

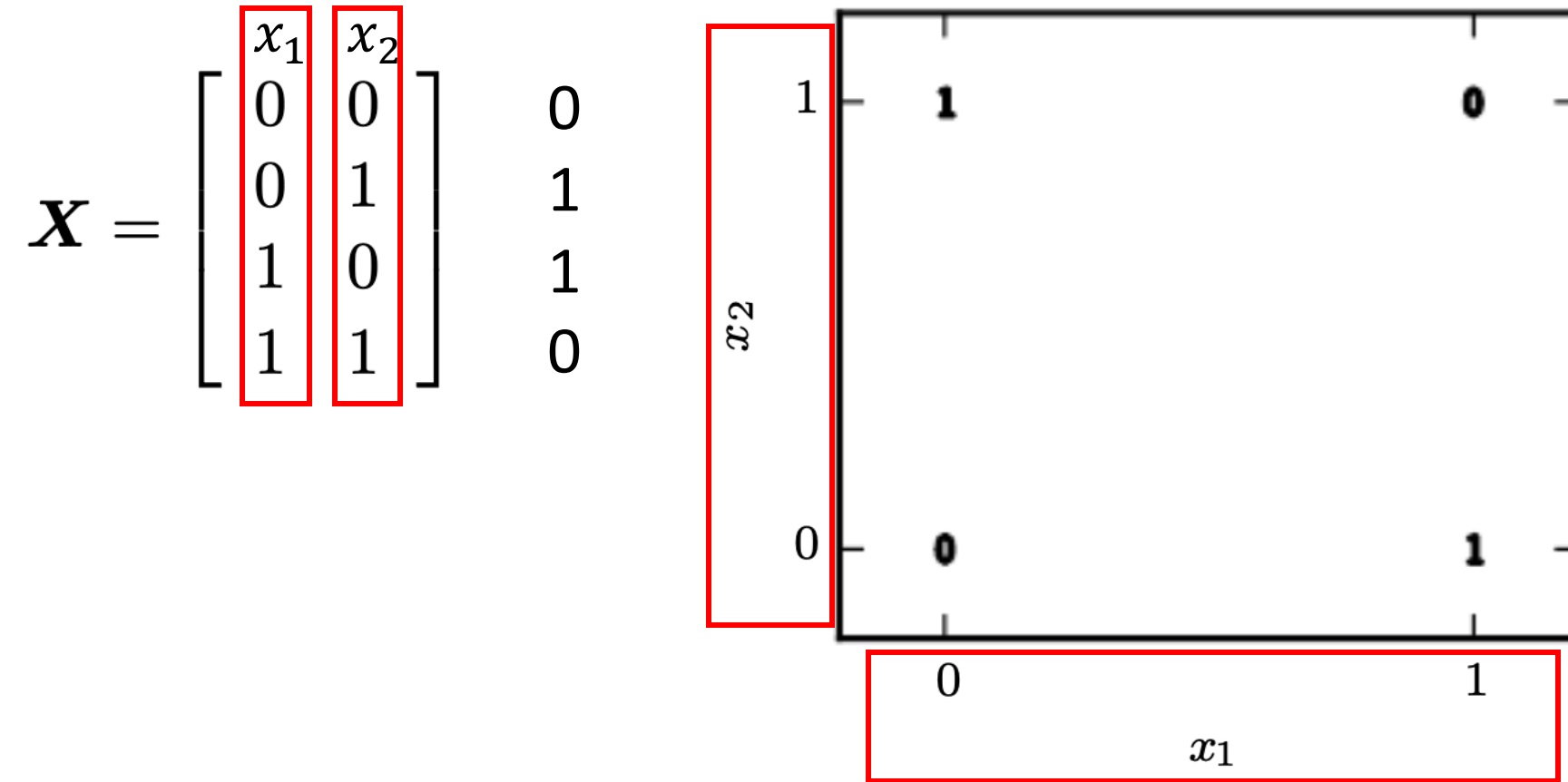
# Compositing Units in Neural Networks

CPT\_S 434/534 Neural network design and application

# Today's class

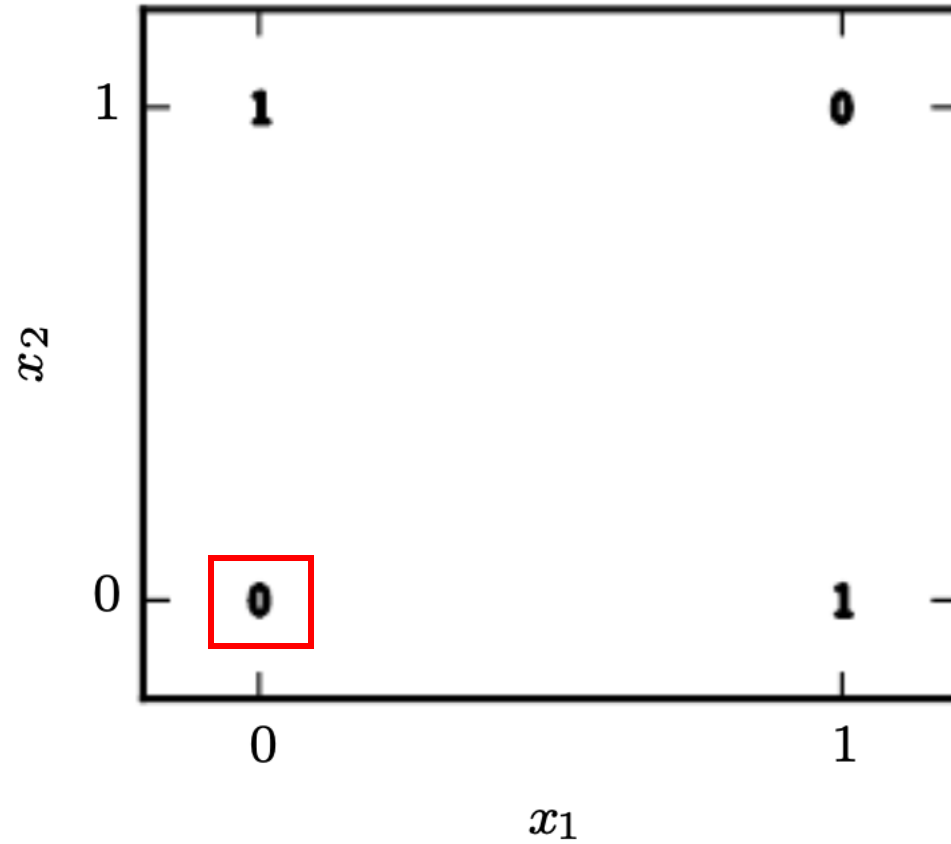
- What makes feedforward network different from linear model
  - Understanding its structure
- Units for neural networks
  - Output units  $\rightarrow$  cost function
    - $f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$
  - Hidden units
    - $f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$

# Learning XOR function



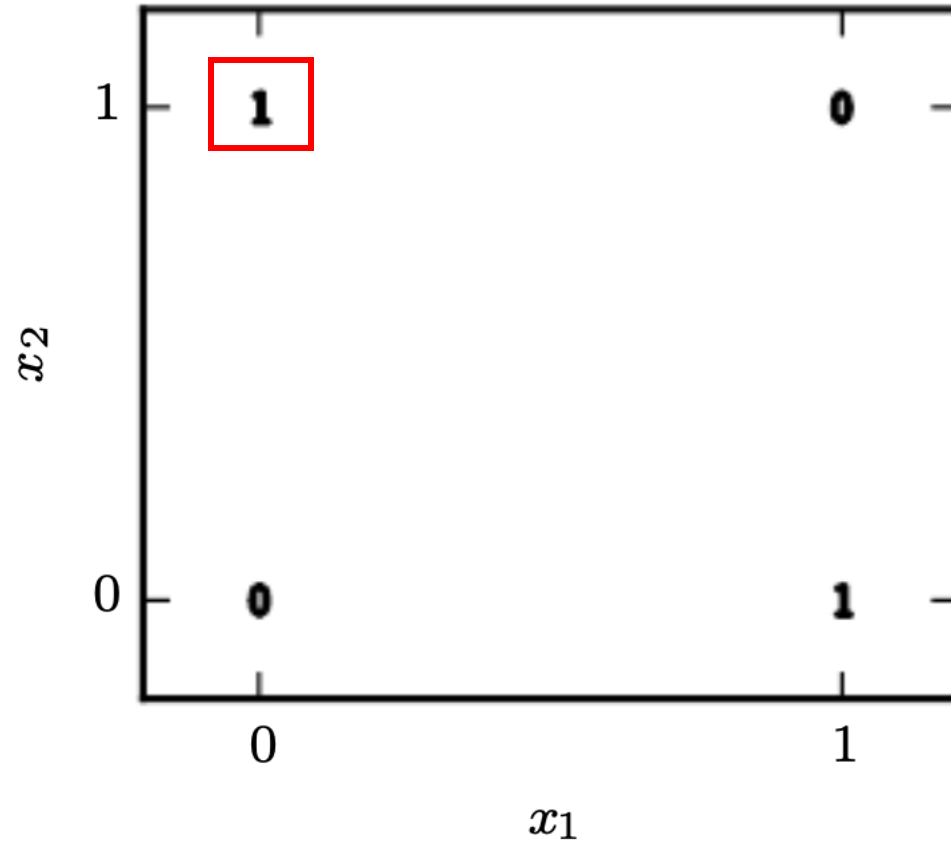
# Learning XOR function

$$\mathbf{X} = \begin{array}{c} \begin{array}{cc|c} x_1 & x_2 & \\ \hline 0 & 0 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{array} \end{array}$$



# Learning XOR function

$$\mathbf{X} = \begin{array}{cc|c} & x_1 & x_2 & \\ \hline & 0 & 0 & 0 \\ \hline & 0 & 1 & 1 \\ \hline & 1 & 0 & 1 \\ \hline & 1 & 1 & 0 \\ \hline \end{array}$$

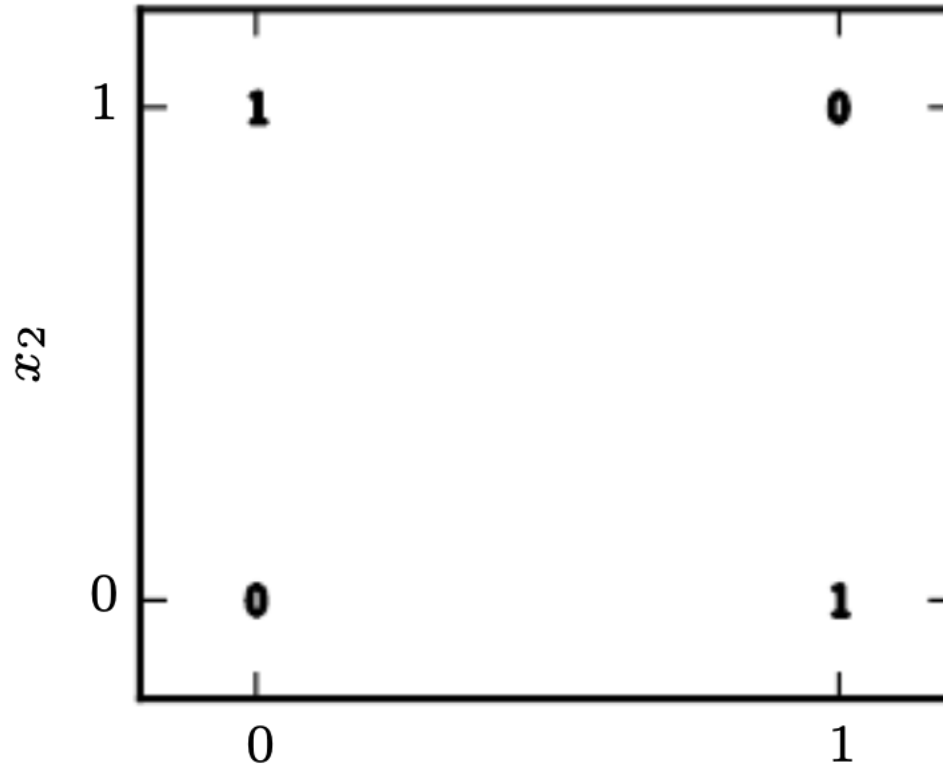


# Learning XOR function

$$\mathbf{X} = \begin{array}{cc|c} & x_1 & x_2 & \\ \hline & 0 & 0 & 0 \\ & 0 & 1 & 1 \\ & 1 & 0 & 1 \\ & 1 & 1 & 0 \end{array}$$

Q: Can we use a linear model to separate 0/1 classes?

$$\mathbf{W}^\top \mathbf{x} + \mathbf{c}$$

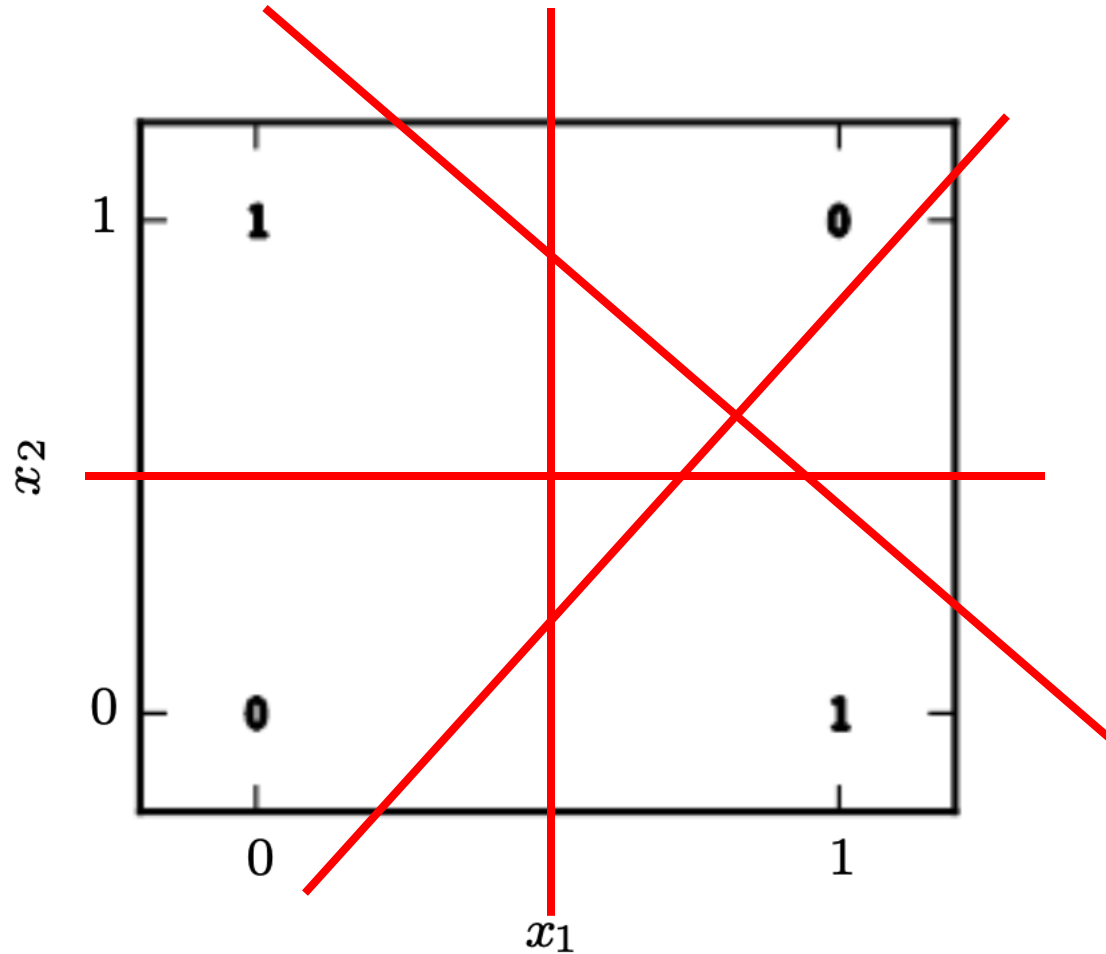


# Learning XOR function

$$\mathbf{X} = \begin{array}{cc|c} & x_1 & x_2 & \\ \hline & 0 & 0 & 0 \\ & 0 & 1 & 1 \\ & 1 & 0 & 1 \\ & 1 & 1 & 0 \end{array}$$

