# Introduction to Deep Learning

## 7. Multilayer Perceptron

**MGMT 735**

**Slides From Mu Li and Alex Smola**
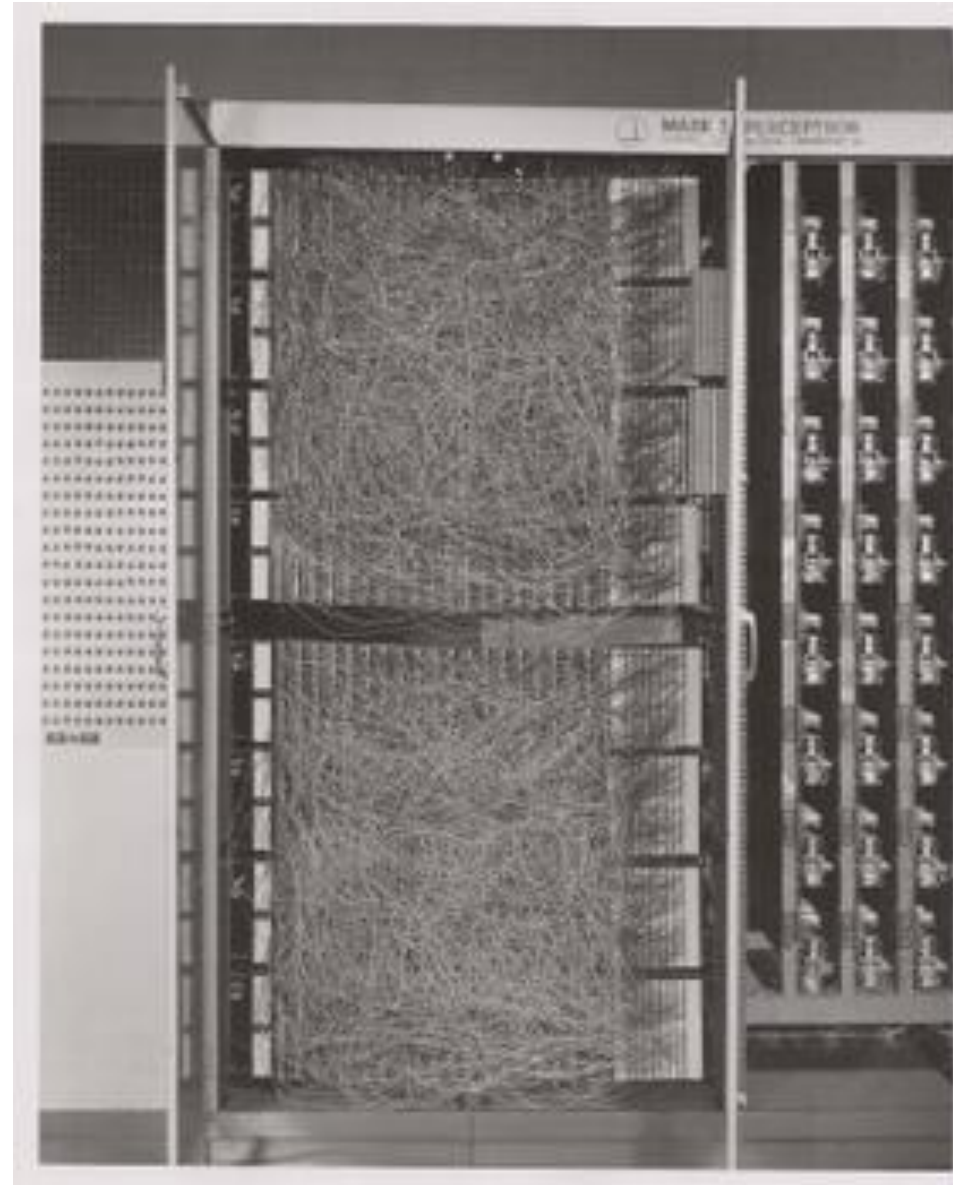
**courses.d2l.ai/berkeley-stat-157**

aws

# Outline

- **Single Layer Perceptron**
  - Decision Boundary
  - XOR is hard
- **Multilayer Perceptron**
  - Layers
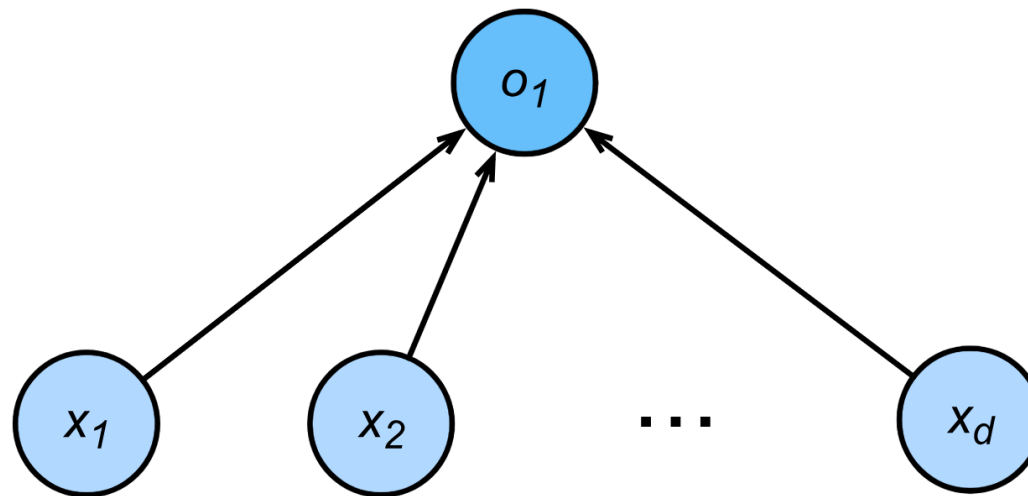  - Nonlinearities
  - Computational Cost

# Perceptron

Mark I Perceptron, 1960
([wikipedia.org](wikipedia.org))

# Perceptron

- Given input $\mathbf{x}$, weight $\mathbf{w}$ and bias $b$, perceptron outputs:

