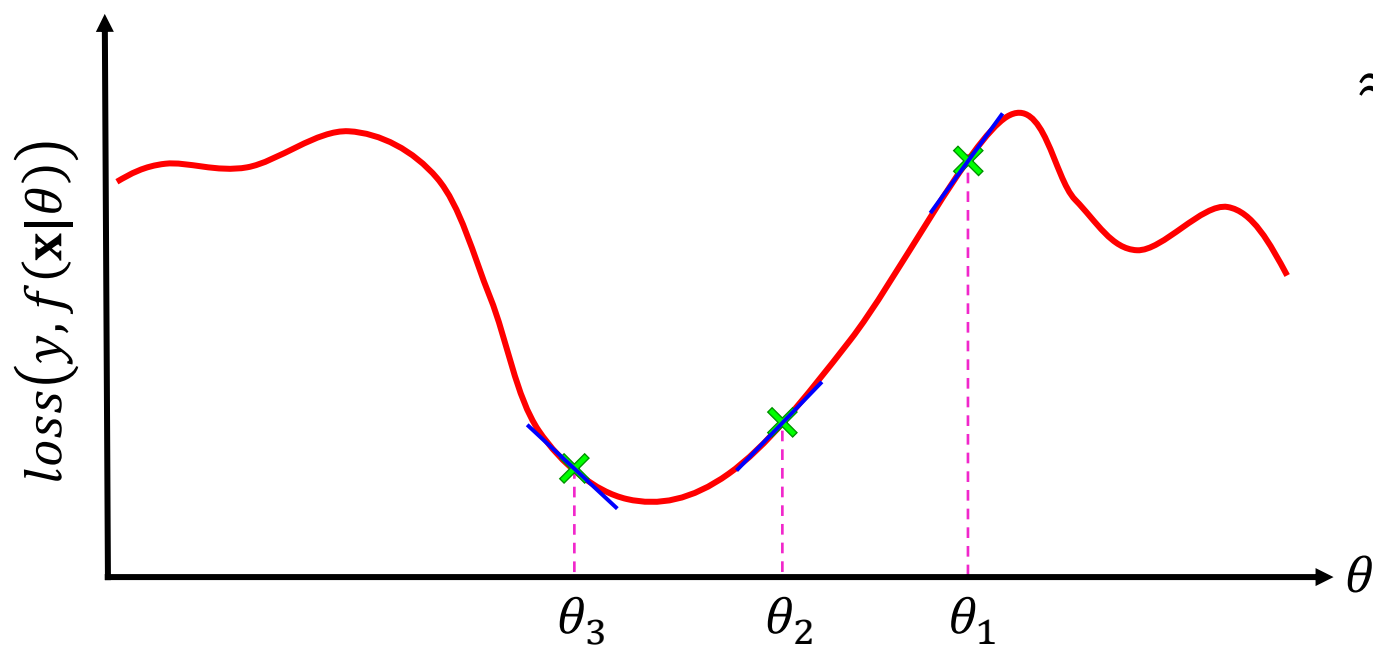
# بهینه‌ساز گرادیان کاهشی

- با یک نقطه اولیه شروع می‌کنیم و در هر گام در جهتی حرکت می‌کنیم که منجر به کاهش تابع شود

$$L(\theta + \Delta\theta) = L(\theta) + \Delta\theta \frac{\partial L(\theta)}{\partial \theta} + \frac{(\Delta\theta)^2}{2!} \frac{\partial^2 L(\theta)}{\partial \theta^2} + \dots$$

$$\approx L(\theta) + \Delta\theta \frac{\partial L(\theta)}{\partial \theta}$$

- در خلاف جهت گرادیان حرکت می‌کنیم



# دسته‌بندی باینری

- در بسیاری از مسائل یادگیری ماشین نیاز به پیش‌بینی مقدار یک متغیر باینری است
- شبکه‌های عصبی باید تنها یک مقدار را پیش‌بینی کنند  $P(y = 1 | \mathbf{x}) \in [0,1]$
- بهتر است خروجی محدود به بازه  $[0,1]$  باشد

$y \in \{0, 1\}$	Two Class Classification	
	1 or Positive Class	0 or Negative Class
Email	Spam	Not Spam
Tumor	Malignant	Benign
Transaction	Fraudulent	Not Fraudulent

e.g.  $P(y = 1 | \mathbf{x}) = \max\{0, \min\{1, \mathbf{w}^T \mathbf{x} + b\}\}$

- بهینه‌سازهای مبتنی بر گرادیان نمی‌توانند پارامترهای مطلوب چنین شبکه‌ای را بیابند
- مشتق تابع ضرر برای داده‌هایی که کاملاً اشتباه پیش‌بینی شوند ۰ است

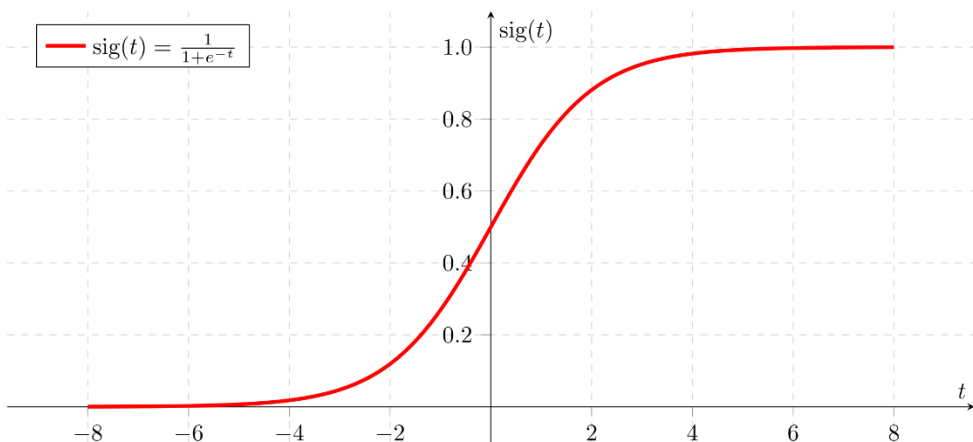
# دسته‌بندی باینری

- بهتر است از تابعی برای محدود کردن خروجی استفاده کنیم که مشتق آن صفر نشود
- تابع سیگموئید برای این منظور پرکاربرد است  $P(y = 1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$
- آیا تابع ضرر MSE در این حالت مناسب است؟

$$l = (y - \sigma(o))^2, \quad o = \mathbf{w}^T \mathbf{x} + b$$

- مثال عددی: اگر  $o = 5.6$  و  $y = 0$

$$\begin{aligned} \frac{dl}{do} &= -2(y - \sigma(o))\sigma(o)(1 - \sigma(o)) = 0.0073 \\ &\approx -2(0 - 1)1(1 - 1) \end{aligned}$$