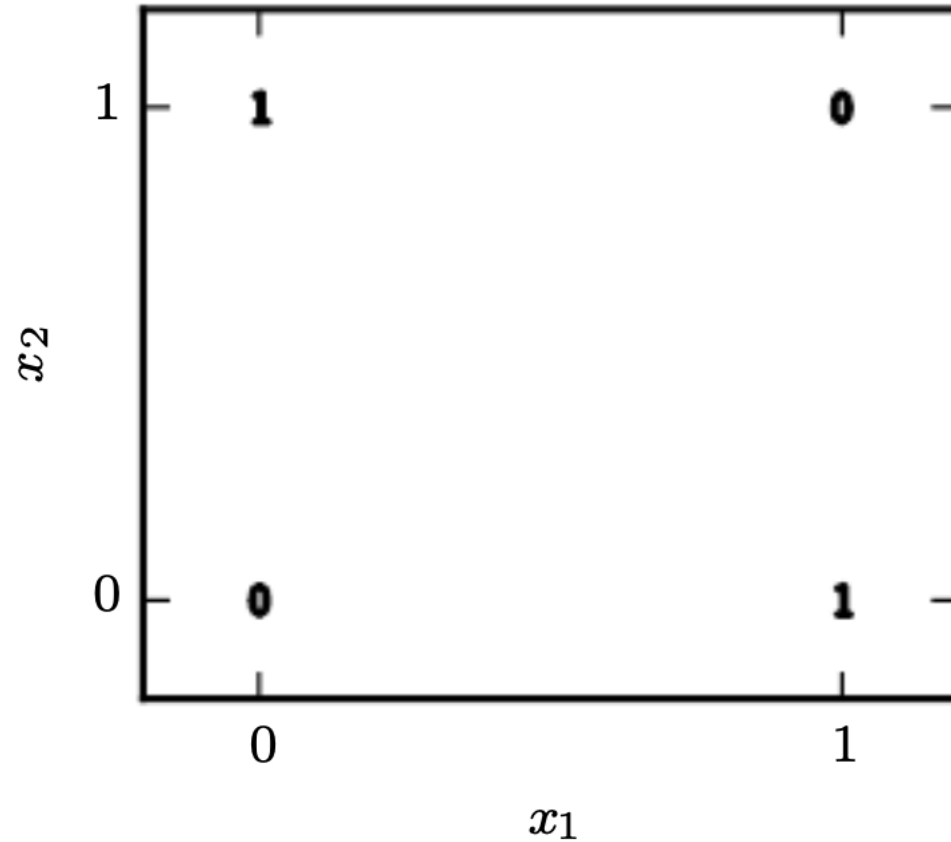
Q: Can we use a linear model to separate 0/1 classes?

$$\mathbf{W}^\top \mathbf{x} + \mathbf{c}$$



# Learning XOR function

$$\mathbf{X} = \begin{array}{cc|c} & x_1 & x_2 & \\ \hline & 0 & 0 & 0 \\ & 0 & 1 & 1 \\ & 1 & 0 & 1 \\ & 1 & 1 & 0 \end{array}$$

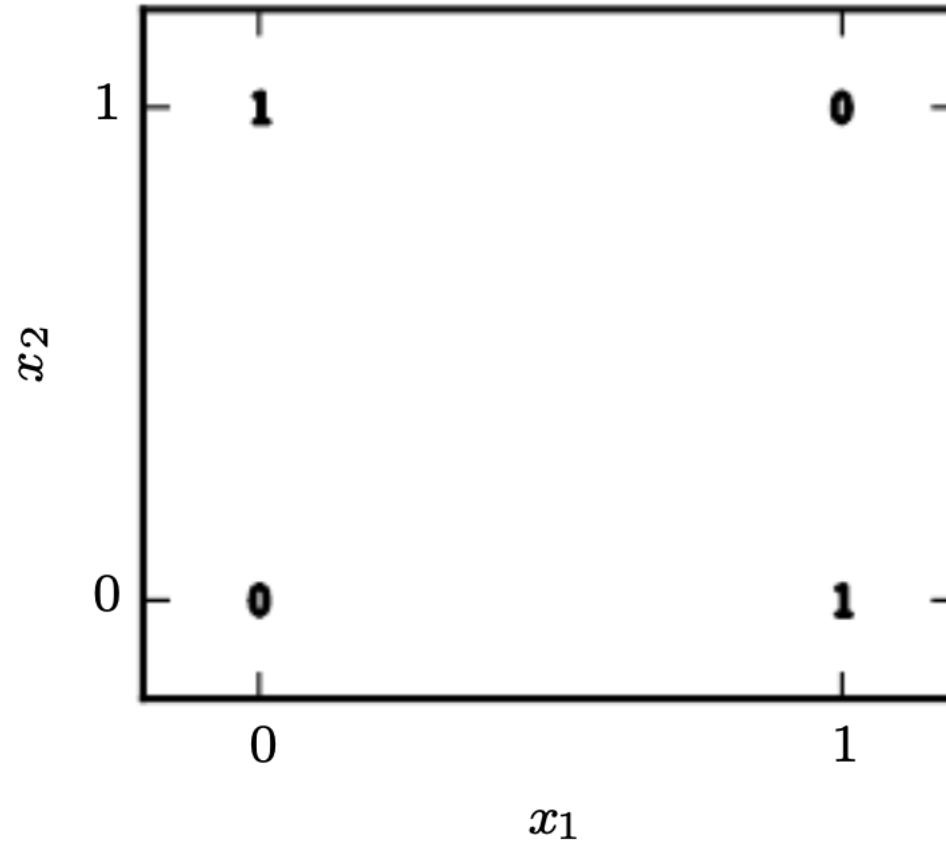


What if we use a nonlinear function as:  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + b$ .



# Learning XOR function

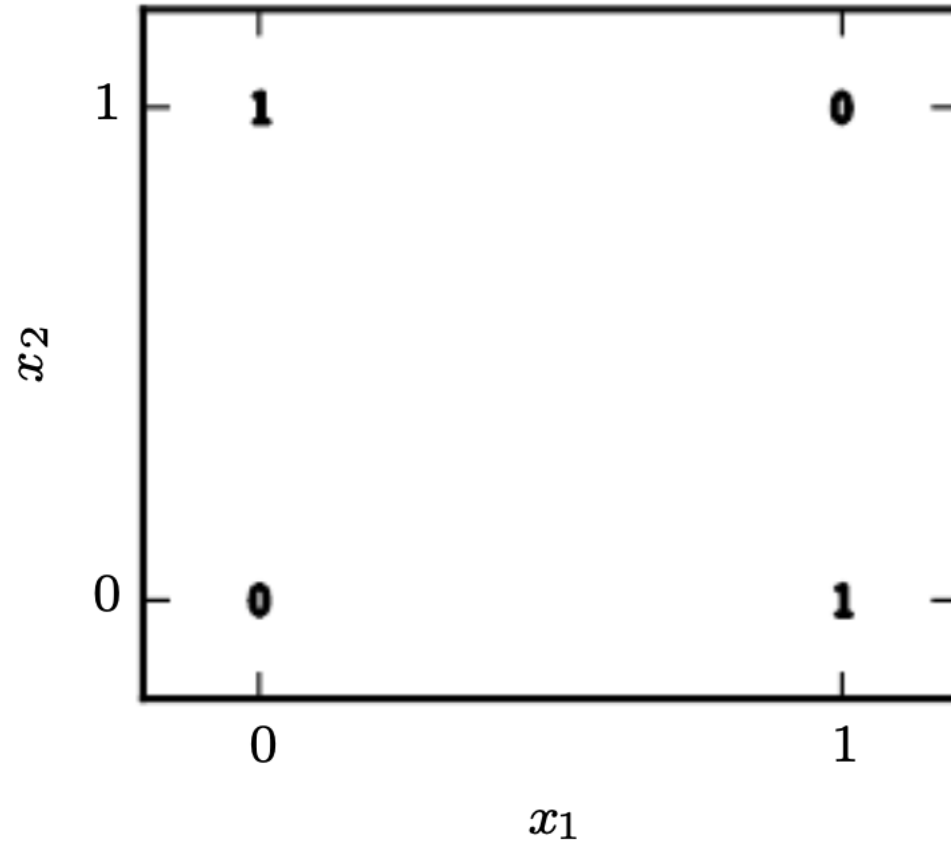
$$\mathbf{X} = \begin{array}{cc|c} & x_1 & x_2 & \\ \hline & 0 & 0 & 0 \\ & 0 & 1 & 1 \\ & 1 & 0 & 1 \\ & 1 & 1 & 0 \end{array}$$



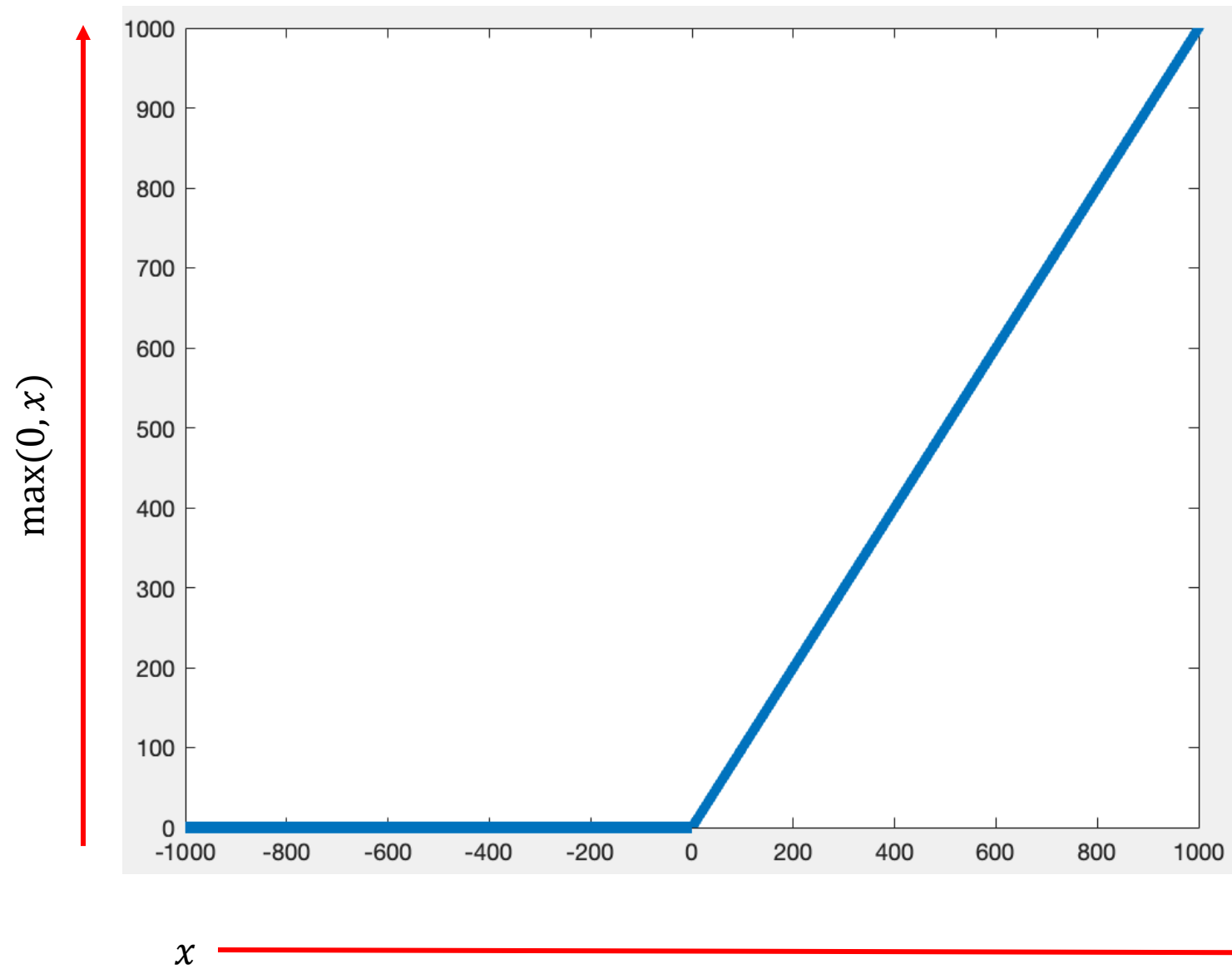
What if we use a nonlinear function as:  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + b$ .