$$o = \sigma\left(\langle \mathbf{w}, \mathbf{x} \rangle + b\right) \qquad \sigma(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$
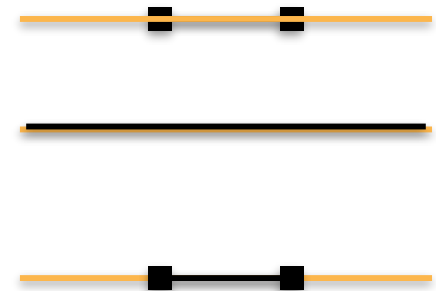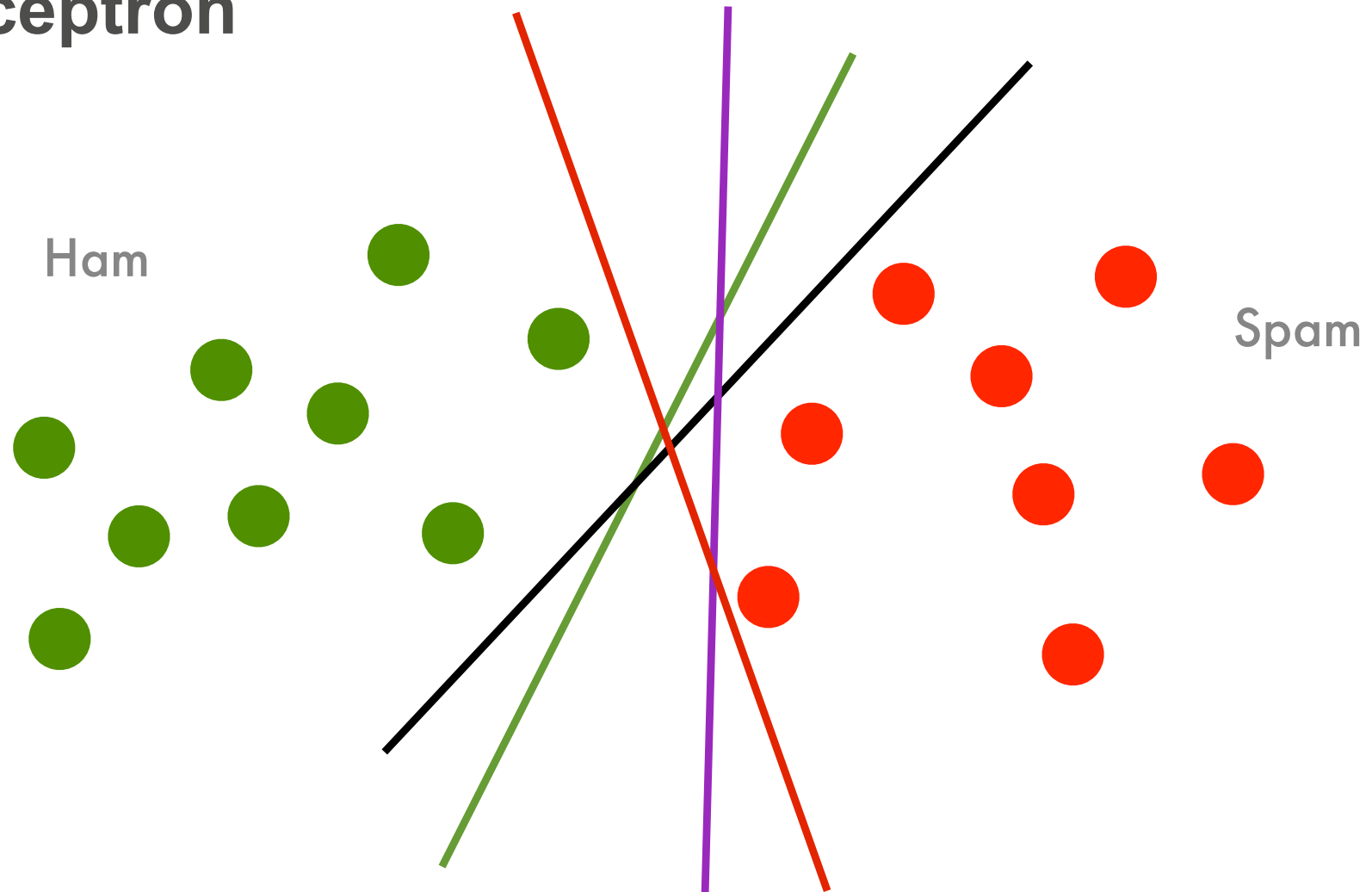
# Perceptron

- Given input $\mathbf{x}$, weight $\mathbf{w}$ and bias $b$, perceptron outputs:

$$o = \sigma\left(\langle \mathbf{w}, \mathbf{x} \rangle + b\right) \qquad \sigma(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

- Binary classification (0 or 1)
  - Vs. scalar real value for regression
  - Vs. probabilities for logistic regression

# Perceptron

Ham

Spam

**Training the Perceptron**

$\textbf{initialize } w = 0 \text{ and } b = 0$

$\textbf{repeat}$

   $\textbf{if } y_i \left[ \langle w, x_i \rangle + b \right] \leq 0 \textbf{ then}$

      $w \leftarrow w + y_i x_i \text{ and } b \leftarrow b + y_i$

   $\textbf{end if}$
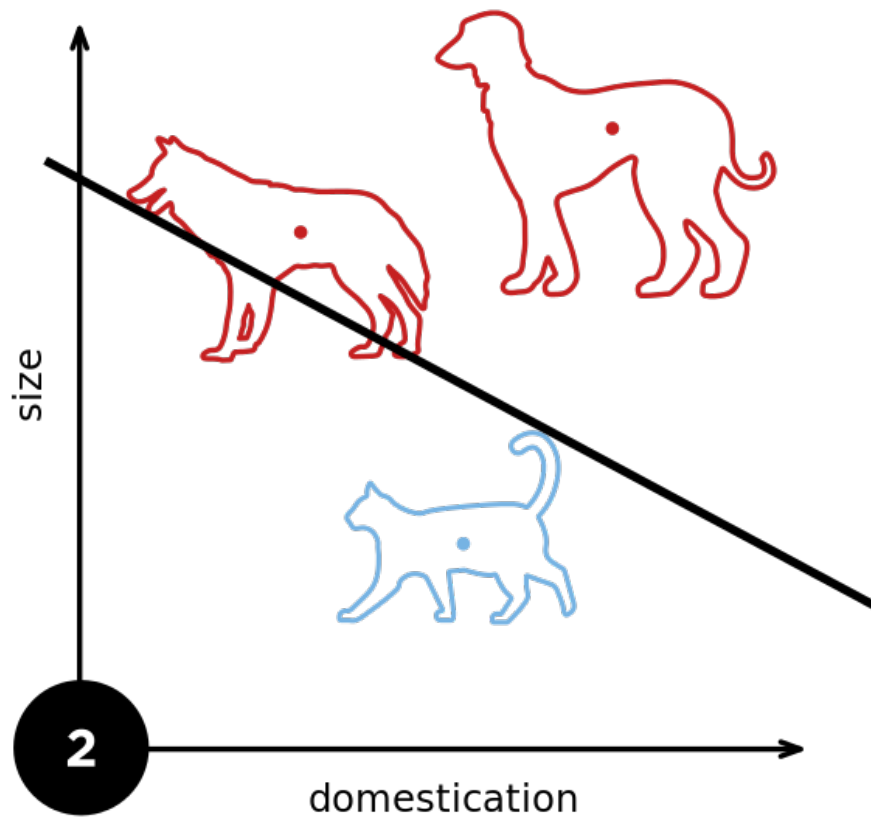
$\textbf{until } \text{all classified correctly}$