# Maximum Likelihood

$$-\log P(\mathbf{y} | \mathbf{X}) = -\sum_{i=1}^n \log P(y^{(i)} | \mathbf{x}^{(i)})$$

- $y$  یک متغیر باینری است و احتمال  $P(y | \mathbf{x})$  به صورت زیر قابل بیان است

$$P(y | \mathbf{x}) = \begin{cases} \hat{y} & \text{if } y = 1 \\ 1 - \hat{y} & \text{if } y = 0 \end{cases} = \hat{y}^y (1 - \hat{y})^{1-y} \quad \text{Bernoulli distribution}$$

$$-\log P(\mathbf{y} | \mathbf{X}) = -\sum_{i=1}^n y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

- تابع ضرر binary cross-entropy نامیده می‌شود که برای مسائل دسته‌بندی باینری پرکاربرد است

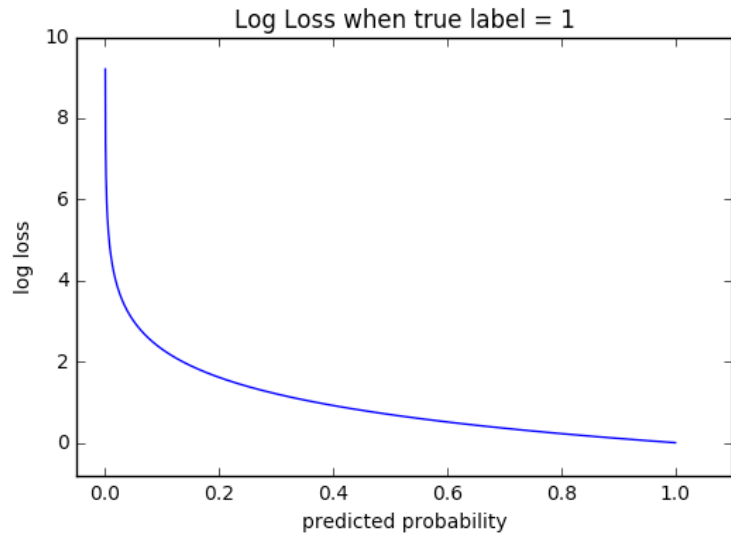
# دسته‌بندی باینری

- بهتر است از تابعی برای محدود کردن خروجی استفاده کنیم که مشتق آن صفر نشود
- تابع سیگموئید برای این منظور پرکاربرد است  $P(y = 1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$
- با استفاده از ML تابع ضرر binary cross-entropy بدست می‌آید

$$l = -y \log \sigma(o) - (1 - y) \log(1 - \sigma(o))$$

- مثال عددی: اگر  $o = 5.6$  و  $y = 0$

$$\frac{dl}{do} = - \frac{-\sigma(o)(1 - \sigma(o))}{1 - \sigma(o)} = \sigma(o) = 0.9963$$



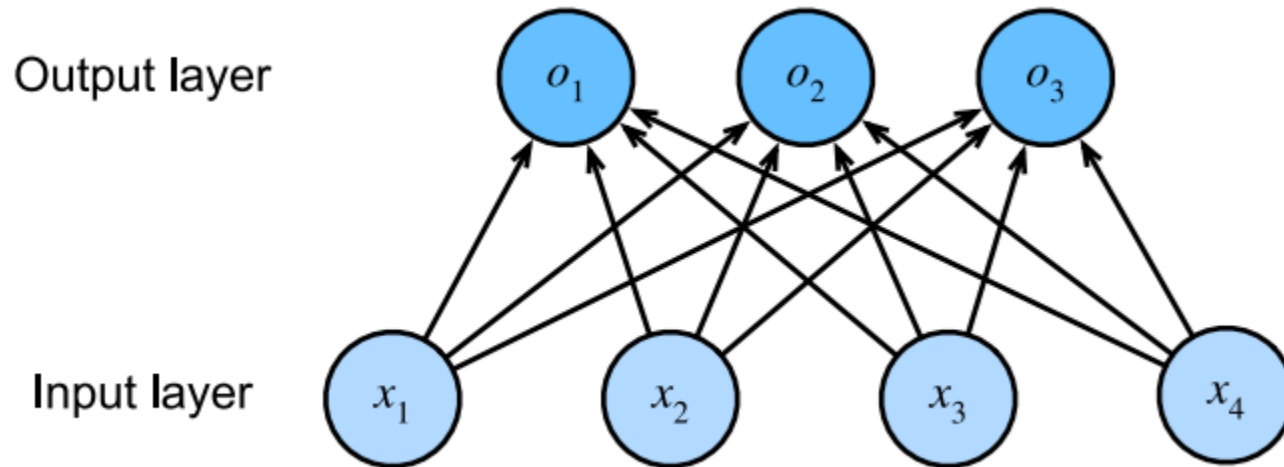
# دسته‌بندی چند کلاسه

$$o_1 = w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + w_{14}x_4 + b_1$$

$$o_2 = w_{21}x_1 + w_{22}x_2 + w_{23}x_3 + w_{24}x_4 + b_2$$

$$o_3 = w_{31}x_1 + w_{32}x_2 + w_{33}x_3 + w_{34}x_4 + b_3$$

- تعیین دسته یک نمونه جدید از میان چندین دسته دارای کاربردهای بسیار فراوانی است
- می‌توان برای هر کلاس یک نرون قرار داد تا احتمال تعلق نمونه به آن کلاس را تخمین بزند
- با استفاده از یک لایه خطی احتمال غیرنرمالیزه  $o$  را پیش‌بینی می‌کنیم که **logit** نامیده می‌شوند
- چطور نرمالیزه کنیم؟