# Learning XOR function

$$\mathbf{X} = \begin{array}{cc|c} & x_1 & x_2 & \\ \hline & 0 & 0 & 0 \\ & 0 & 1 & 1 \\ & 1 & 0 & 1 \\ & 1 & 1 & 0 \end{array}$$



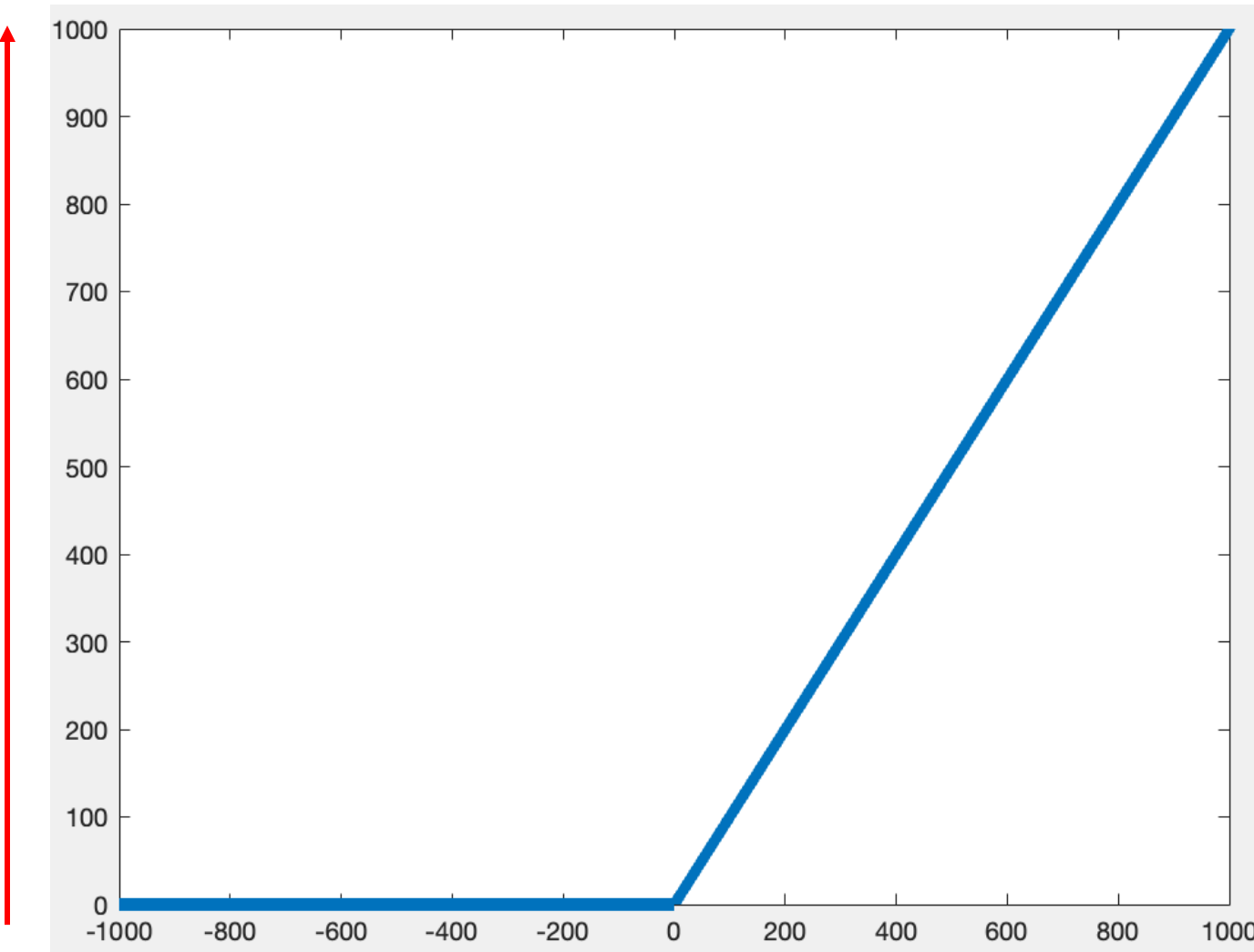
What if we use a nonlinear function as:  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + b$ .



Piecewise linear



$\max(0, x)$



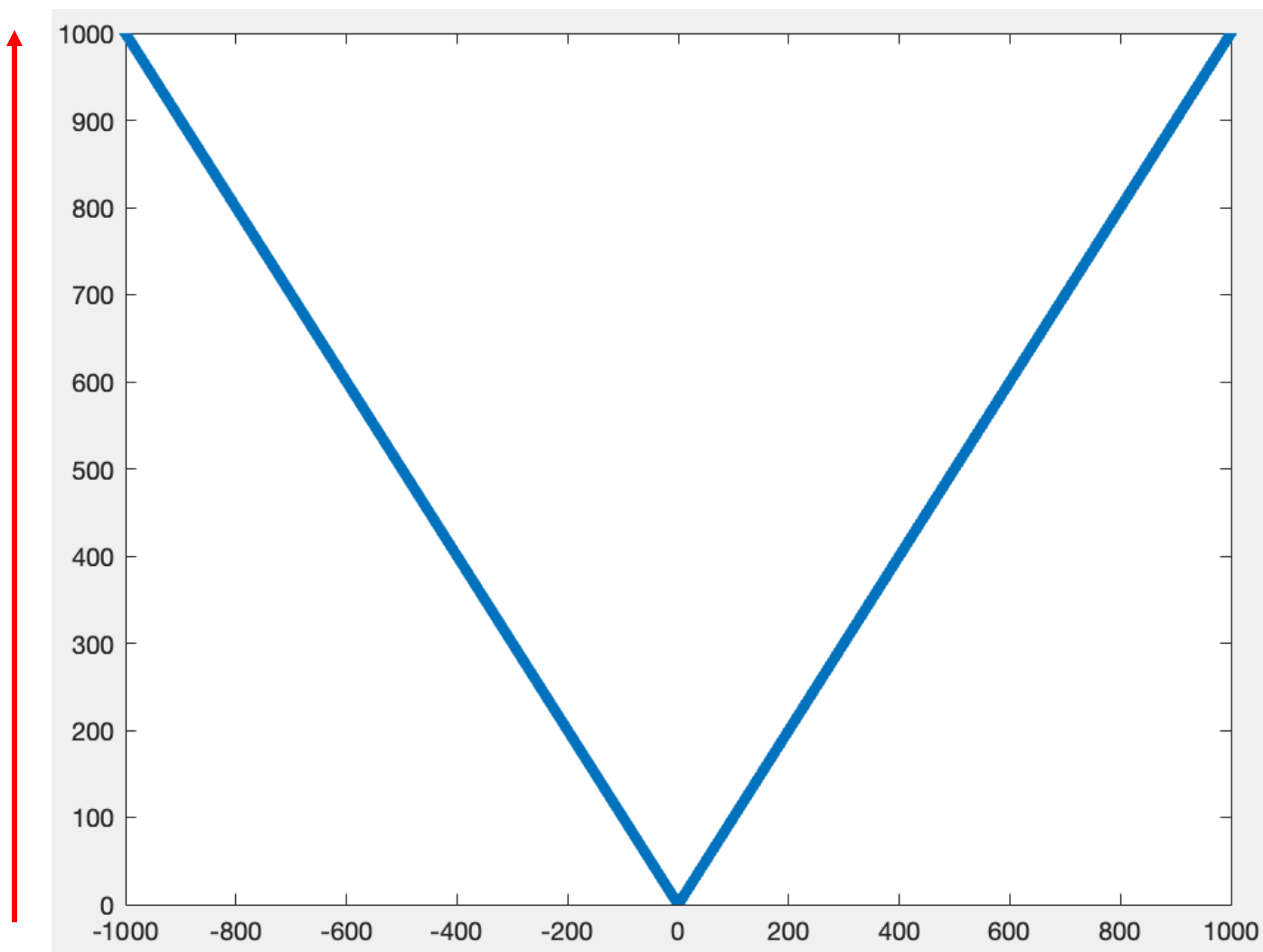
$x$



Piecewise linear



$\text{abs}(x)$

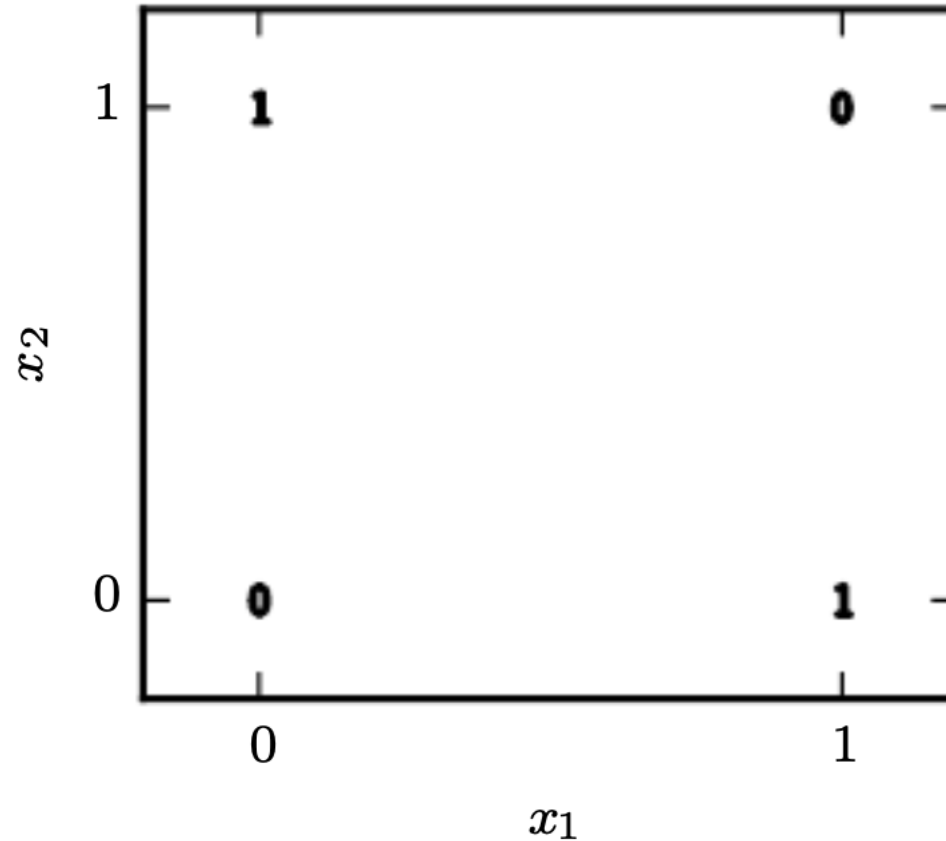


$x$



# Learning XOR function

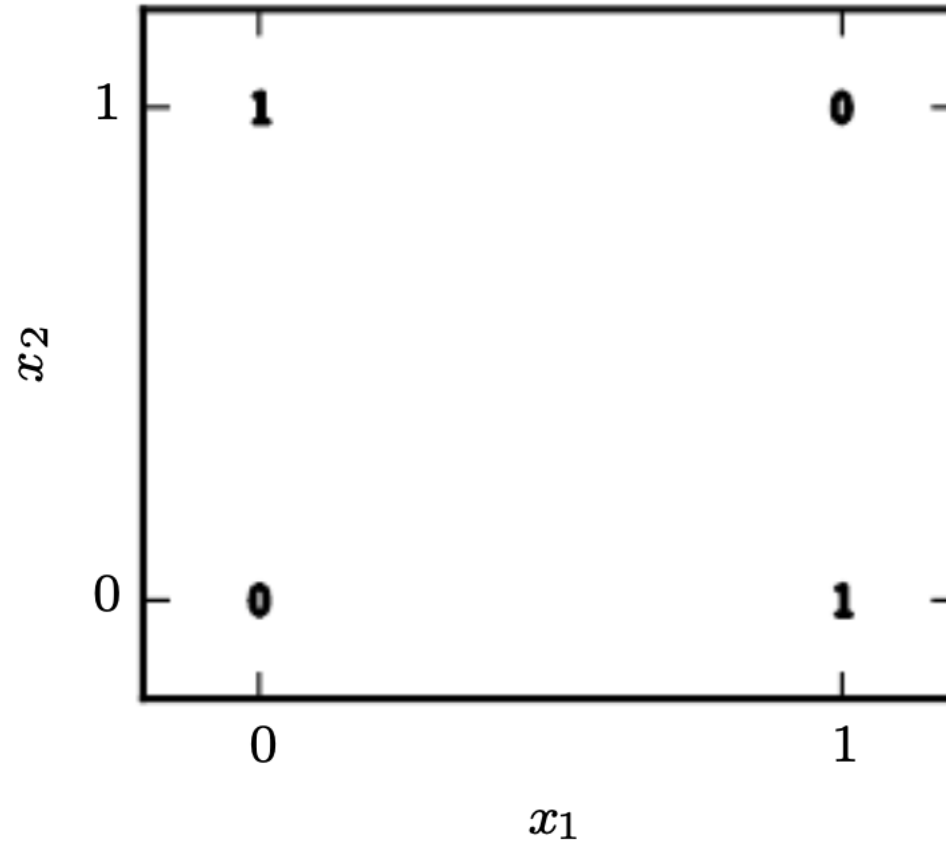
$$\mathbf{X} = \begin{array}{cc} & \begin{matrix} x_1 & x_2 \end{matrix} \\ \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} & \begin{matrix} 0 \\ 1 \\ 1 \\ 0 \end{matrix} \end{array}$$



What if we use a nonlinear function as:  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + b$ .

# Learning XOR function

$$\mathbf{X} = \begin{array}{cc|c} & x_1 & x_2 & \\ \hline & 0 & 0 & 0 \\ & 0 & 1 & 1 \\ & 1 & 0 & 1 \\ & 1 & 1 & 0 \end{array}$$

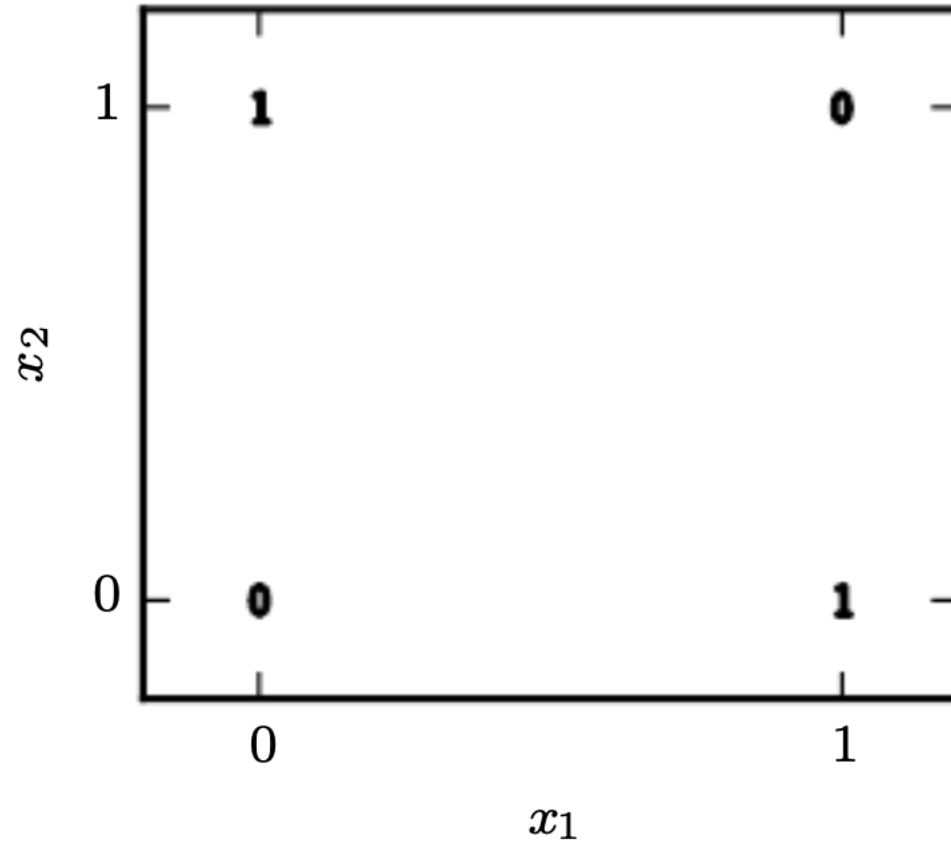


This output  $\rightarrow$  input of a linear function

What if we use a nonlinear function as:  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + b$ .