Equals to SGD (batch size is 1) with the following loss

$$\ell(y, \mathbf{x}, \mathbf{w}) = \max(0, -y\langle \mathbf{w}, \mathbf{x} \rangle)$$
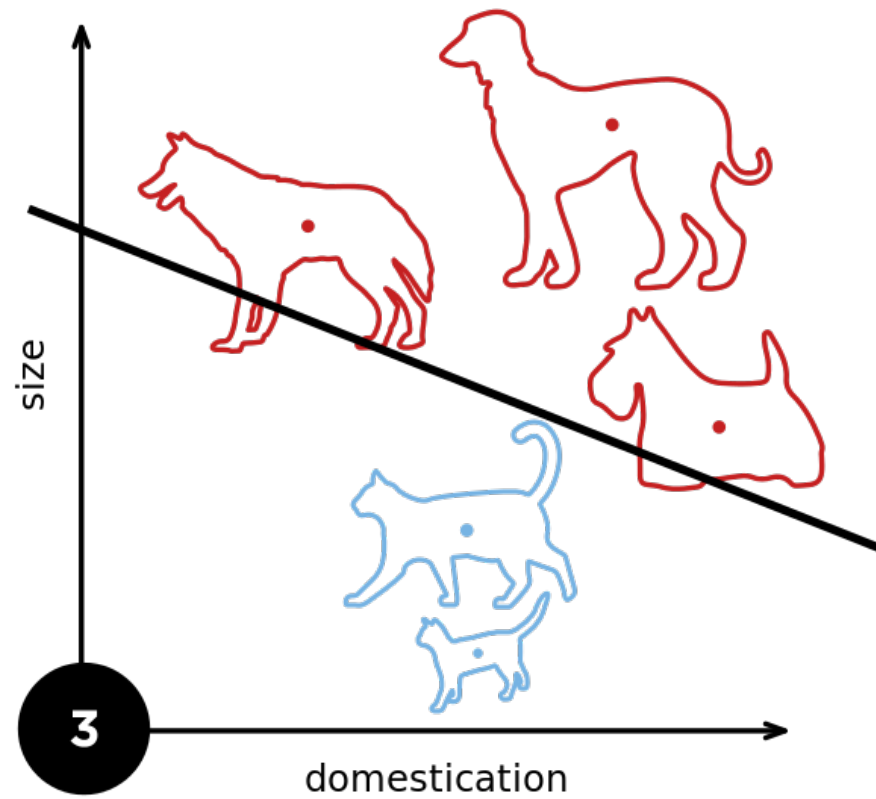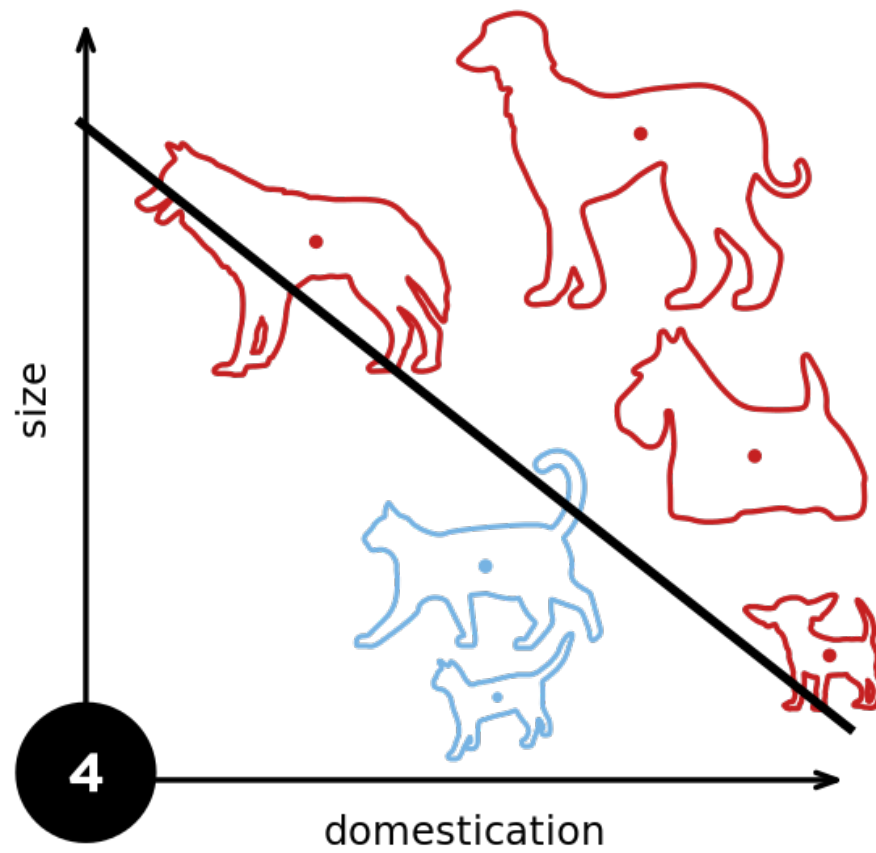
size

domestication

1

From wikipedia

size

domestication

**2**

From wikipedia

From wikipedia

size
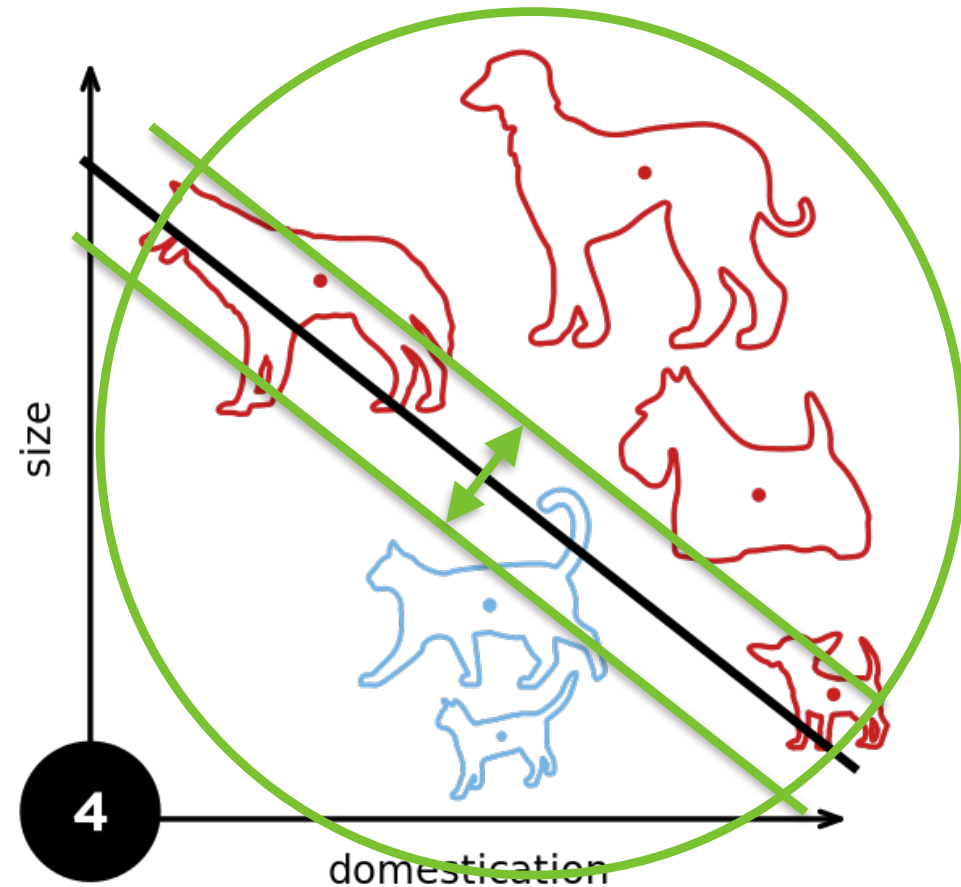
domestication

4

From wikipedia

# Convergence Theorem

- Radius r enclosing the data
- Margin $\rho$ separating the classes

$$y(\mathbf{x}^\top \mathbf{w} + b) \geq \rho$$

for $\|\mathbf{w}\|^2 + b^2 \leq 1$

- Guaranteed that perceptron will converge after

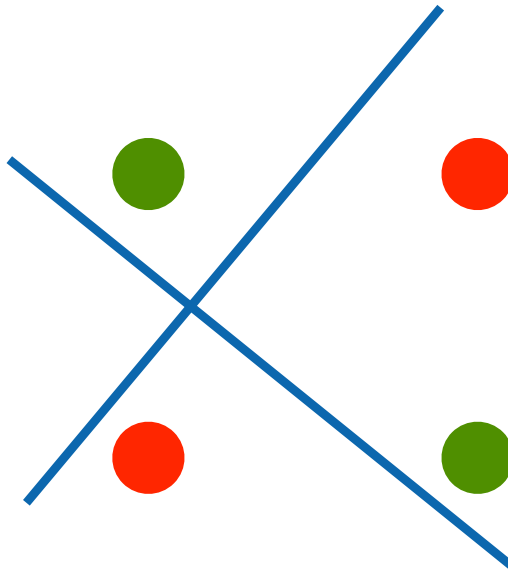$$\frac{r^2 + 1}{\rho^2} \text{ steps}$$

# Consequences

- Only need to store errors. This gives a compression bound for perceptron.