## دسته‌بندی چند کلاسه

- می‌خواهیم احتمال پسین مربوط به هر کلاس را تخمین می‌زنیم  $P(y = i | \mathbf{x}) \in [0,1]$
- برای آنکه مقادیر خروجی از جنس احتمال باشند (هر کدام نامنفی و مجموع برابر با ۱) می‌توانیم از تابع فعال‌سازی Softmax استفاده کنیم که تعمیم تابع Sigmoid است

$$\mathbf{o} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$$

$$\hat{y}_i = \text{softmax}(\mathbf{o})_i = \frac{\exp(o_i)}{\sum_{k=1}^q \exp(o_k)}$$

- تابع ضرر متناسب با این تابع فعال‌سازی، categorical cross-entropy است

$$l(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^q y_i \log \hat{y}_i$$

## مشتق softmax

$$\begin{aligned}l(\mathbf{y}, \hat{\mathbf{y}}) &= - \sum_{i=1}^q y_i \log \frac{\exp(o_i)}{\sum_{k=1}^q \exp(o_k)} = \sum_{i=1}^q y_i \log \sum_{k=1}^q \exp(o_k) - \sum_{i=1}^q y_i o_i \\&= \log \sum_{k=1}^q \exp(o_k) \sum_{i=1}^q y_i - \sum_{i=1}^q y_i o_i = \log \sum_{k=1}^q \exp(o_k) - \sum_{i=1}^q y_i o_i\end{aligned}$$

$$\partial_{o_i} l(\mathbf{y}, \hat{\mathbf{y}}) = \frac{\exp(o_i)}{\sum_{k=1}^q \exp(o_k)} - y_i = \text{softmax}(\mathbf{o})_i - y_i$$

- مشابه با رگرسیون خطی، گرادیان تابع ضرر نسبت به خروجی بخش خطی برابر با میزان اختلاف پیش‌بینی با مقدار واقعی است



# پرسپترون چندلایه

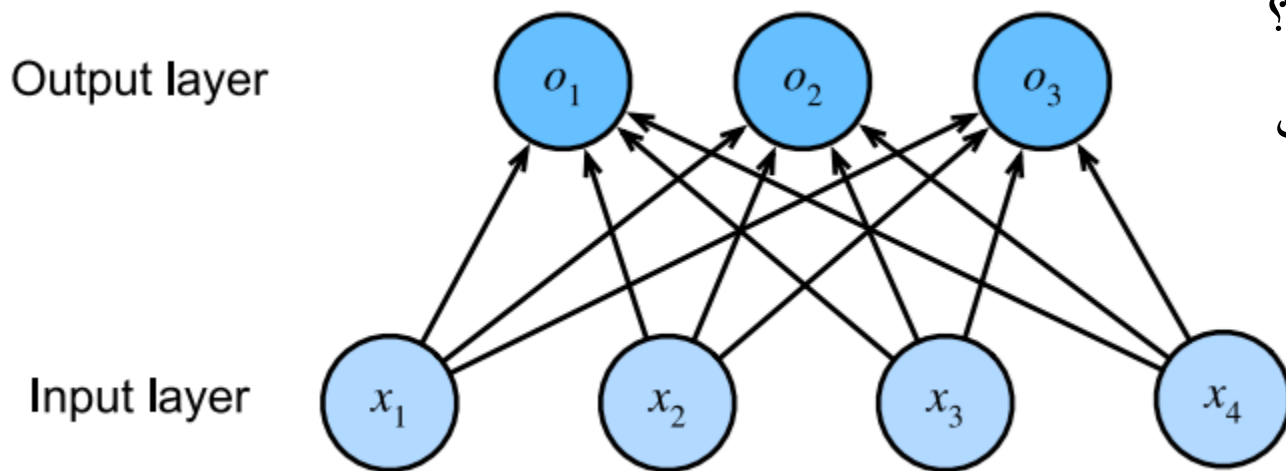
Multilayer Perceptron

# مدل‌های خطی

- بر فرض یکنواختی دلالت دارد

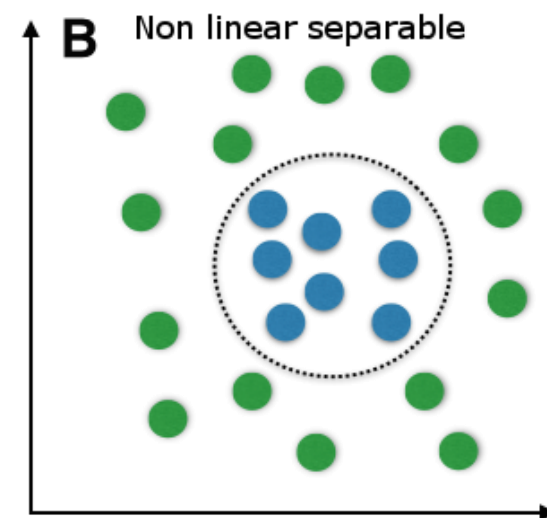
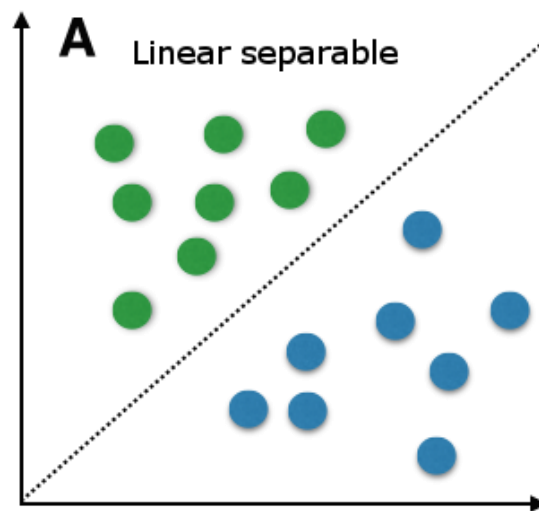
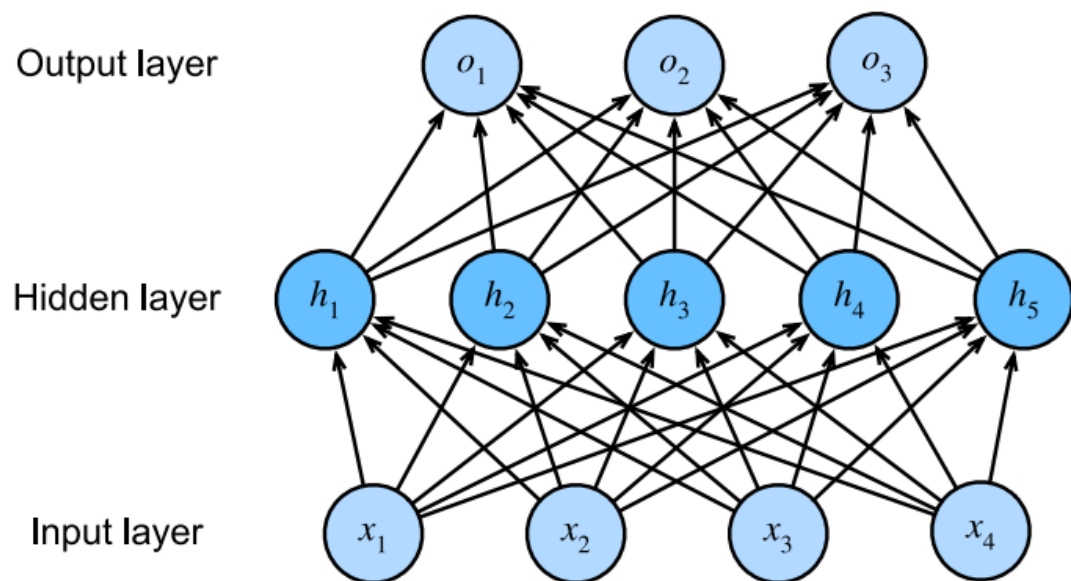
- هر افزایشی در ویژگی همواره منجر به افزایش (کاهش) خروجی می‌شود اگر وزن مربوطه + (-) باشد
- مثال: احتمال بازپرداخت اقساط بر حسب ویژگی‌هایی شامل درآمد
  - اثر افزایش درآمد از ۰ به ۵ میلیون متناظر با افزایش از ۱۰۰ به ۱۰۵ میلیون است
- مثال: احتمال مرگ بر حسب ویژگی‌هایی شامل دمای بدن
  - هر چه دما بالاتر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