# Learning XOR function

$$\mathbf{X} = \begin{array}{cc} & \begin{matrix} x_1 & x_2 \end{matrix} \\ \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} & \begin{matrix} 0 \\ 1 \\ 1 \\ 0 \end{matrix} \end{array}$$



$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix},$$
$$\mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix},$$
$$\mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix},$$

What if we use a nonlinear function as:  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + \text{ReLU}$



# Learning XOR function

$$\mathbf{X} = \begin{array}{c} x_1 \quad x_2 \\ \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \end{array} \quad \begin{array}{c} 0 \\ 1 \\ 1 \\ 0 \end{array} \rightarrow \mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix},$$

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

$$\mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix},$$

What if we use a nonlinear function as:  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + \text{Ⓢ}$

# Learning XOR function

$$\mathbf{X} = \begin{array}{cc|c} & x_1 & x_2 & \\ \hline 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \end{array} \rightarrow \mathbf{XW} = \begin{array}{c|c} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{array}$$
$$\mathbf{W} = \begin{array}{c|c} 1 & 1 \\ 1 & 1 \end{array},$$
$$\mathbf{c} = \begin{array}{c} 0 \\ -1 \end{array},$$
$$\mathbf{w} = \begin{array}{c} 1 \\ -2 \end{array},$$

What if we use a nonlinear function as:  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + \text{Ⓢ}$

# Learning XOR function

$$\mathbf{X} = \begin{array}{cc|c} & x_1 & x_2 & \\ \hline 0 & 0 & & 0 \\ 0 & 1 & & 1 \\ 1 & 0 & & 1 \\ 1 & 1 & & 0 \end{array} \rightarrow \mathbf{XW} = \begin{array}{c|c} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{array}$$

$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix},$$

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

$$\mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix},$$

What if we use a nonlinear function as:  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + \text{Ⓢ}$

# Learning XOR function

$$\mathbf{X} = \begin{array}{cc} & \begin{matrix} x_1 & x_2 \end{matrix} \\ \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} & \begin{matrix} 0 \\ 1 \\ 1 \\ 0 \end{matrix} \end{array} \rightarrow \mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix},$$

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

$$\mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix},$$

What if we use a nonlinear function as:  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + \text{Ⓢ}$

# Learning XOR function

$$\mathbf{X} = \begin{array}{cc} & \begin{matrix} x_1 & x_2 \end{matrix} \\ \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} & \begin{matrix} 0 \\ 1 \\ 1 \\ 0 \end{matrix} \end{array} \rightarrow \mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$
$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix},$$
$$\mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix},$$
$$\mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix},$$

What if we use a nonlinear function as:  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + \text{Ⓢ}$

# Learning XOR function

$$\mathbf{X} = \begin{array}{cc} & \begin{matrix} x_1 & x_2 \end{matrix} \\ \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} & \begin{matrix} 0 \\ 1 \\ 1 \\ 0 \end{matrix} \end{array} \rightarrow \mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} + \mathbf{c}$$

$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix},$$

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

$$\mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix},$$

What if we use a nonlinear function as:  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + \text{ⓧ}$

# Learning XOR function

$$\mathbf{X} = \begin{array}{cc|c} & x_1 & x_2 & \\ \hline & 0 & 0 & 0 \\ & 0 & 1 & 1 \\ & 1 & 0 & 1 \\ & 1 & 1 & 0 \end{array} \rightarrow \mathbf{XW} = \begin{array}{cc|c} & & & \\ \hline & 0 & 0 & \\ & 1 & 1 & \\ & 1 & 1 & \\ & 2 & 2 & \end{array} + \mathbf{c}$$
$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix},$$
$$\mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix},$$
$$\mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix},$$

What if we use a nonlinear function as:  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + \text{Ⓢ}$

# Learning XOR function

$$\mathbf{X} = \begin{array}{cc} & \begin{matrix} x_1 & x_2 \end{matrix} \\ \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} & \begin{matrix} 0 \\ 1 \\ 1 \\ 0 \end{matrix} \end{array} \rightarrow \mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} + \mathbf{c}$$

$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix},$$

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

$$\mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix},$$

What if we use a nonlinear function?

$$f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + \text{Ⓢ}$$



# Learning XOR function

$$\mathbf{X} = \begin{array}{cc} & \begin{matrix} x_1 & x_2 \end{matrix} \\ \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} & \begin{matrix} 0 \\ 1 \\ 1 \\ 0 \end{matrix} \end{array} \rightarrow \mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} + \mathbf{c} \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix} \quad \begin{matrix} \mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \\ \mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \\ \mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, \end{matrix}$$

What if we use a nonlinear function?

$$f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + \text{ⓧ}$$

# Learning XOR function

$$\mathbf{X} = \begin{array}{cc} & \begin{matrix} x_1 & x_2 \end{matrix} \\ \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} & \begin{matrix} 0 \\ 1 \\ 1 \\ 0 \end{matrix} \end{array} \rightarrow \mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} + \mathbf{c} \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix} \quad \mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \\
 \mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \\
 \mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix},$$

The inner linear model

What if we use a nonlinear function as:  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + \text{ⓧ}$

# Learning XOR function

$$\mathbf{X} = \begin{matrix} & \begin{matrix} x_1 & x_2 \end{matrix} \\ \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} & \begin{matrix} 0 \\ 1 \\ 1 \\ 0 \end{matrix} \end{matrix} \rightarrow \mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} + \mathbf{c} \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$$

$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

$$\begin{matrix} \text{max}(0, \cdot) \\ \downarrow \\ \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix} \end{matrix}$$

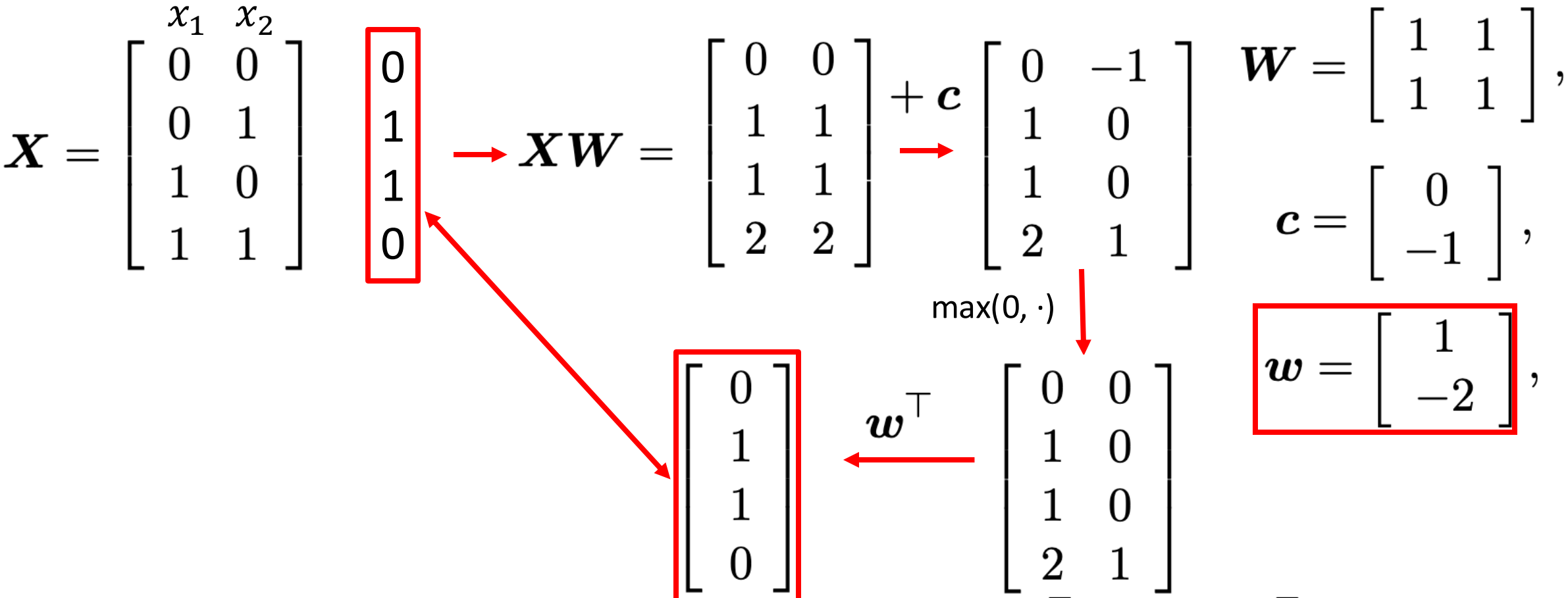
What if we use a nonlinear function as:  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + \text{ⓧ}$

# Learning XOR function

$$\mathbf{X} = \begin{matrix} & \begin{matrix} x_1 & x_2 \end{matrix} \\ \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} & \begin{matrix} 0 \\ 1 \\ 1 \\ 0 \end{matrix} \end{matrix} \rightarrow \mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} \xrightarrow{+ \mathbf{c}} \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix} \quad \mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \\
 \mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \\
 \xrightarrow{\max(0, \cdot)} \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, \\
 \xrightarrow{\mathbf{w}^\top} \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

What if we use a nonlinear function as:  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + b$

# Learning XOR function

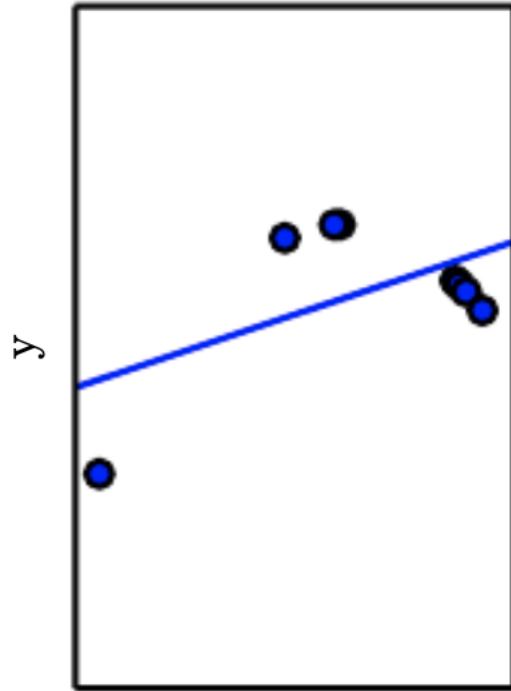


What if we use a nonlinear function as:  $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + \text{[red circle with slash]}$

# Learning XOR function

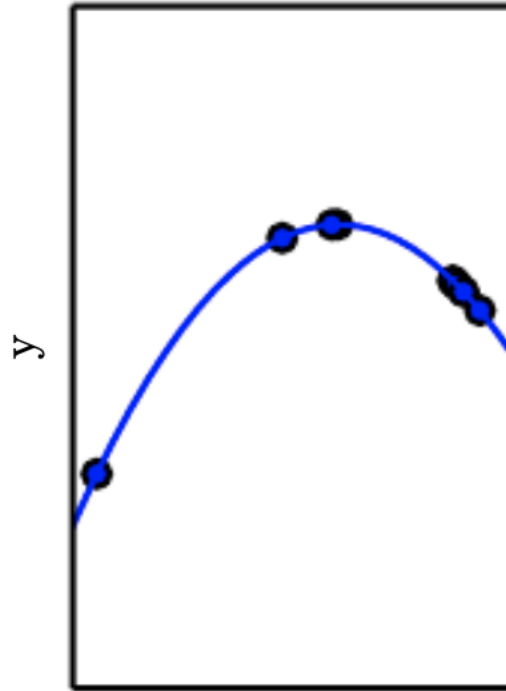
$$\mathbf{W}^\top \mathbf{x} + \mathbf{c} \quad \text{v.s.} \quad \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\}$$

# Learning XOR function



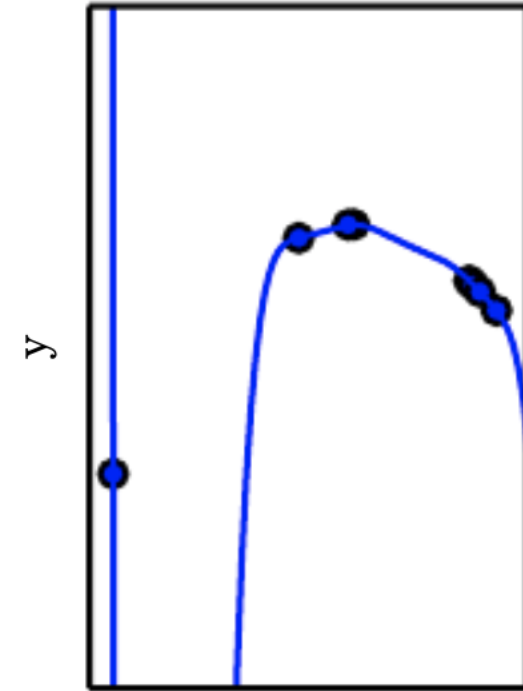
$x_0$

Linear model  
 $f(w; x) = w_1 x^1 + w_0$



$x_0$

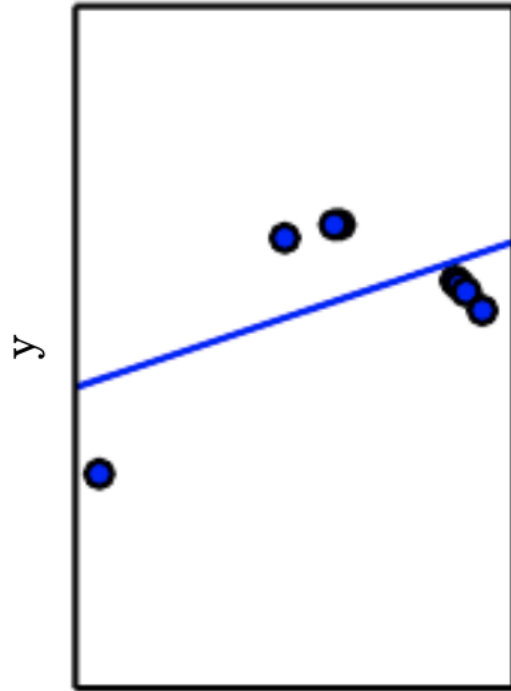
Quadratic model  
 $f(w; x) = w_2 x^2 + w_1 x^1 + w_0$