- Fails with noisy data

do NOT train your avatar with perceptrons
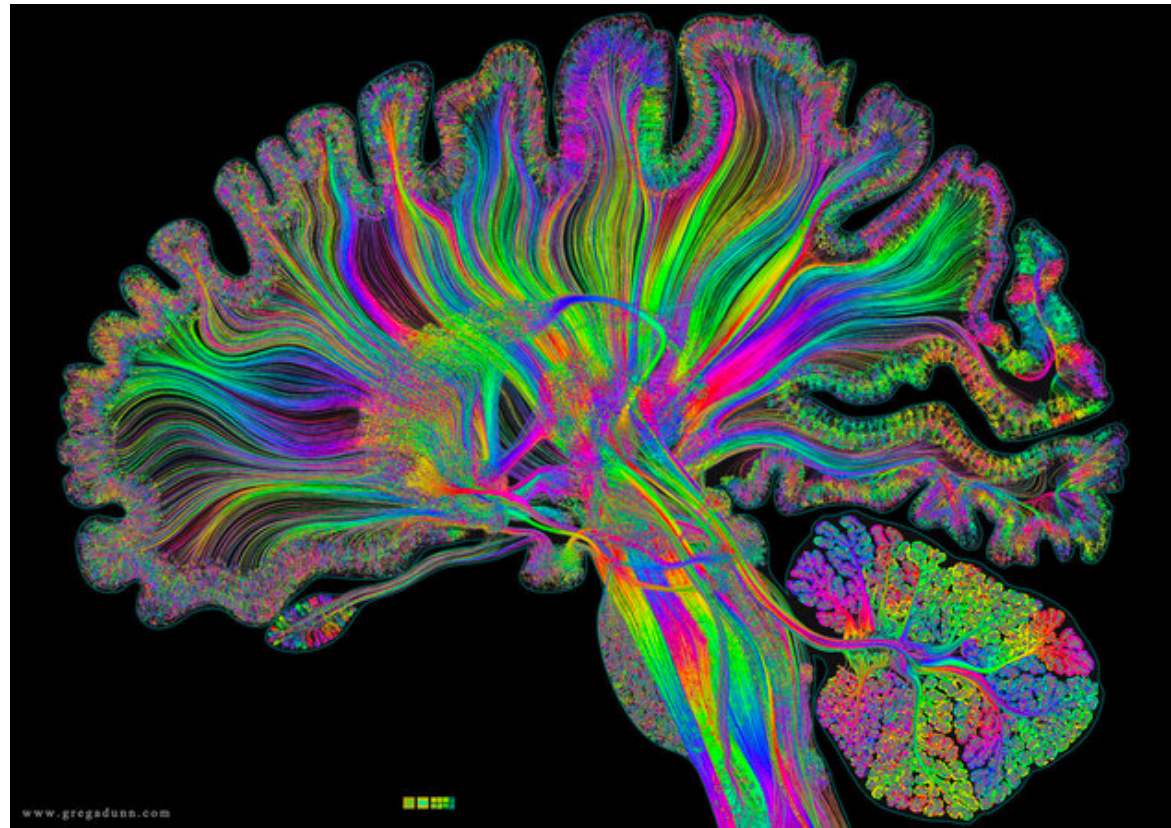
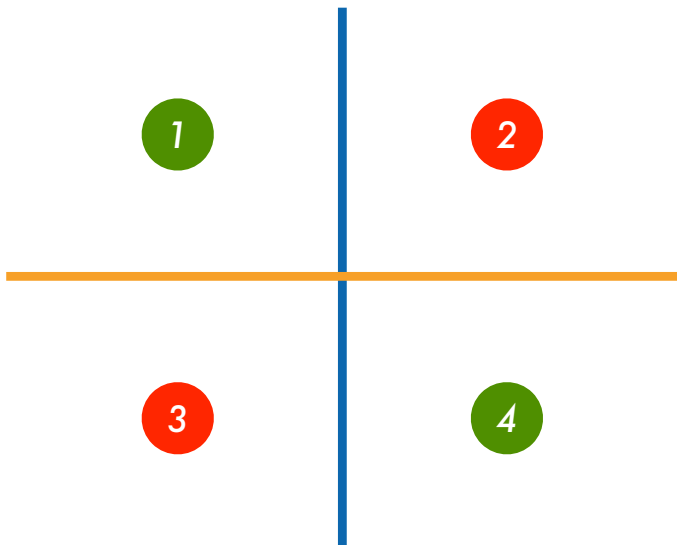Black & White

# XOR Problem (Minsky & Papert, 1969)

The perceptron cannot learn an XOR function
(neurons can only generate linear separators)