- استخراج ویژگی قبل از لایه خطی بسیار مهم است



# پرسپترون چند لایه

- 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

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

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

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

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

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

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

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

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

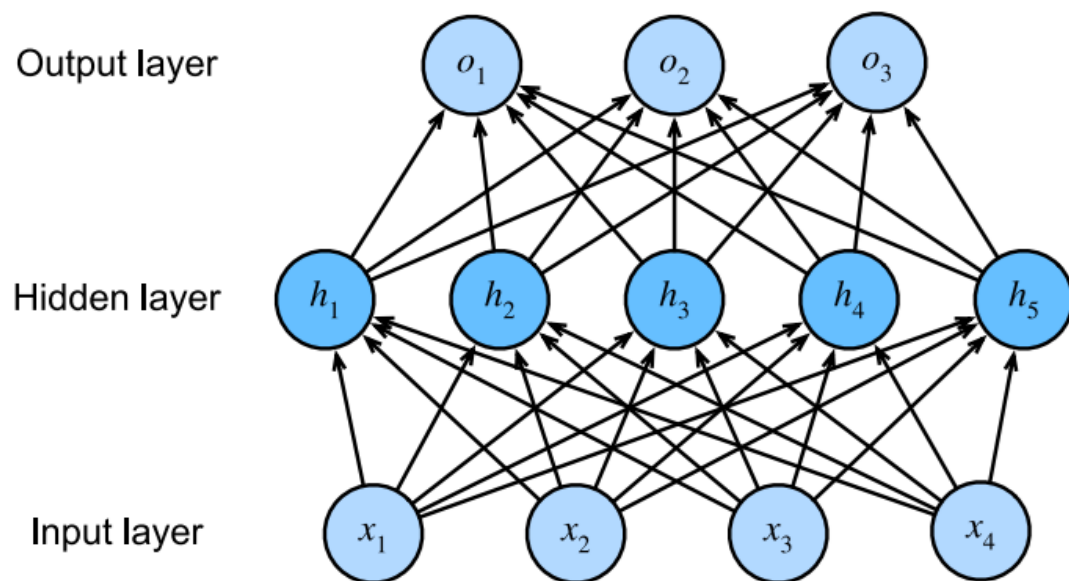
$$\mathbf{b}^{(2)} \in \mathbb{R}^{1 \times q}$$

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

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

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

• 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:

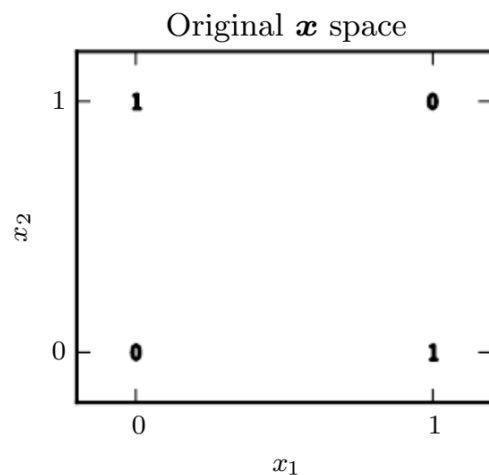
$$\hat{y}^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)} + b$$

• مدل خطی:

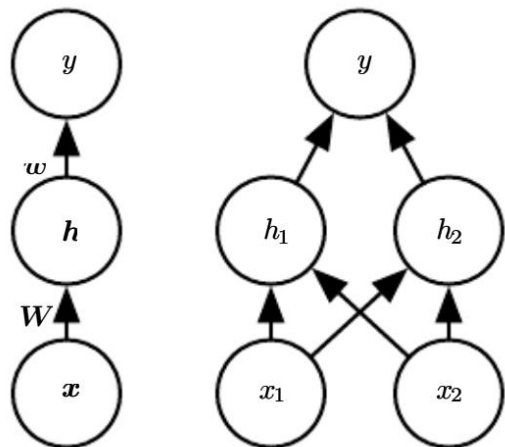
• با استفاده از بهینه‌سازی تحلیلی به پاسخ زیر می‌رسیم