$x_0$

Polynomial model (9 degree)  
 $f(w; x) = \sum_{i=1}^9 w_i x^i + w_0$

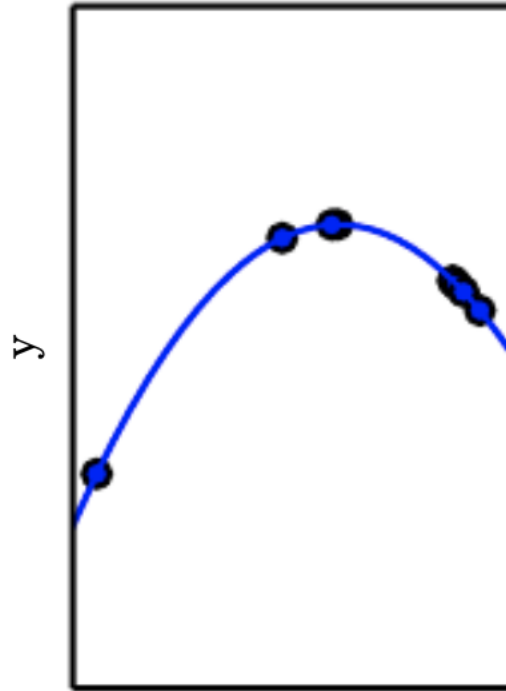
# Learning XOR function



$x_0$

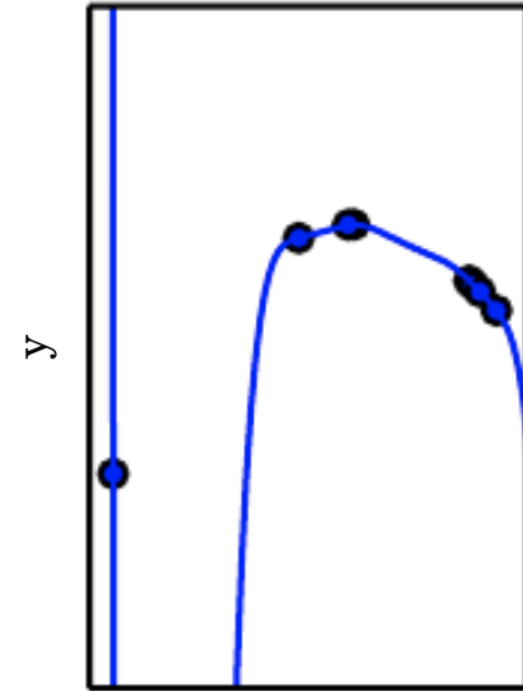
Linear model  
 $f(w; x) = w_1 x^1 + w_0$

$$= \sum_{i=2}^9 0 \cdot x^i + \sum_{i=1}^1 w_i x^i + w_0$$



$x_0$

Quadratic model  
 $f(w; x) = w_2 x^2 + w_1 x^1 + w_0$



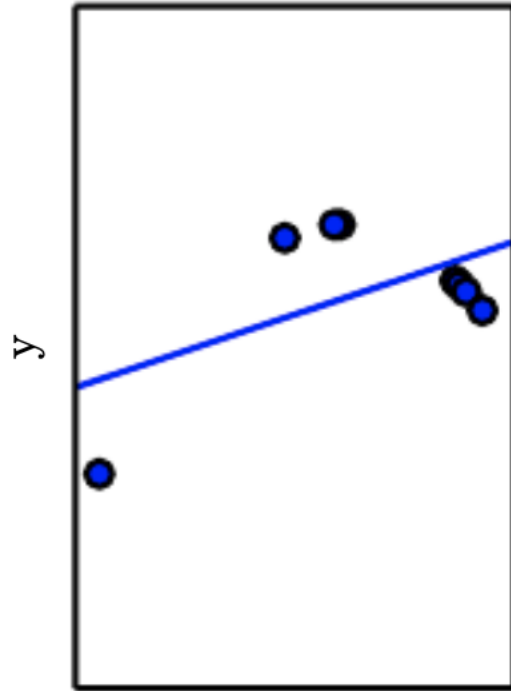
$x_0$

Polynomial model (9 degree)

$$f(w; x) = \sum_{i=1}^9 w_i x^i + w_0$$

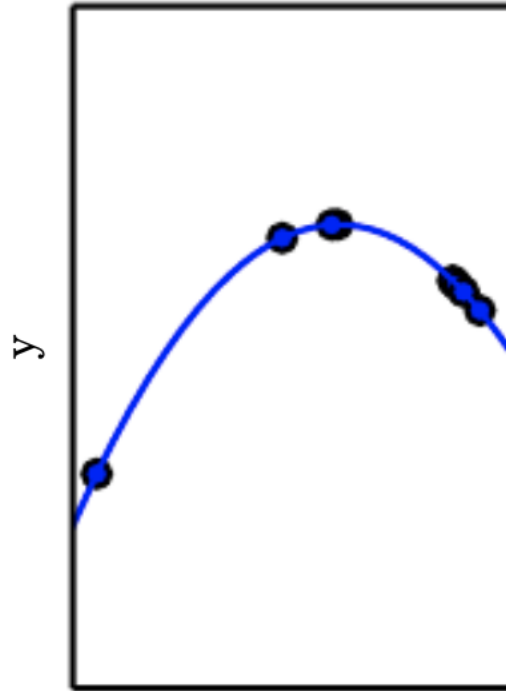


# Learning XOR function



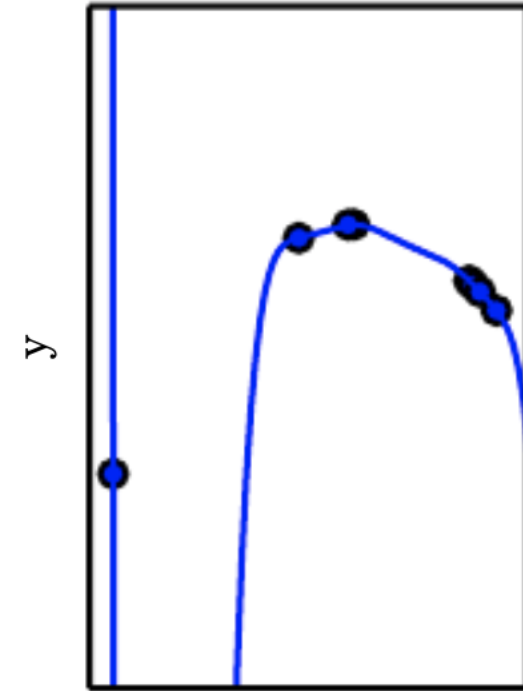
$x_0$

Linear model  
 $f(w; x) = w_1 x^1 + w_0$



$x_0$

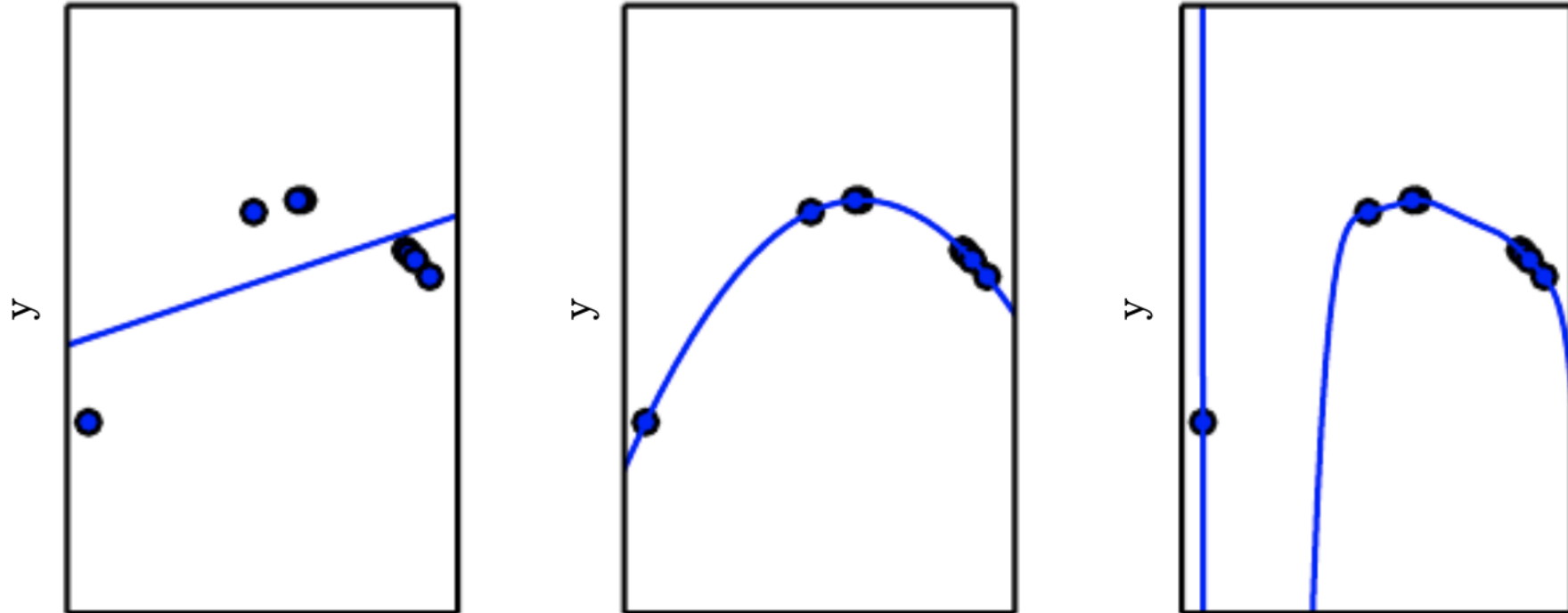
Quadratic model  
 $f(w; x) = w_2 x^2 + w_1 x^1 + w_0$   
 $= \sum_{i=3}^9 0 \cdot x^i + \sum_{i=1}^2 w_i x^i + w_0$



$x_0$

Polynomial model (9 degree)  
 $f(w; x) = \sum_{i=1}^9 w_i x^i + w_0$

# Learning XOR function



	Linear function	Quadratic function	9-degree polynomial function
Linear function	Yes	No	No
Quadratic function	Yes	Yes	No
9-degree polynomial function	Yes	Yes	Yes

Best capacity

# Learning XOR function

$$\mathbf{W}^\top \mathbf{x} + \mathbf{c}$$

v.s.

$$\mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\}$$

nonlinear model: better capacity

	Linear function	Quadratic function	9-degree polynomial function
Linear function	Yes	No	No
Quadratic function	Yes	Yes	No
9-degree polynomial function	Yes	Yes	Yes

Best capacity

# Learning XOR function

$$\mathbf{W}^\top \mathbf{x} + \mathbf{c}$$



Raw features

v.s.

$$\mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\}$$

nonlinear model: better capacity  
(Activation layer)

# Learning XOR function

$$\mathbf{W}^\top \mathbf{x} + \mathbf{c}$$



Raw features

v.s.

$$\mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\}$$



nonlinear model: better capacity

Raw features

# Learning XOR function

$$\mathbf{W}^\top \mathbf{x} + \mathbf{c}$$

Raw features

v.s.

$$\mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\}$$

nonlinear model: better capacity

Raw features

Learned features ← output of inner function

# Learning XOR function

$$\mathbf{W}^\top \mathbf{x} + \mathbf{c}$$

Raw features

v.s.

$$\mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\}$$

nonlinear model: better capacity

Raw features

Learned features ← output of inner function

If  $\mathbf{W}$  can be learned (determined by training data)

# What makes feedforward network different from linear model

- Nonlinear functions in hidden layers

$$f_3(f_2(f_1(x))) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\}$$


$f_1(x) = W'x + c$



# What makes feedforward network different from linear model

- Nonlinear functions in hidden layers


$$f_3(f_2(f_1(x))) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\}$$


$$f_2(x) = \max(0, x)$$

# What makes feedforward network different from linear model

- Nonlinear functions in hidden layers

$$f_3(f_2(f_1(x))) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\}$$

$$f_3(x) = \mathbf{w}'\mathbf{x}$$


# What makes feedforward network different from linear model

- Nonlinear functions in hidden layers

$$f_3(f_2(f_1(x))) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\}$$

- Q: Why composition makes nonconvexity?

# What makes feedforward network different from linear model

- Nonlinear functions in hidden layers

$$f_3(f_2(f_1(x))) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\}$$

- Q: Why composition makes nonconvexity?

$$f(x) = g(h(x)) \quad \text{where } g(x) = h(x) = x^2$$

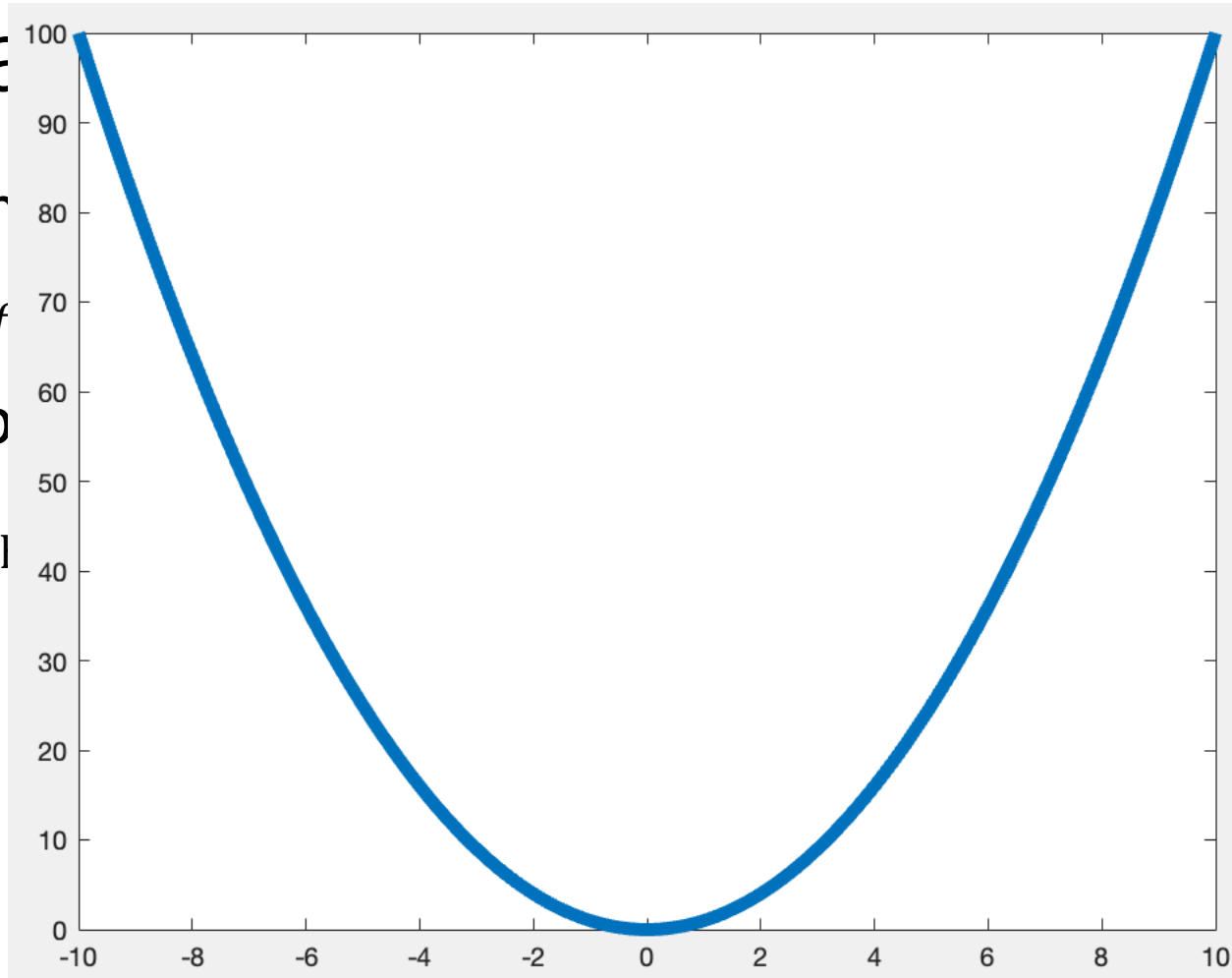
# What makes feedforward network different from linear