# Multilayer Perceptron
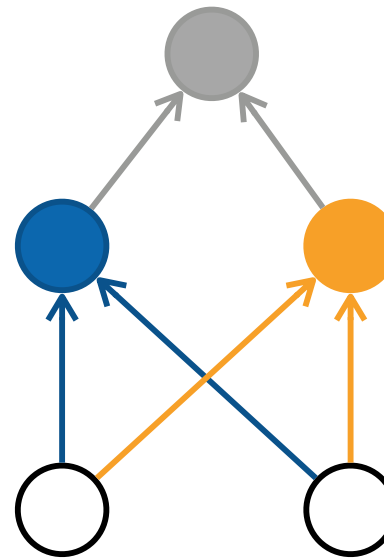


www.gregadunn.com
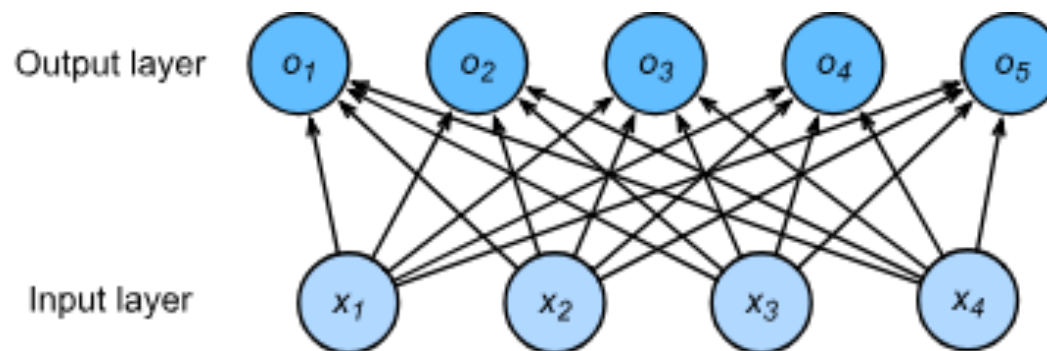
# Learning XOR

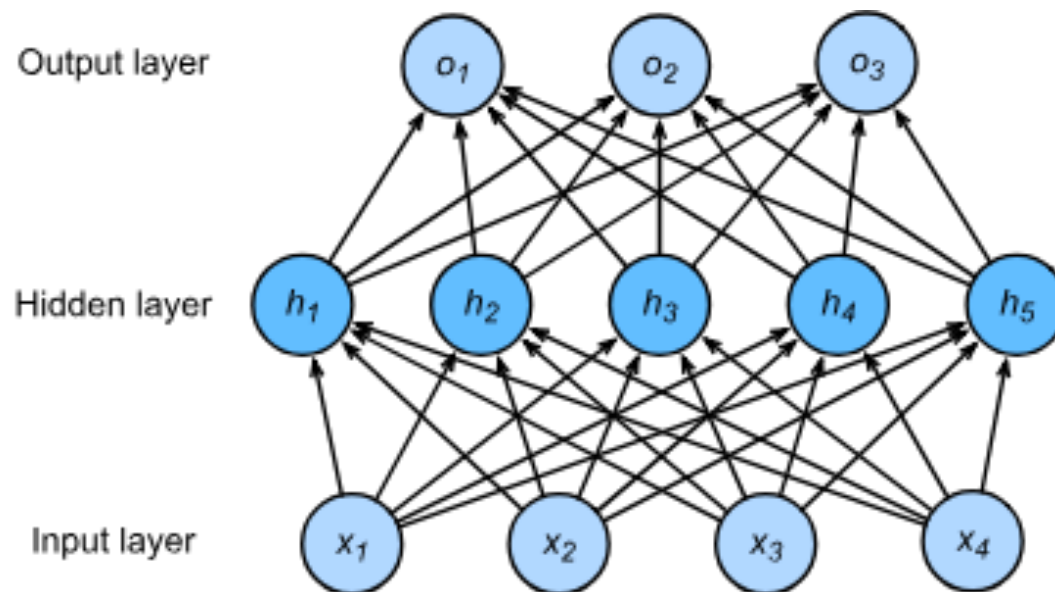| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| | + | - | + | - |
| | + | + | - | - |
| product | + | - | - | + |

# Single Hidden Layer

# Single Hidden Layer



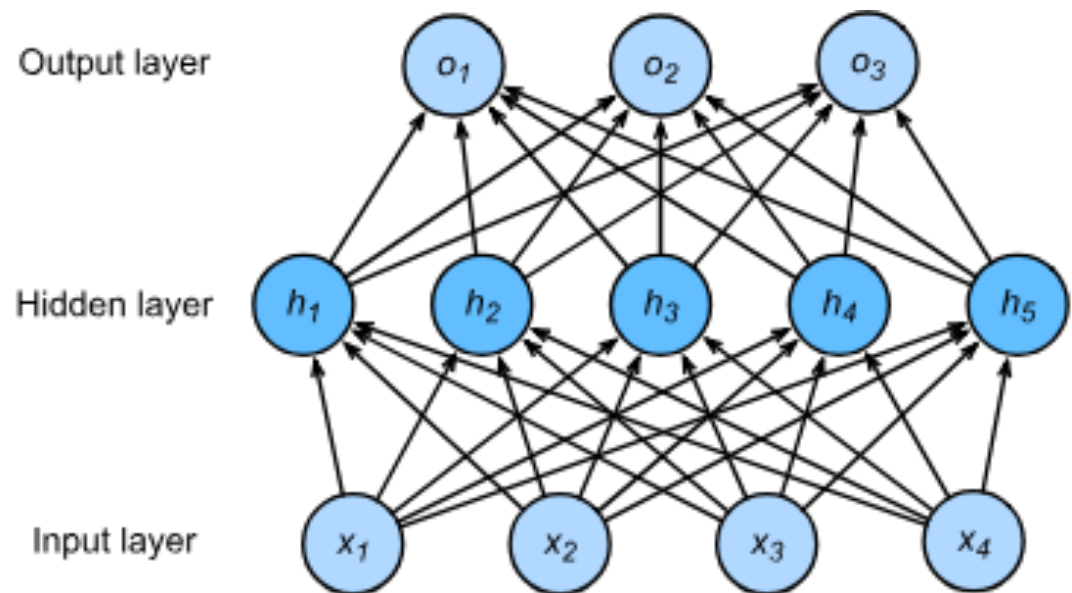Hyperparameter - size m of hidden layer

# Single Hidden Layer

- Input $\mathbf{x} \in \mathbb{R}^n$
- Hidden $\mathbf{W}_1 \in \mathbb{R}^{m \times n}, \mathbf{b}_1 \in \mathbb{R}^m$
- Output $\mathbf{w}_2 \in \mathbb{R}^m, b_2 \in \mathbb{R}$

$$\mathbf{h} = \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$$

$$\mathbf{o} = \mathbf{w}_2^T \mathbf{h} + b_2$$

$\sigma$ is an element-wise activation function

Output layer $o_1$ $o_2$ $o_3$

Hidden layer $h_1$ $h_2$ $h_3$ $h_4$ $h_5$
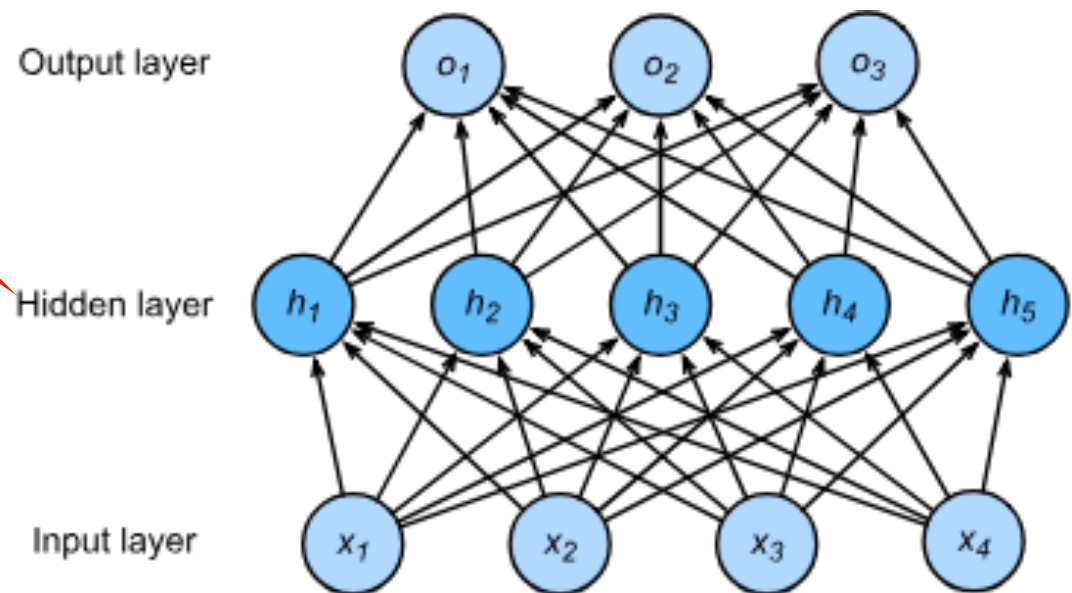
Input layer $x_1$ $x_2$ $x_3$ $x_4$