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

-  $L = 0.25$  خواهد شد



# مثال



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

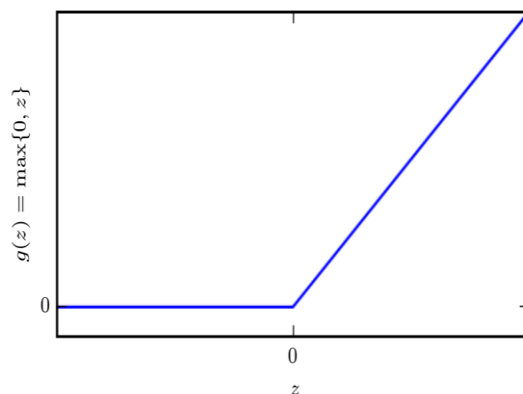
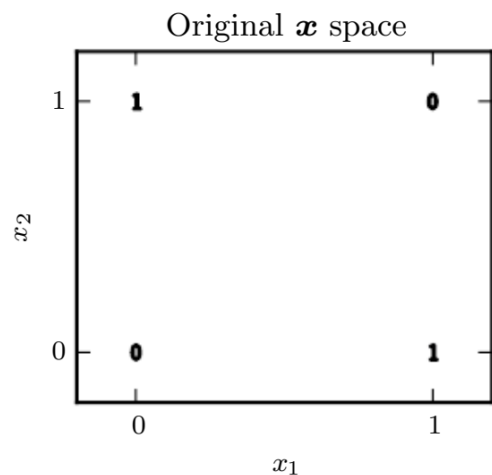
• تابع ضرر MSE:

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

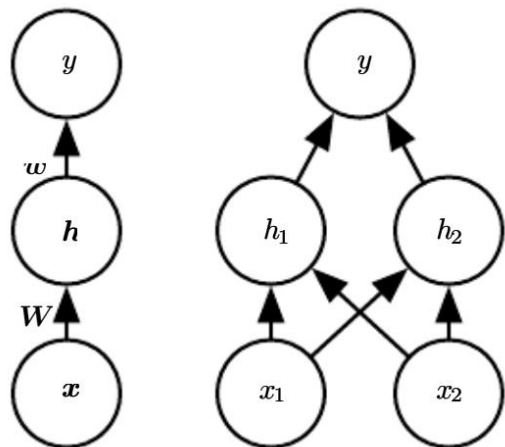
• مدل ۲ لایه:

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

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



# مثال

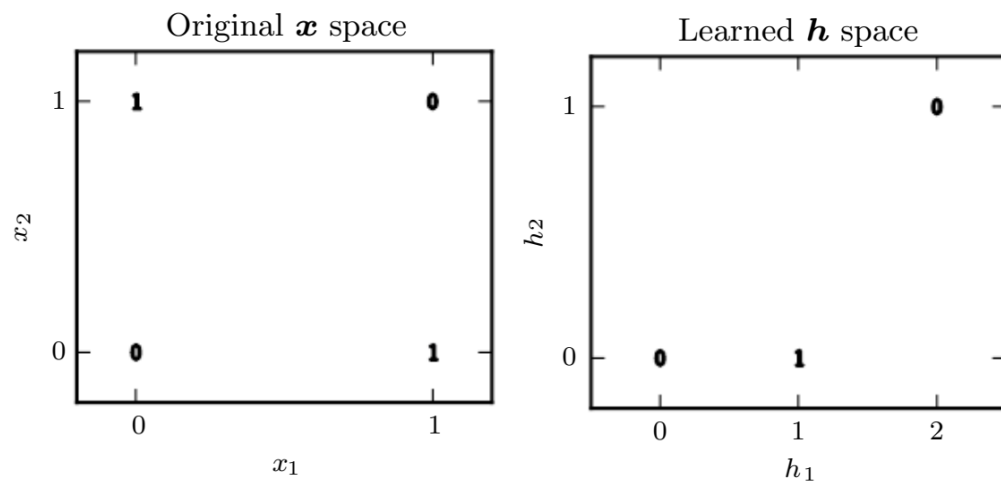


$$L = \frac{1}{4} \sum_{i=1}^4 (y^{(i)} - \hat{y}^{(i)})^2$$

• 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

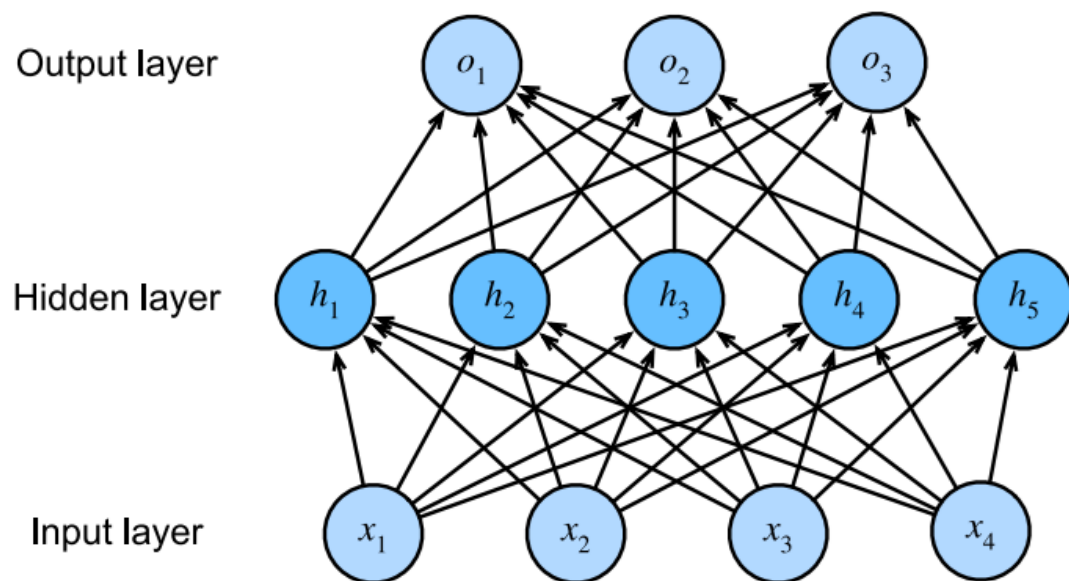
- یک شبکه MLP با تابع فعال‌سازی مناسب در لایه‌های میانی و تعداد نورون کافی می‌تواند هر تابعی را پیاده‌سازی کند