- Nonlinear function

$$f_3(f)$$

- Q: Why compute

$$f(x) = g(l)$$



$$f(x) = x^2$$

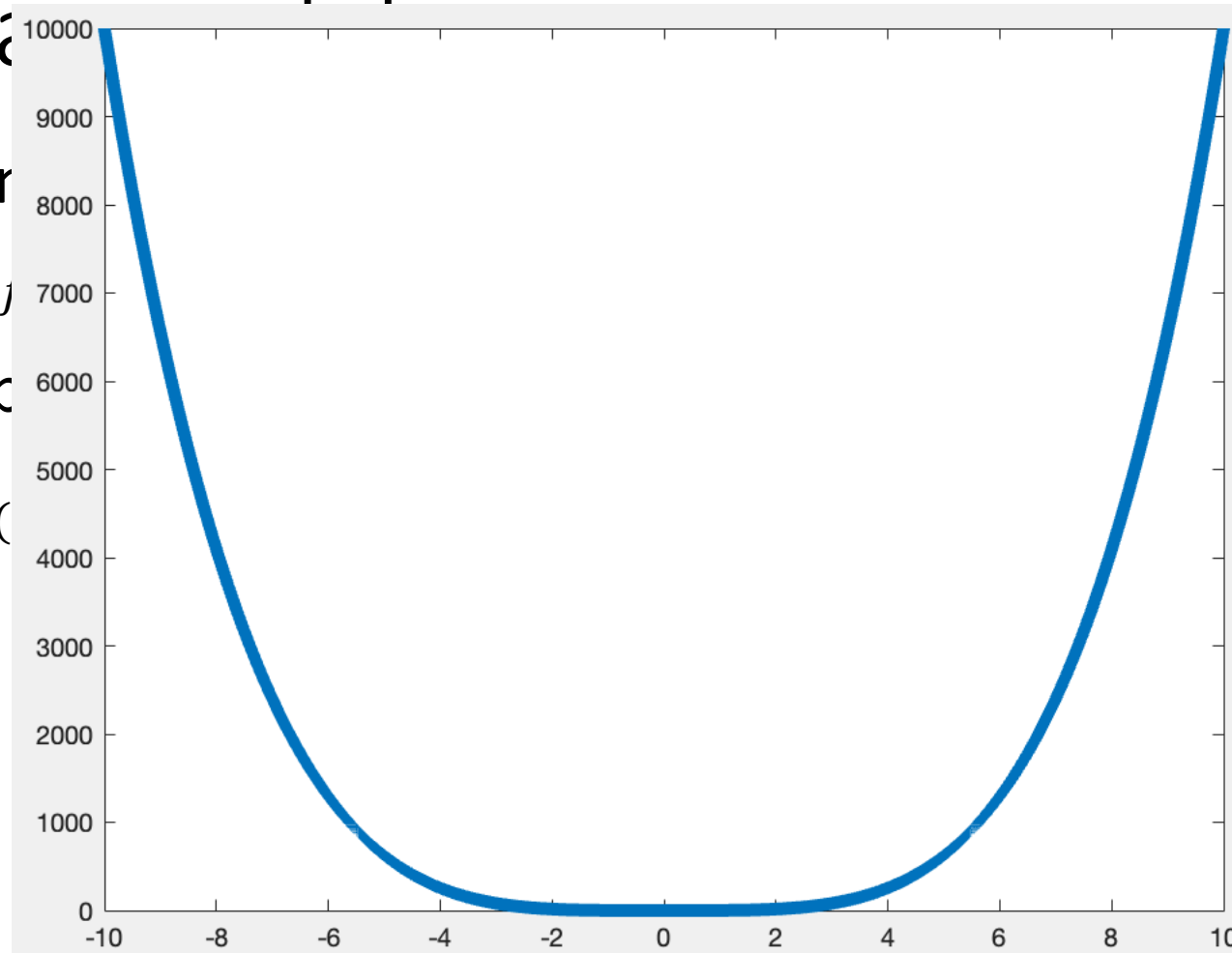
# What makes feedforward network different from linear?

- Nonlinear function

$$f_3(j)$$

- Q: Why compute

$$f(x) = g(h(x))$$



$$f(x) = g(h(x)) \quad \text{where } g(x) = h(x) = x^2$$

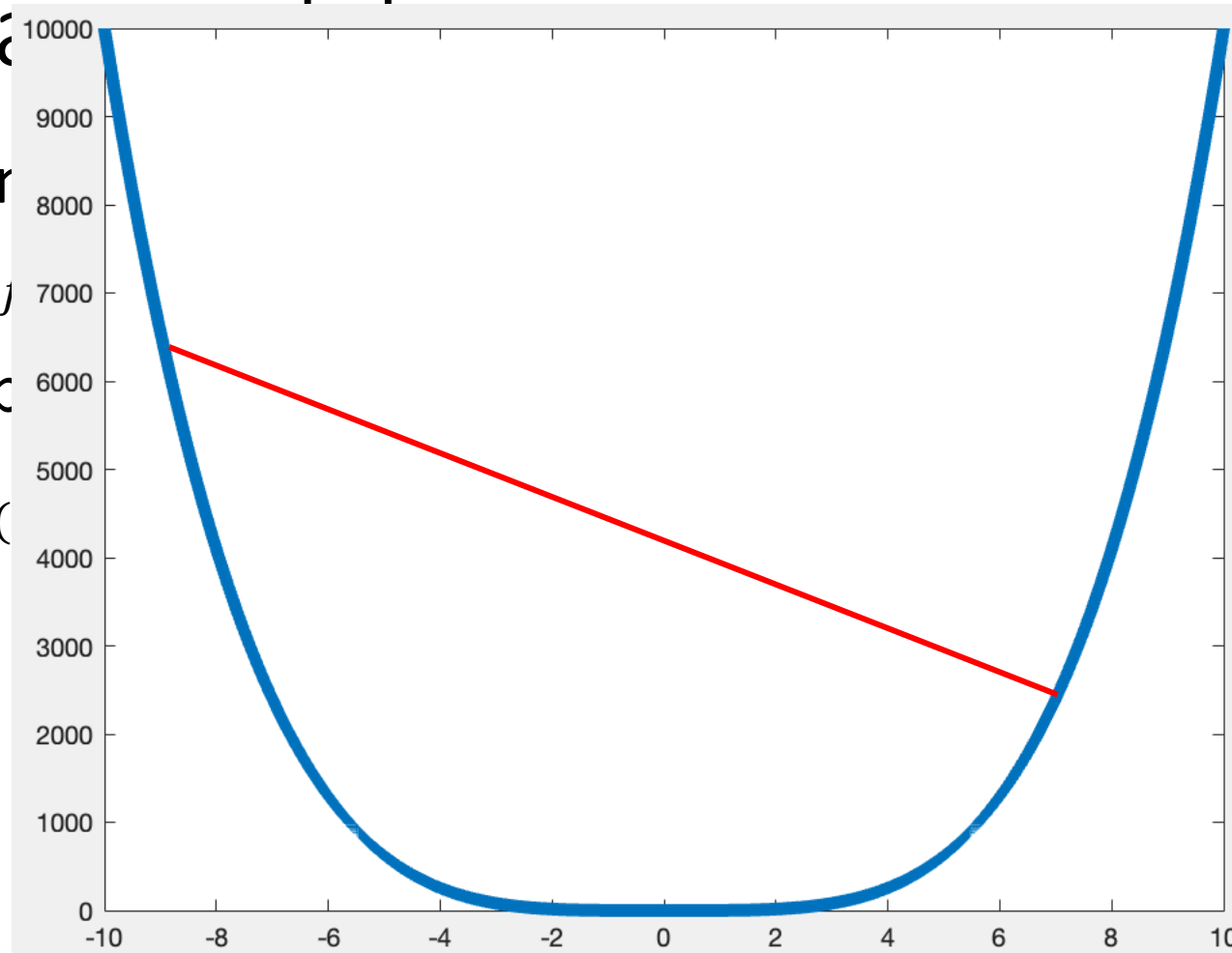
# What makes feedforward network different from linear?

- Nonlinear function

$$f_3(j)$$

- Q: Why compute

$$f(x) = g(h(x))$$

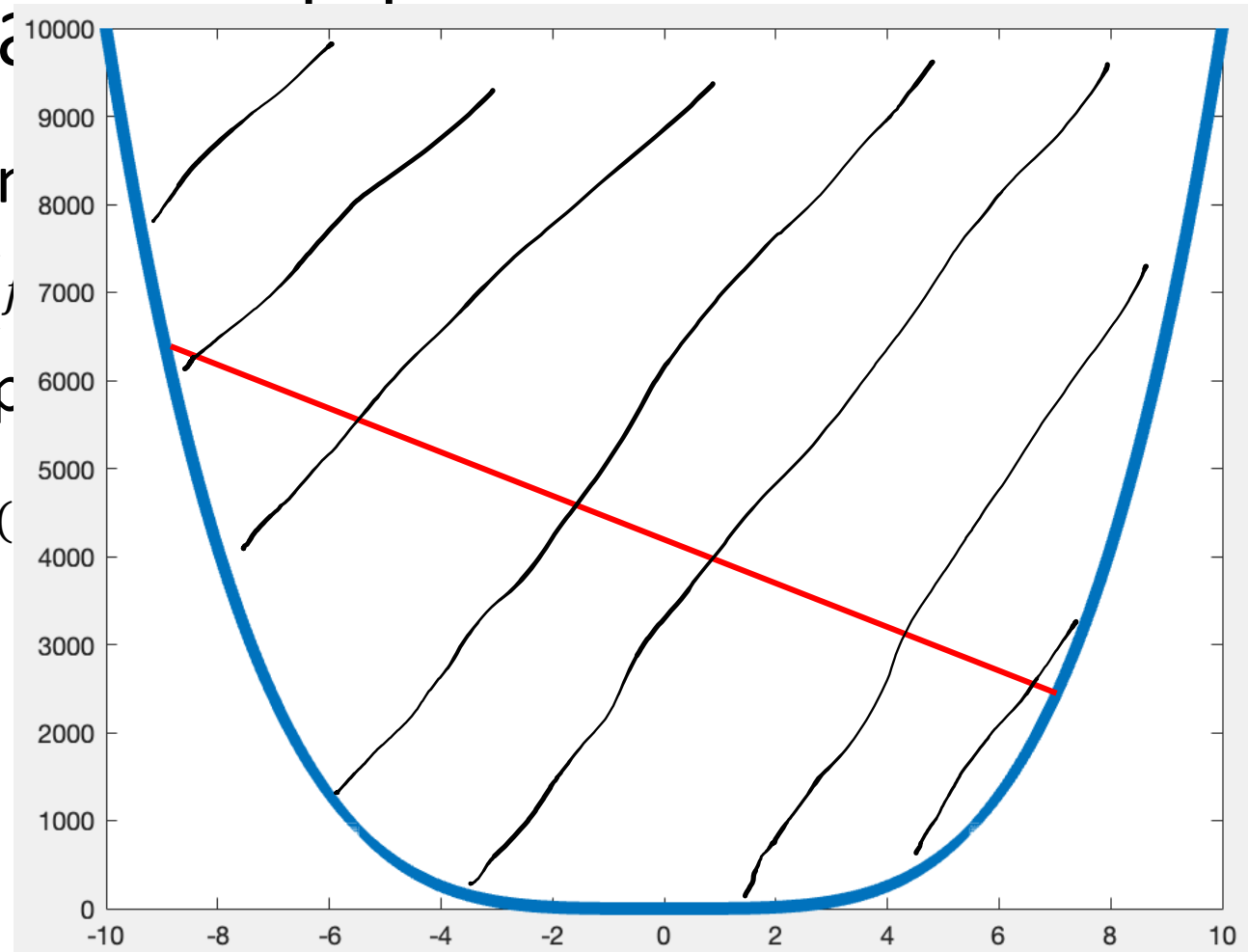


$$f(x) = g(h(x)) \quad \text{where } g(x) = h(x) = x^2$$

# What makes feedforward network different from linear?

- Nonlinear function
- Q: Why compute

$$f_3(j)$$
$$f(x) = g(h(x))$$



$$f(x) = g(h(x)) \quad \text{where } g(x) = h(x) = x^2$$



# What makes feedforward network different from linear model

- Nonlinear functions in hidden layers

$$f_3(f_2(f_1(x))) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\}$$

- Q: Why composition makes nonconvexity?