# Single Hidden Layer

Why do we need an a nonlinear activation?

$$\mathbf{h} = \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$$

$$\mathbf{o} = \mathbf{w}_2^T \mathbf{h} + b_2$$

$\sigma$ is an element-wise activation function

Output layer

Hidden layer

Input layer

$o_1$ $o_2$ $o_3$

$h_1$ $h_2$ $h_3$ $h_4$ $h_5$

$x_1$ $x_2$ $x_3$ $x_4$

# Single Hidden Layer

Why do we need an a nonlinear activation?

$$\mathbf{h} = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1$$

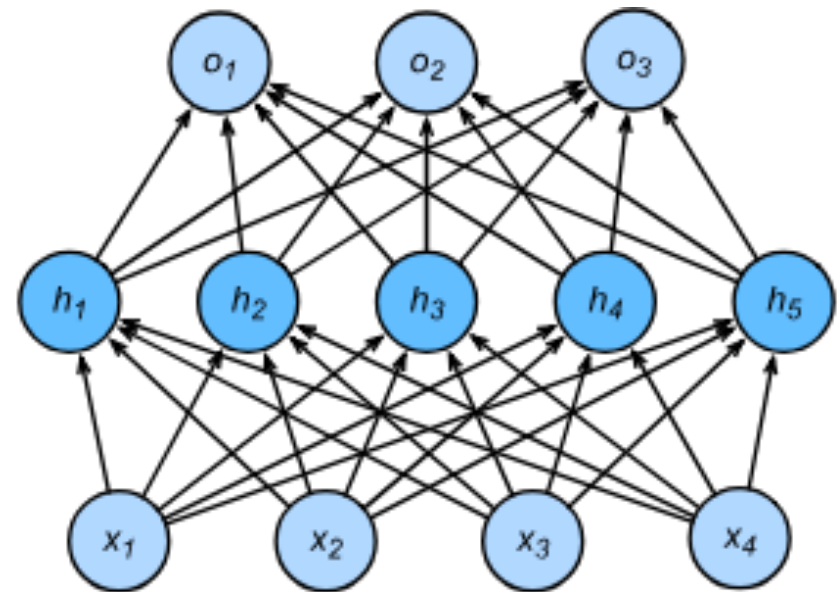$$\mathbf{o} = \mathbf{w}_2^T \mathbf{h} + b_2$$

hence $o = \mathbf{w}_2^\top \mathbf{W}_1 \mathbf{x} + b'$

Linear …
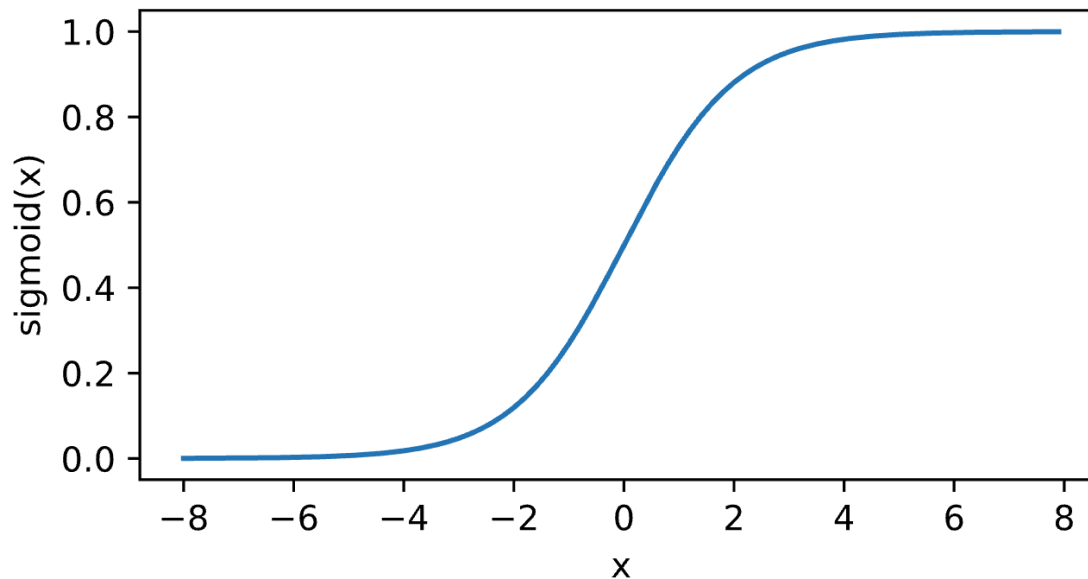
Output layer

Hidden layer

Input layer

# Sigmoid Activation

Map input into (0, 1), a soft version of $\quad \sigma(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$
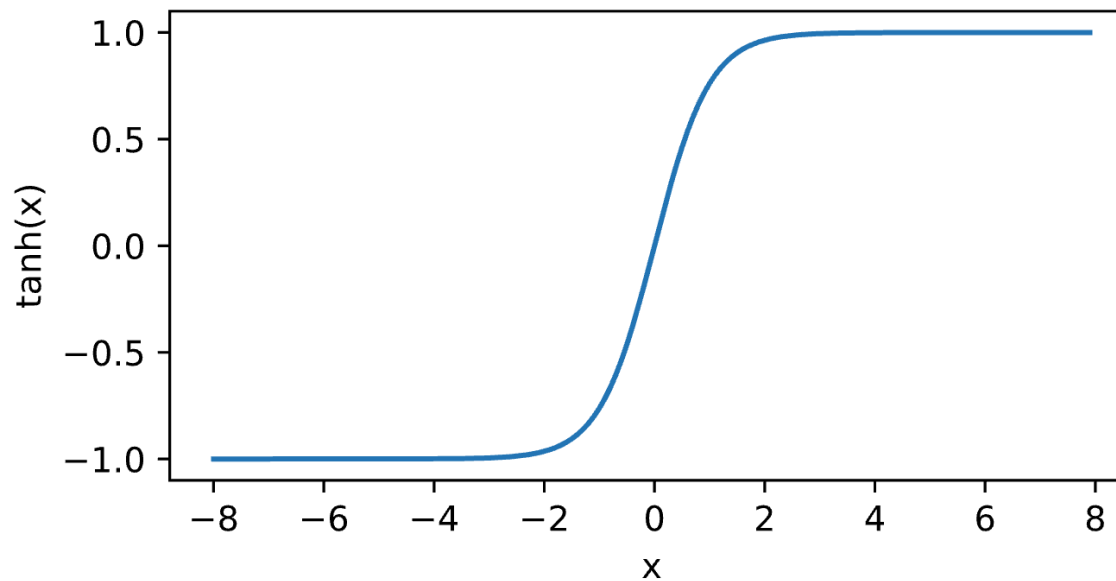
$$\text{sigmoid}(x) = \frac{1}{1 + \exp(-x)}$$