- اصطلاحاً Universal Approximator است

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

- قسمت سخت کار آموزش وزن‌های مدل است

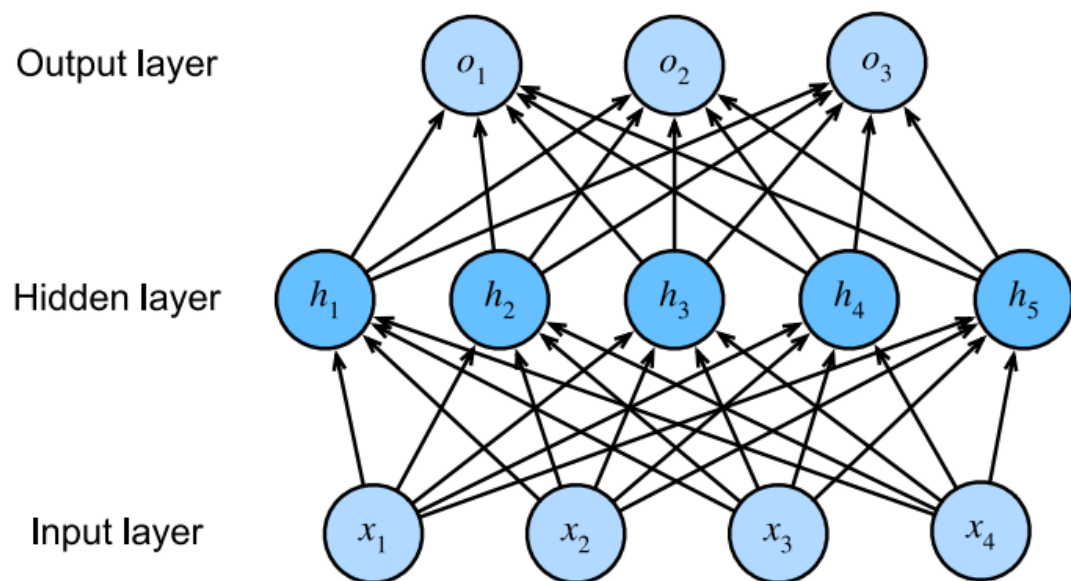
- برای تقریب بسیاری از توابع می‌توانیم از شبکه‌های عمیق‌تر (بجای عریض‌تر) استفاده کنیم که بسیار فشرده‌تر و کاراتر باشد





# توابع فعال سازی

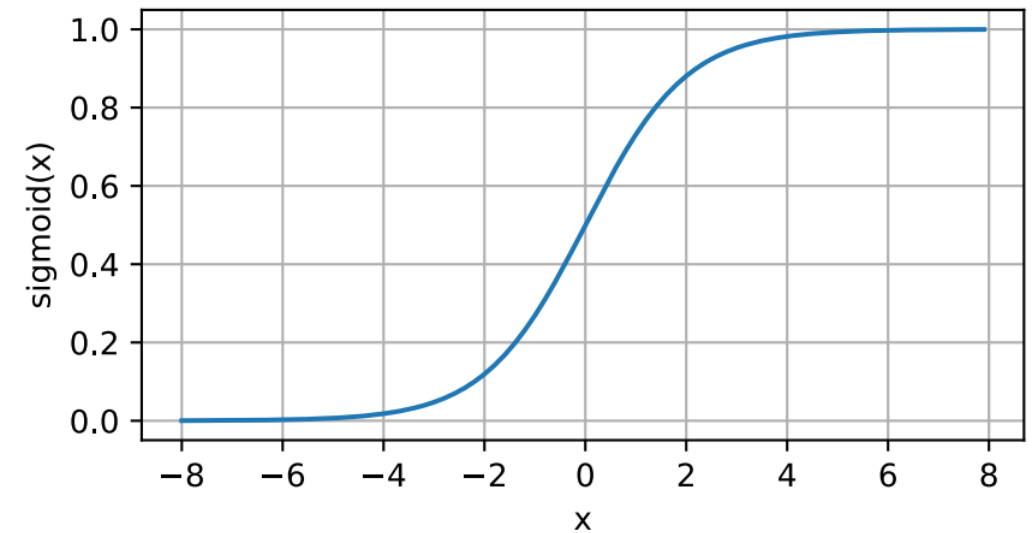
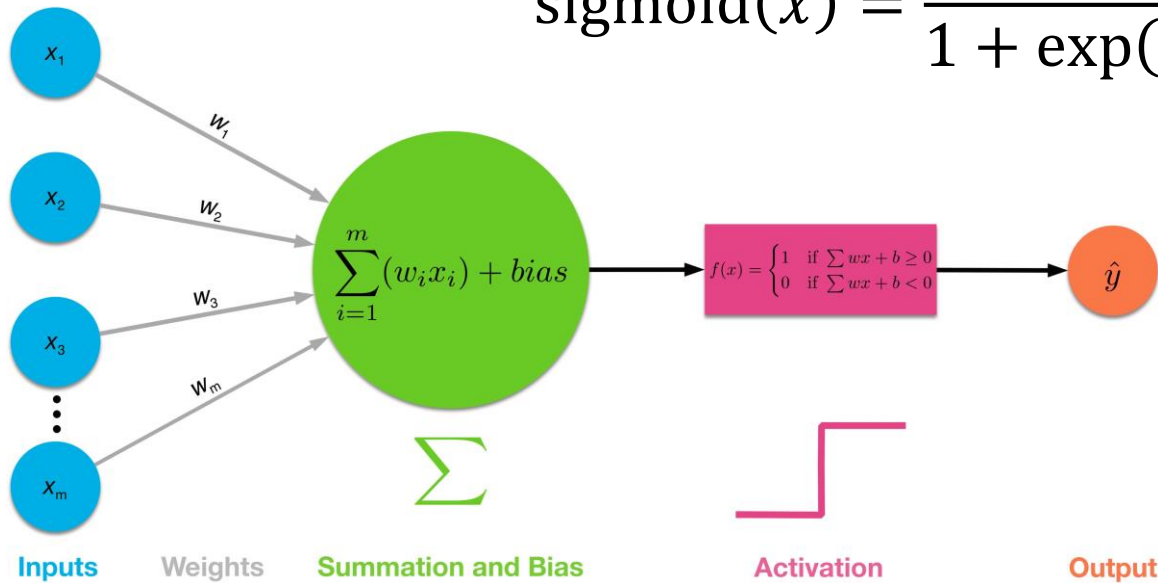
- توابع فعال سازی برای هر نورون مشخص می کنند که خروجی بخشی خطی در چه شرایطی منجر به فعال شدن نورون شود
- به خصوص در شبکه های عمیق بسیار با اهمیت هستند