$$f(x) = g(h(x)) \quad \text{where } g(x) = h(x) = x^2$$

$$f(x) = g(h(x)) \quad \text{where } g(x) = h(x) = \exp(-x)$$

# What makes feedforward network different from linear

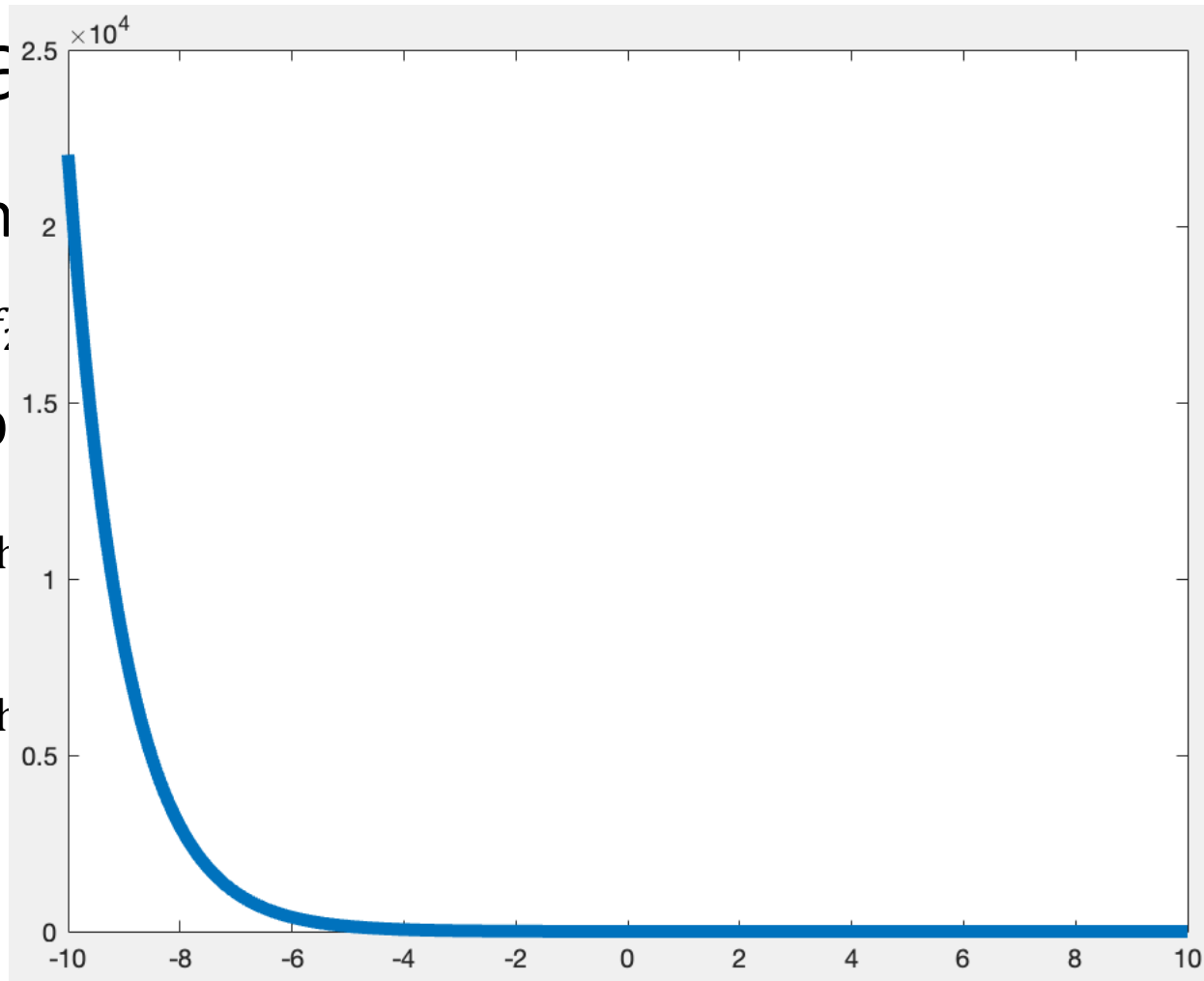
- Nonlinear function

$$f_3(f_2(f_1(x)))$$

- Q: Why compute

$$f(x) = g(h(x))$$

$$f(x) = g(h(x))$$



$$f(x) = \exp(-x)$$

# What makes feedforward network different from linear

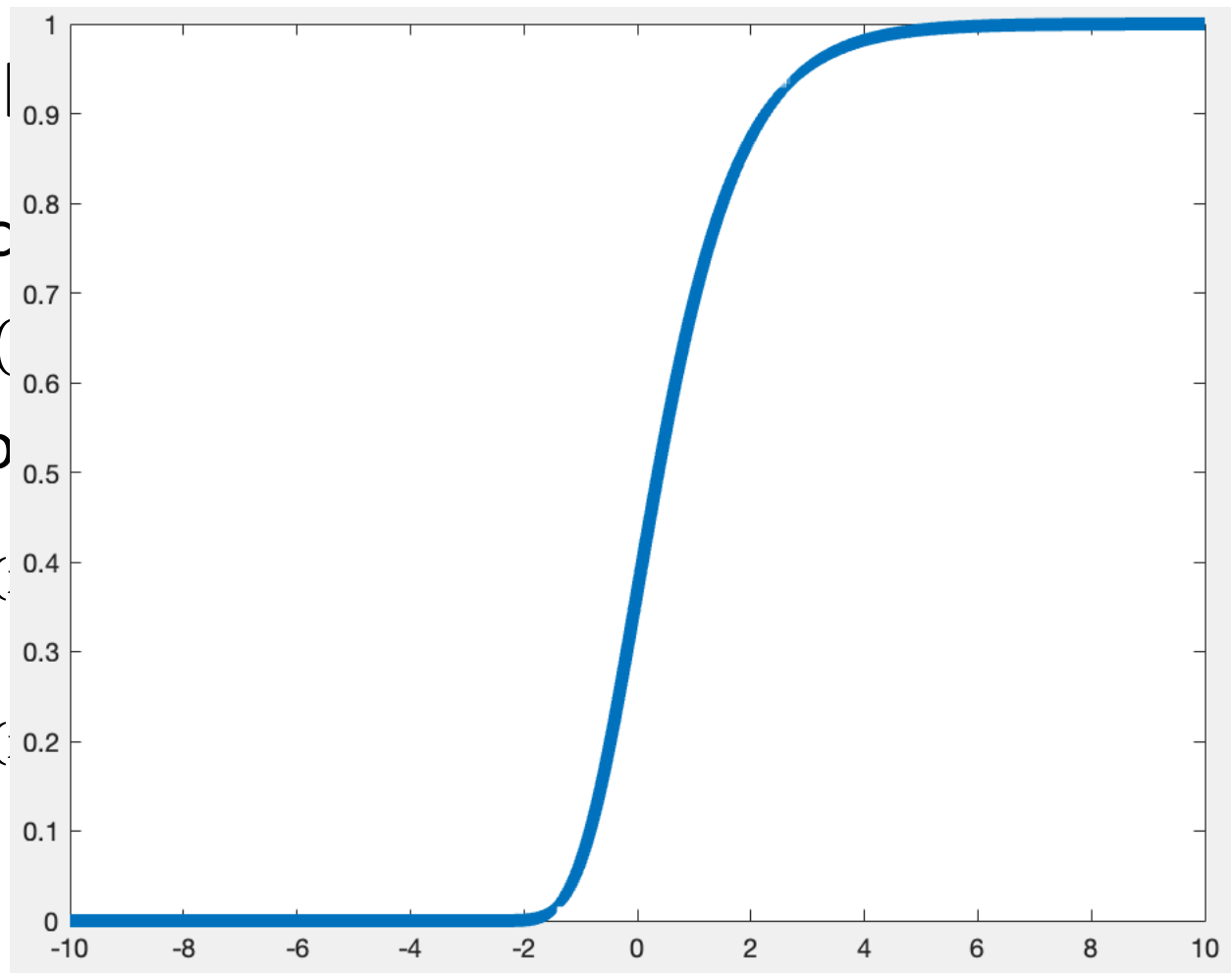
- Nonlinear func

$$f_3(f_2(\dots))$$

- Q: Why compo

$$f(x) = g(h(x))$$

$$f(x) = g(h(x))$$



$$f(x) = g(h(x)) \quad \text{where } g(x) = h(x) = \exp(-x)$$

# What makes feedforward network different from linear model

- Nonlinear functions in hidden layers

$$f_3(f_2(f_1(x))) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\}$$

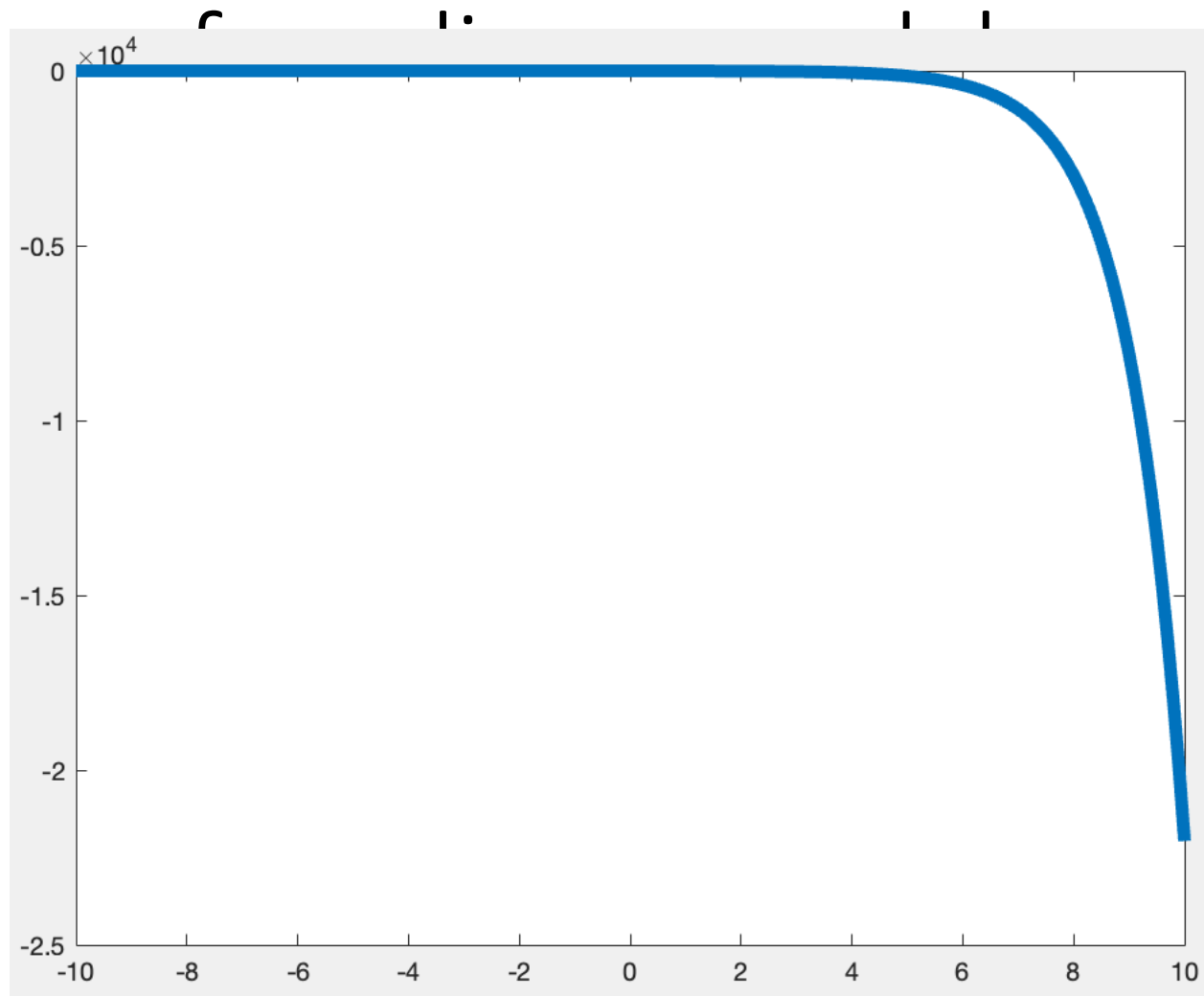
- Q: Why composition makes nonconvexity?

$$f(x) = g(h(x)) \quad \text{where } g(x) = h(x) = x^2$$

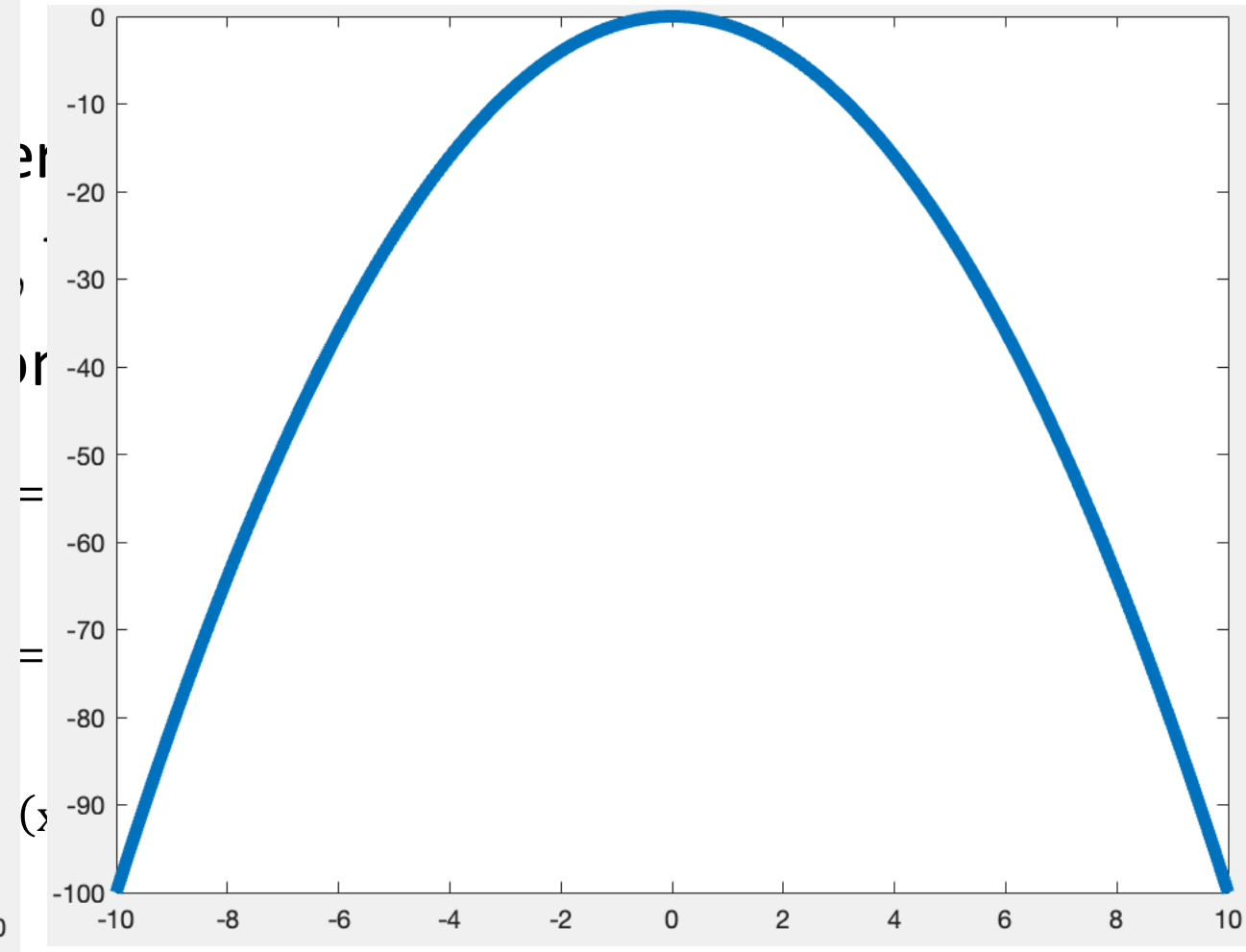
$$f(x) = g(h(x)) \quad \text{where } g(x) = h(x) = \exp(-x)$$