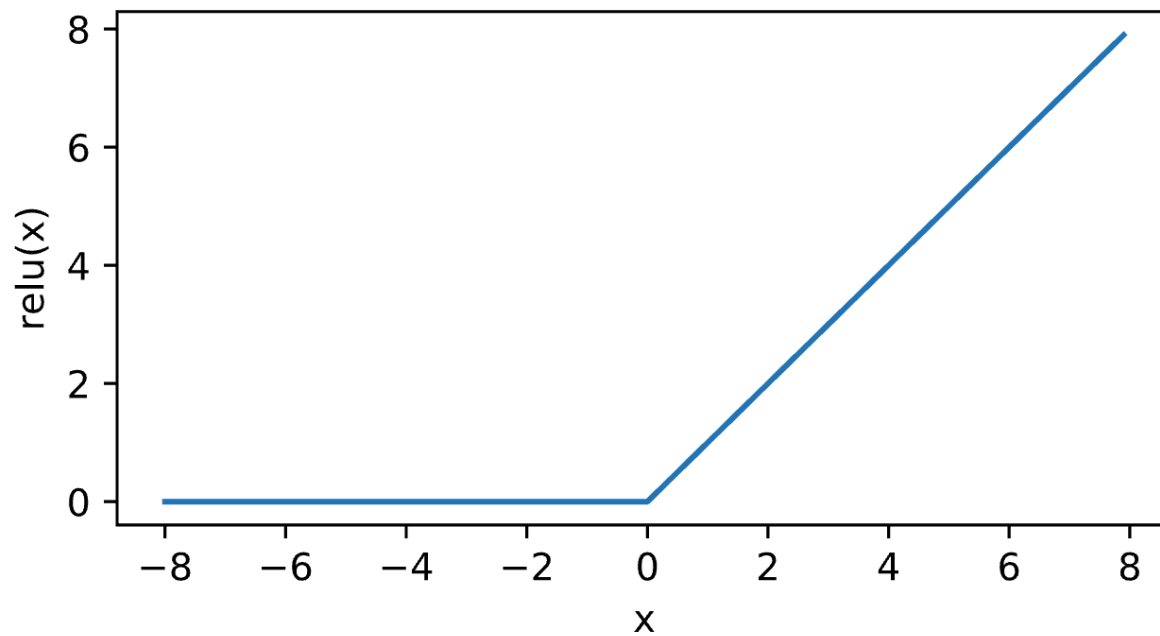# Tanh Activation

Map inputs into (-1, 1)

$$\tanh(x) = \frac{1 - \exp(-2x)}{1 + \exp(-2x)}$$
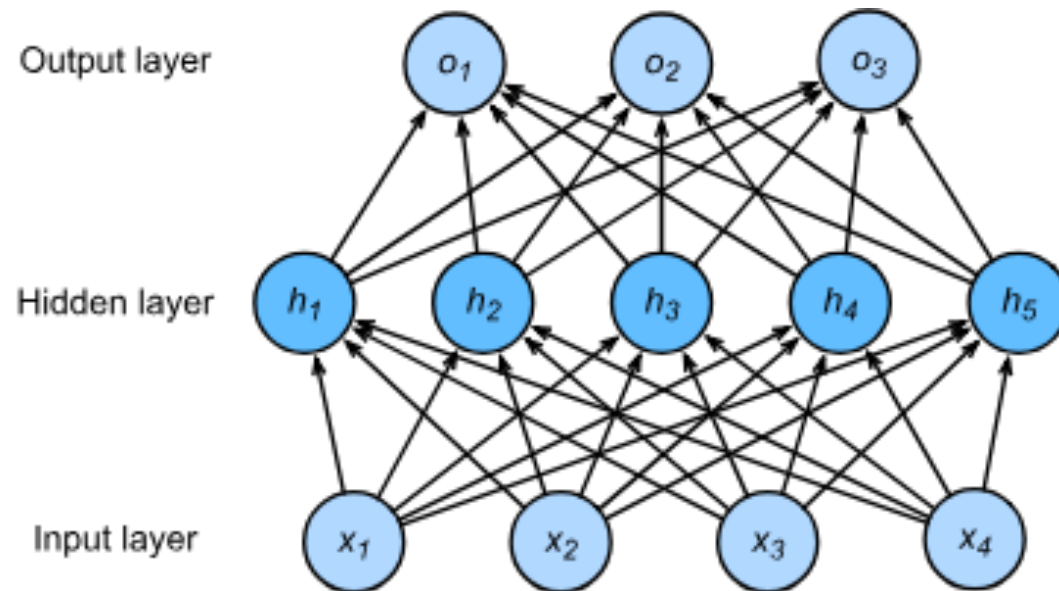
# ReLU Activation

ReLU: rectified linear unit

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

# Multiclass Classification

$$y_1, y_2, \ldots, y_k = \mathrm{softmax}(o_1, o_2, \ldots, o_k)$$



25

# Multiclass Classification

- Input $\mathbf{x} \in \mathbb{R}^n$
- Hidden $\mathbf{W}_1 \in \mathbb{R}^{m \times n}$ and $\mathbf{b}_1 \in \mathbb{R}^m$
- Output $\mathbf{W}_2 \in \mathbb{R}^{m \times d}$ and $\mathbf{b}_2 \in \mathbb{R}^d$

$$\mathbf{h} = \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$$

$$\mathbf{o} = \mathbf{w}_2^T \mathbf{h} + \mathbf{b}_2$$

$$\mathbf{y} = \mathrm{softmax}(o)$$



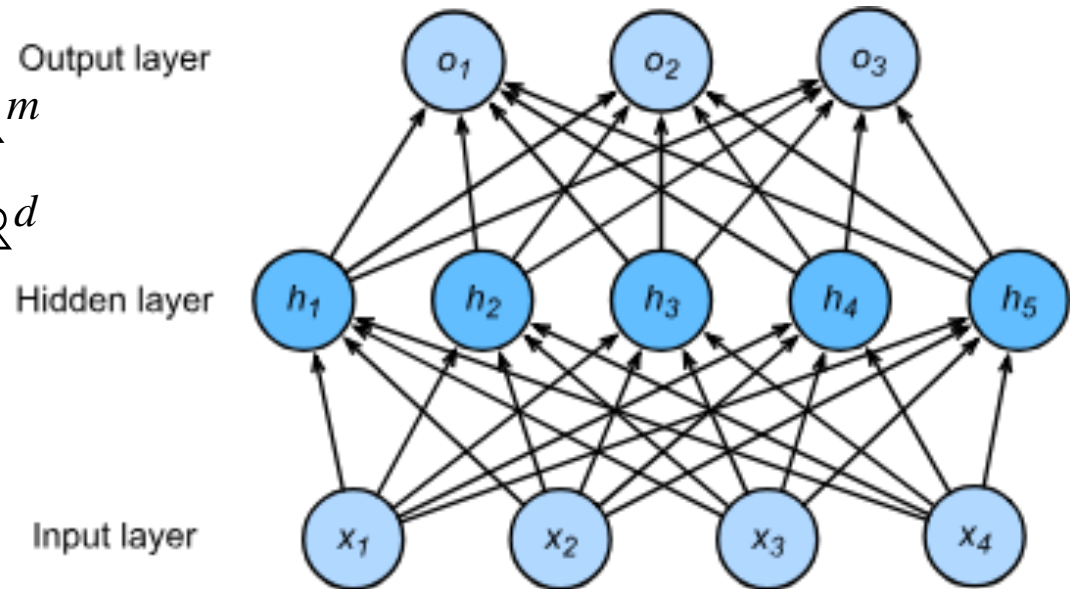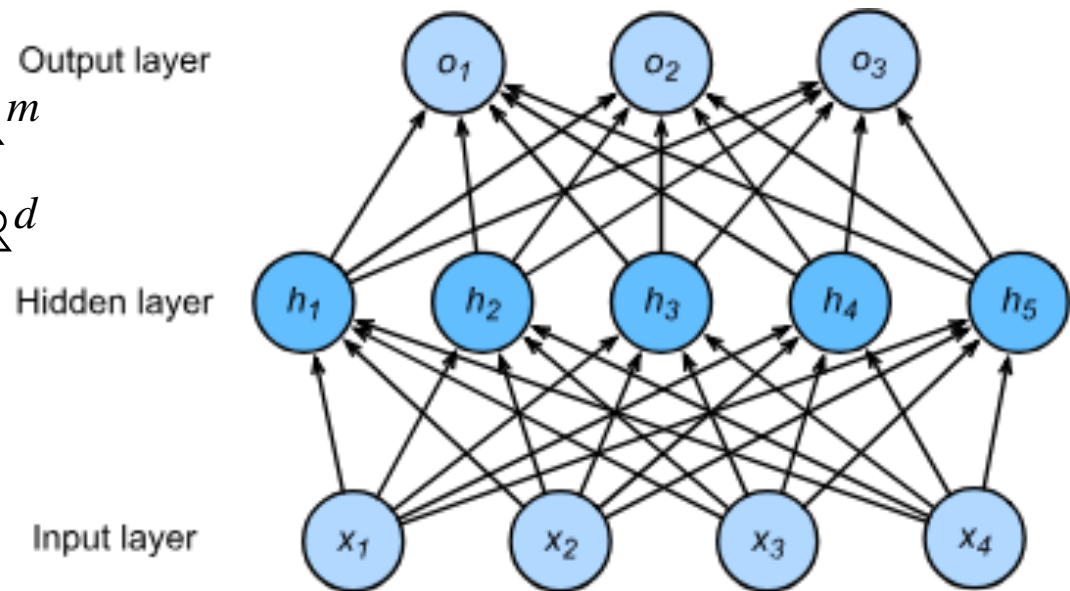Output layer