# 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

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

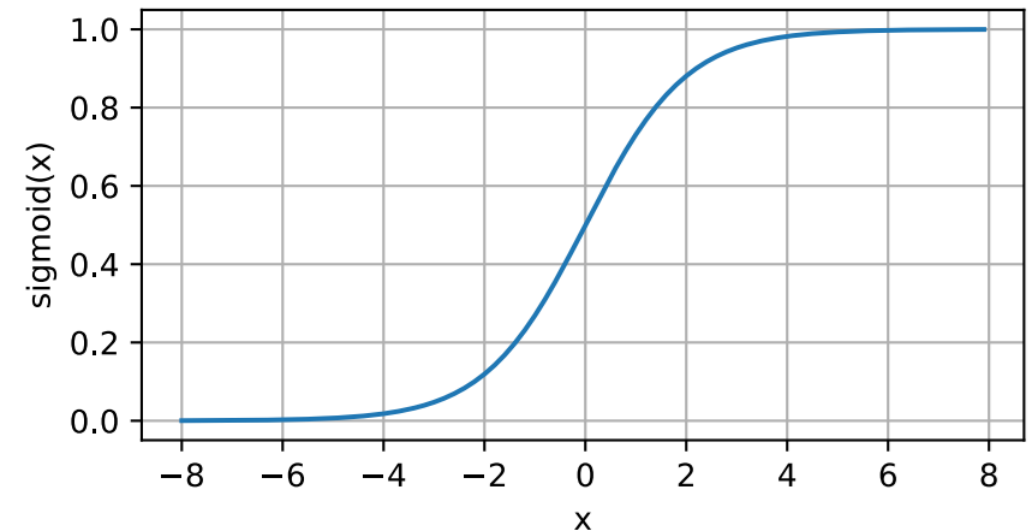
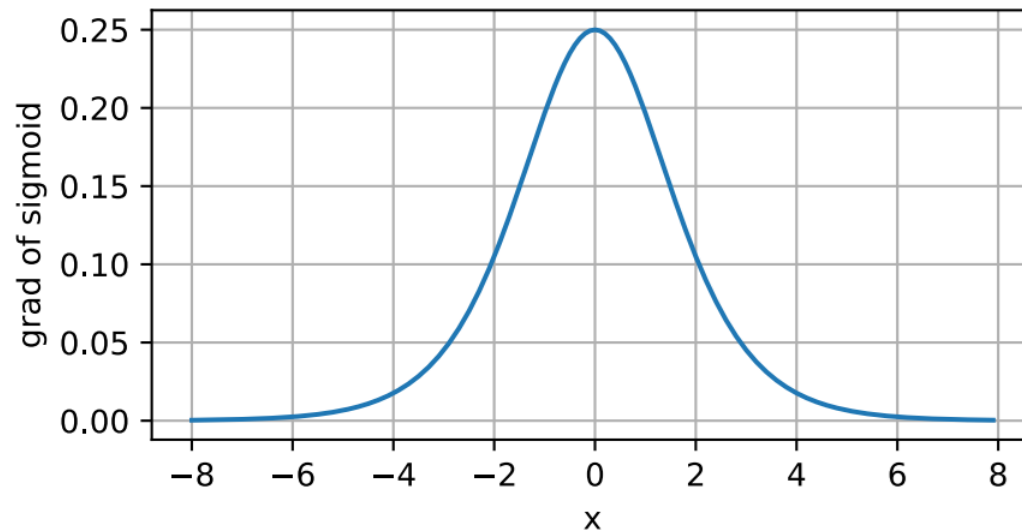


# Sigmoid

- 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

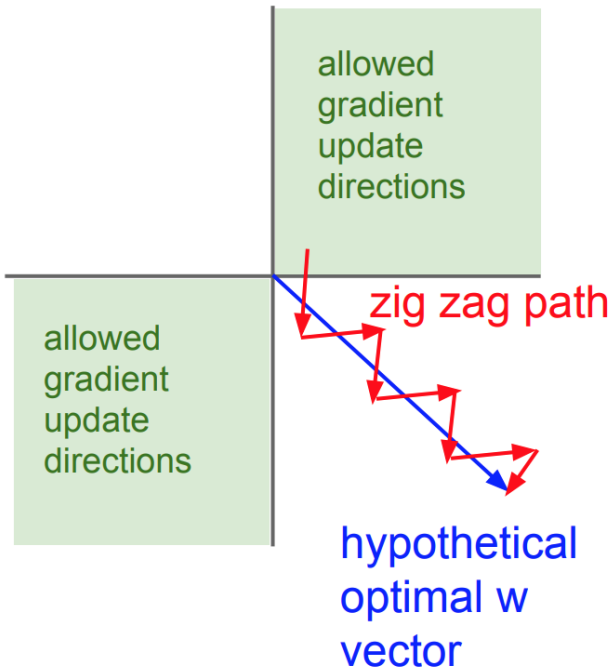
$$z = \sum_k w_k h_k + b$$

- 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

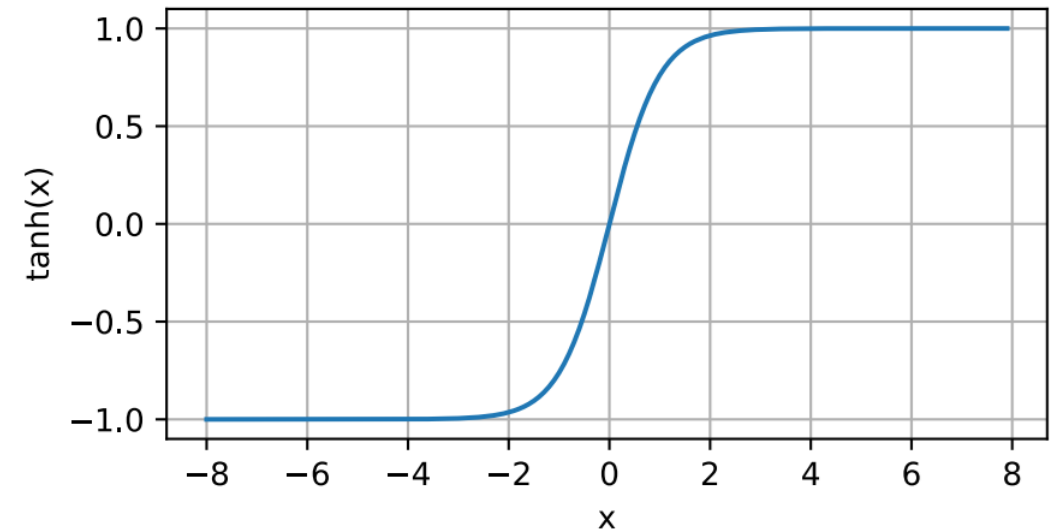
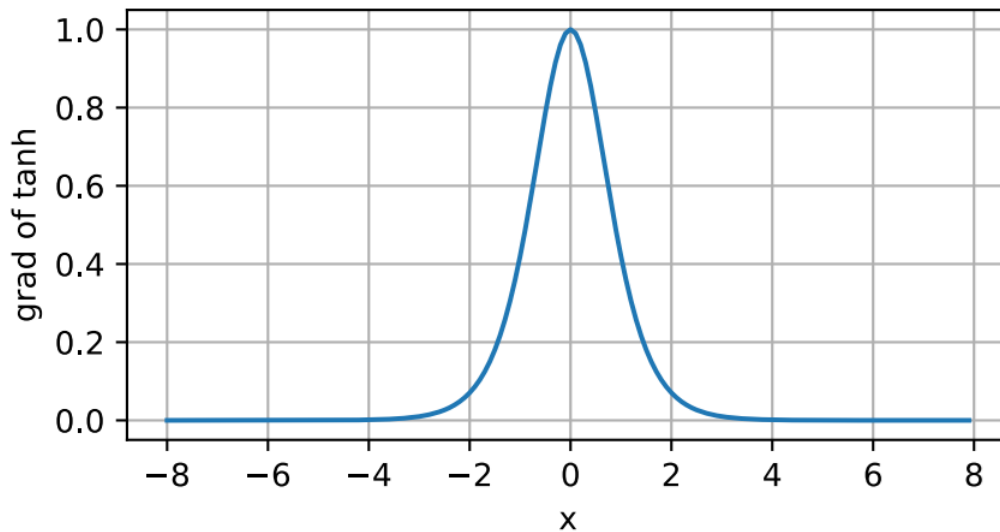


# Tanh

- 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

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

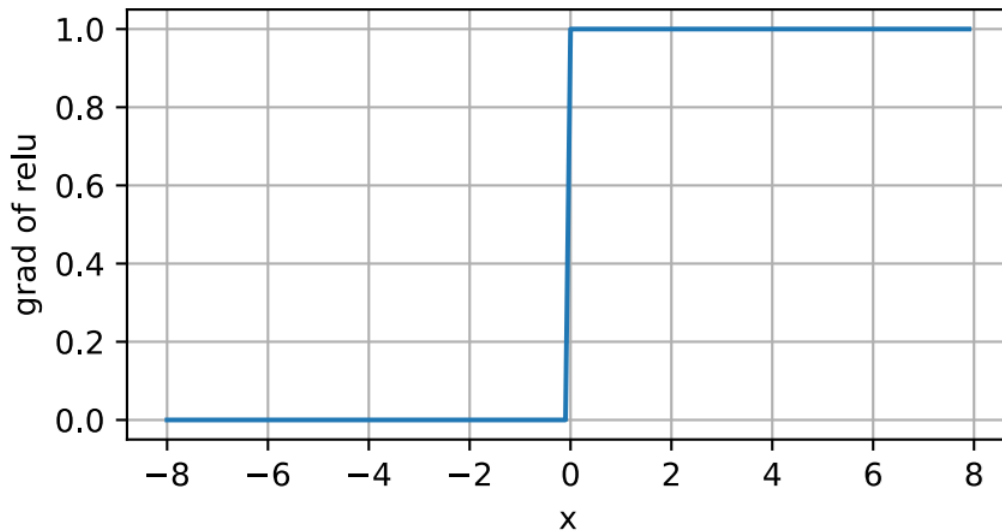
$$\tanh(x) = \frac{1 - \exp(-2x)}{1 + \exp(-2x)} = 2\sigma(2x) - 1$$



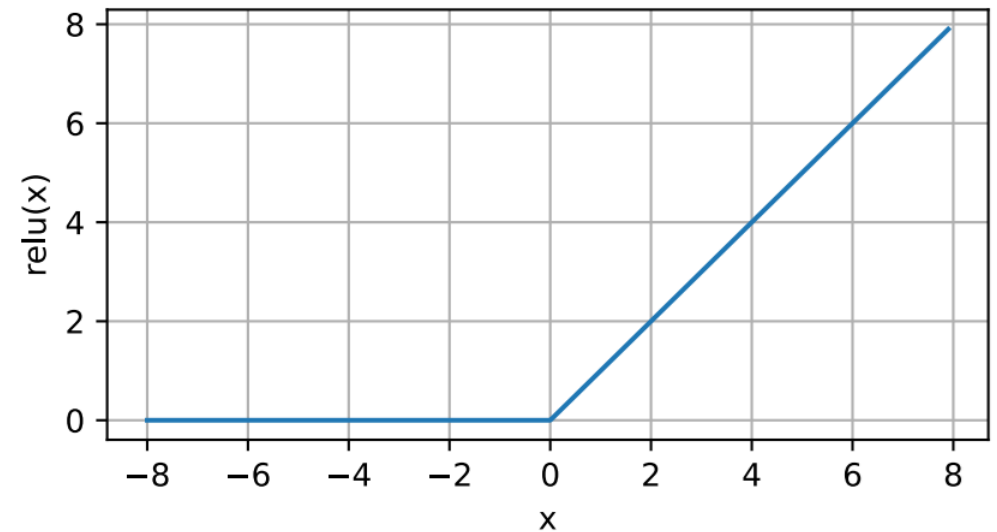
# Rectified Linear Unit (ReLU)

- 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

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



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



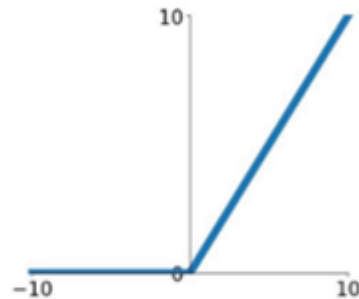
# ReLU با شیب مخالف صفر

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)  $\alpha_i$  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**  
 $\max(0.1x, x)$