Hidden layer

Input layer

# Multiclass Classification

- Input $\mathbf{x} \in \mathbb{R}^n$
- Hidden $\mathbf{W}_1 \in \mathbb{R}^{m \times n}$ and $\mathbf{b}_1 \in \mathbb{R}^m$
- Output $\mathbf{W}_2 \in \mathbb{R}^{m \times d}$ and $\mathbf{b}_2 \in \mathbb{R}^d$

$$\mathbf{h} = \sigma(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1)$$

$$\mathbf{o} = \mathbf{w}_2^T\mathbf{h} + \mathbf{b}_2$$

$$\mathbf{y} = \mathrm{softmax}(o)$$

Output layer

Hidden layer

Input layer

$o_1$ $o_2$ $o_3$

$h_1$ $h_2$ $h_3$ $h_4$ $h_5$

$x_1$ $x_2$ $x_3$ $x_4$

# Multiple Hidden Layers

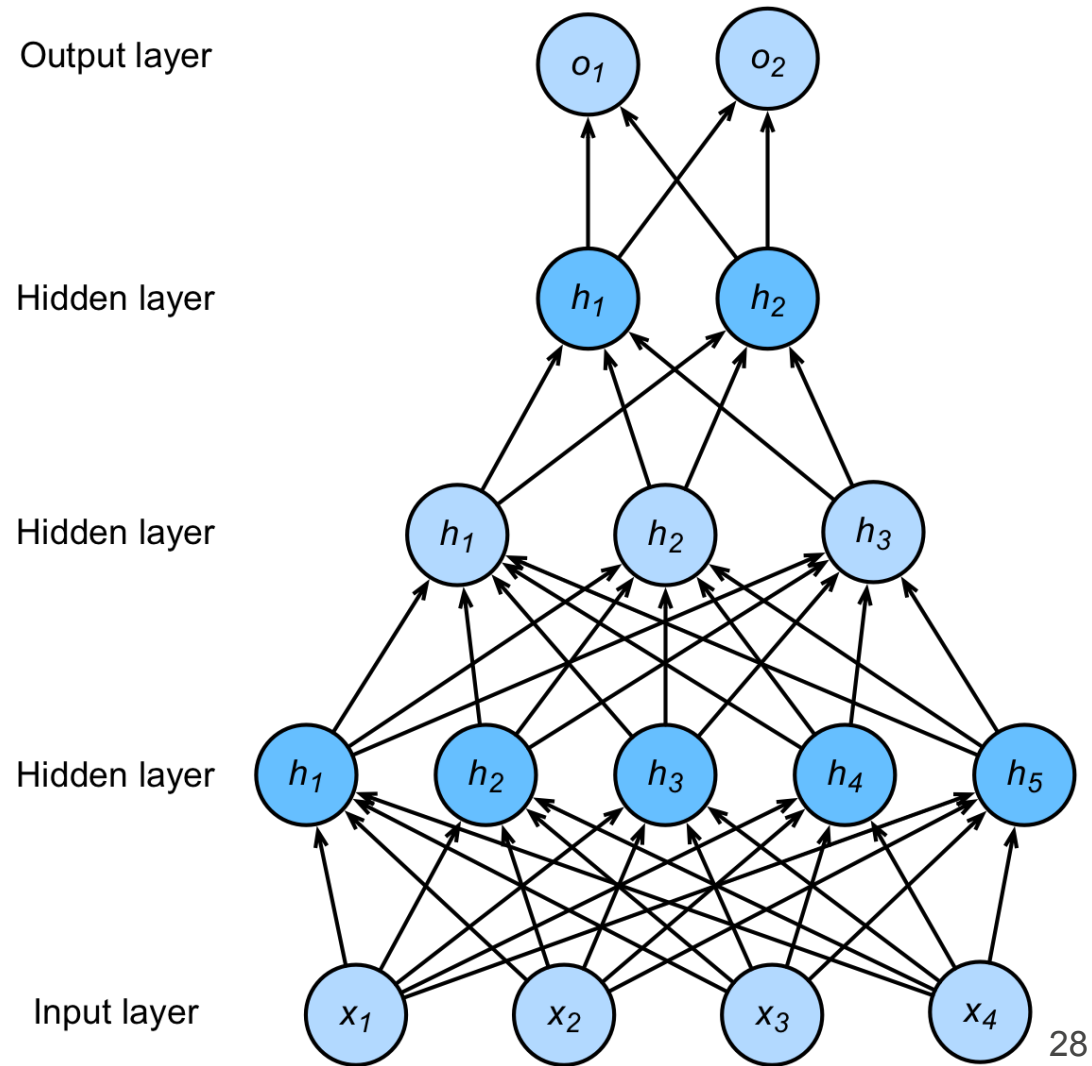$$\mathbf{h}_1 = \sigma(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1)$$

$$\mathbf{h}_2 = \sigma(\mathbf{W}_2\mathbf{h}_1 + \mathbf{b}_2)$$

$$\mathbf{h}_3 = \sigma(\mathbf{W}_3\mathbf{h}_2 + \mathbf{b}_3)$$

$$\mathbf{o} = \mathbf{W}_4\mathbf{h}_3 + \mathbf{b}_4$$

Hyper-parameters
- # of hidden layers
- Hidden size for each layer



Output layer

Hidden layer

Hidden layer

Hidden layer

Input layer

28

# Summary

- Perceptron
  - Simple updates
  - Limited function complexity
- Multilayer Perceptron
  - Multiple layers add more complexity
  - Nonlinearity is needed
  - Simple composition (but architecture search needed)