$$f(x) = g(h(x)) \quad \text{where } g(x) = -\exp(x), h(x) = -x^2$$

# What makes feedforward network different



$$g(x) = -\exp(x)$$



$$h(x) = -x^2$$

# What makes feedforward network different from linear model

- Nonlinear function

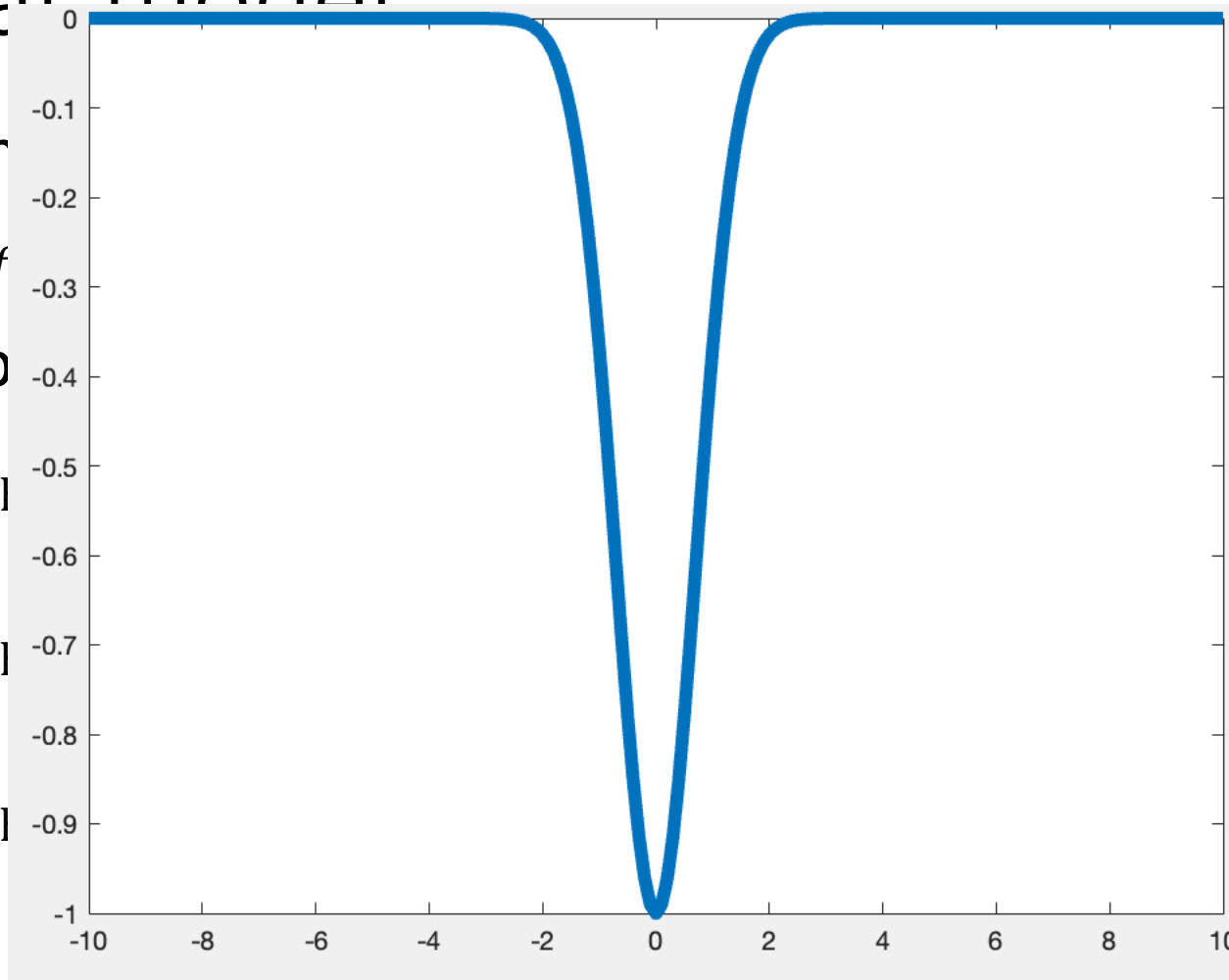
$$f_3(f)$$

- Q: Why compute

$$f(x) = g(l)$$

$$f(x) = g(l)$$

$$f(x) = g(l)$$



$$f(x) = g(h(x)) \quad \text{where } g(x) = -\exp(x), h(x) = -x^2$$

# What makes feedforward network different from linear model

- Nonlinear function

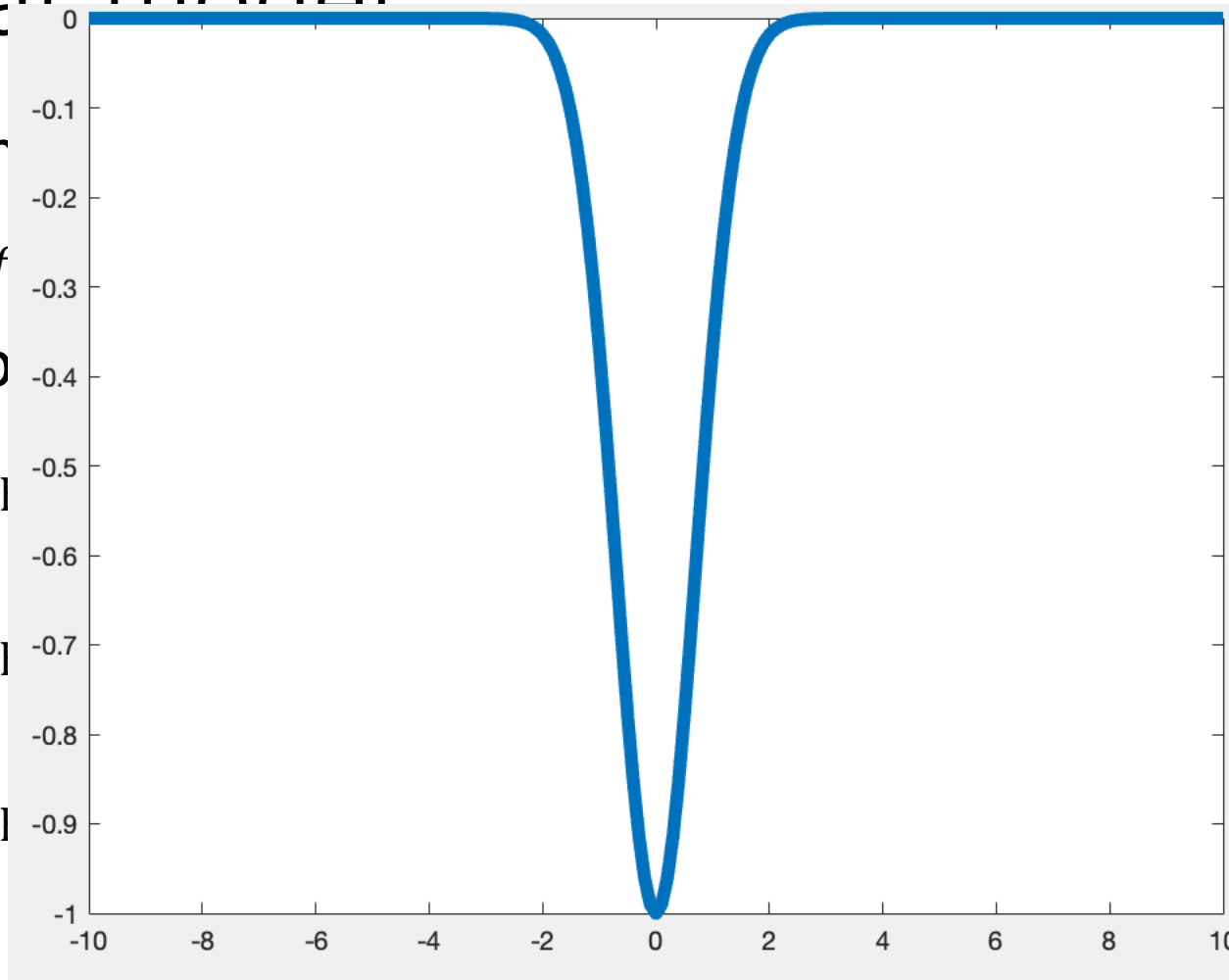
$$f_3(f)$$

- Q: Why compute

$$f(x) = g(l)$$

$$f(x) = g(l)$$

$$f(x) = g(l)$$



$$f(x) = g(h(x)) \quad \text{where } g(x) = -\exp(x), h(x) = -x^2$$

# What makes feedforward network different from linear model

- Nonlinear functions in hidden layers

$$f_3(f_2(f_1(x))) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\}$$

- Q: Why composition makes nonconvexity?

$$f(x) = g(h(x)) \quad \text{where } g(x) = h(x) = x^2$$

$$f(x) = g(h(x)) \quad \text{where } g(x) = h(x) = \exp(-x)$$

$$f(x) = g(h(x)) \quad \text{where } g(x) = -\exp(x), h(x) = -x^2$$

- Special cases



# What makes feedforward network different from linear model

- Nonlinear functions in hidden layers

$$f_3(f_2(f_1(x))) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\}$$

- Q: Why composition makes nonconvexity?

$$f(x) = g(h(x)) \quad \text{where } g(x) = h(x) = x^2$$

$$f(x) = g(h(x)) \quad \text{where } g(x) = h(x) = \exp(-x)$$

$$f(x) = g(h(x)) \quad \text{where } g(x) = -\exp(x), h(x) = -x^2$$

- Special cases

$$f(x) = h(g(x))$$

# What makes feedforward network different from linear model

- Nonlinear functions in hidden layers

$$f_3(f_2(f_1(x))) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\}$$

- Q: Why composition makes nonconvexity?

$f$  is convex if  $h$  is convex and nondecreasing, and  $g$  is convex,  
 $f$  is convex if  $h$  is convex and nonincreasing, and  $g$  is concave,  
 $f$  is concave if  $h$  is concave and nondecreasing, and  $g$  is concave,  
 $f$  is concave if  $h$  is concave and nonincreasing, and  $g$  is convex.

- Special cases

$$f(x) = h(g(x))$$

# Today's class

- What makes feedforward network different from linear model
  - Understanding its structure
- Units for neural networks
  - Output units  $\rightarrow$  cost function
    - $f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$
  - Hidden units
    - $f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$

# Cost function and output units

- How to interact with **groundtruth labels**?
- Likelihood function

$$f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$$

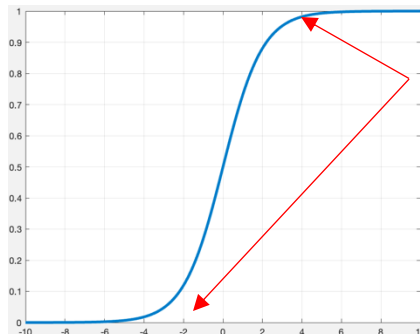
Binary	Gray code	One-hot
000	000	00000001
001	001	00000010
010	011	00000100
011	010	00001000
100	110	00010000
101	111	00100000
110	101	01000000
111	100	10000000

$$L(w) = P_w(X_1 = x_1, \dots, X_n = x_n) = f(w; x_1) \times \dots \times f(w; x_n) = \prod_{i=1}^n f(w; x_i)$$

approximate  
Probability mass function  
(e.g., Bernoulli distribution)

- Maximum likelihood estimation

$$\max_w \log L(w) = \log \prod_{i=1}^n f(w; x_i) = \sum_{i=1}^n \log(f(w; x_i))$$



satisfies probability format

$$f_{sigmoid}(z) = \frac{1}{1 + e^{-z}}$$

$z$ : may be unbounded

# Cost function and output units

$$f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$$

- How to interact with **groundtruth labels**?
- Maximum likelihood estimation

$$p_{\text{data}}(\mathbf{x}) \longleftarrow \text{approximate} \longrightarrow p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$$

# Cost function and output units

$$f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$$

- How to interact with **groundtruth labels**?
- Maximum likelihood estimation

$$p_{\text{data}}(\mathbf{x}) \xleftarrow{\text{approximate}} p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$$

Probability for data



# Cost function and output units

$$f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$$

- How to interact with **groundtruth labels**?
- Maximum likelihood estimation

$$p_{\text{data}}(\mathbf{x}) \xleftarrow{\text{approximate}} p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$$

Probability for data



Q: Where?

One-hot label: dog cat chair  
1 0 0

# Cost function and output units

$$f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$$

- How to interact with **groundtruth labels**?
- Maximum likelihood estimation

$$p_{\text{data}}(\mathbf{x}) \xleftarrow{\text{approximate}} p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$$

Probability for data



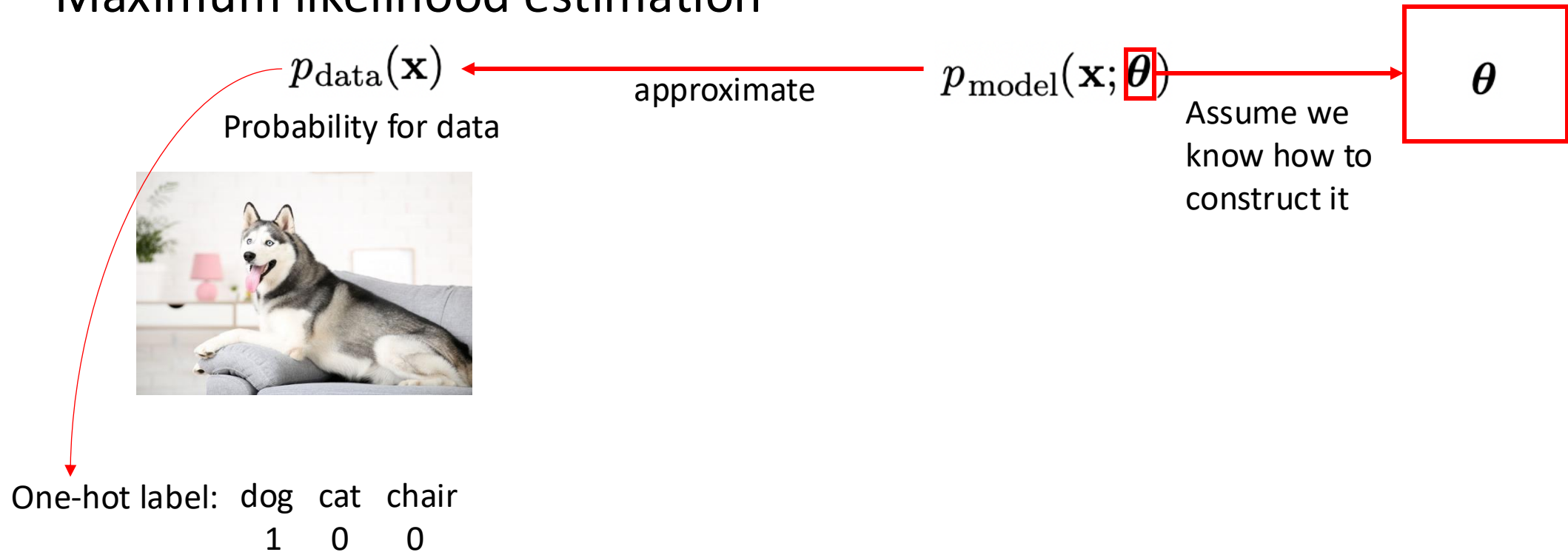
One-hot label: dog cat chair  
1 0 0



# Cost function and output units

$$f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$$

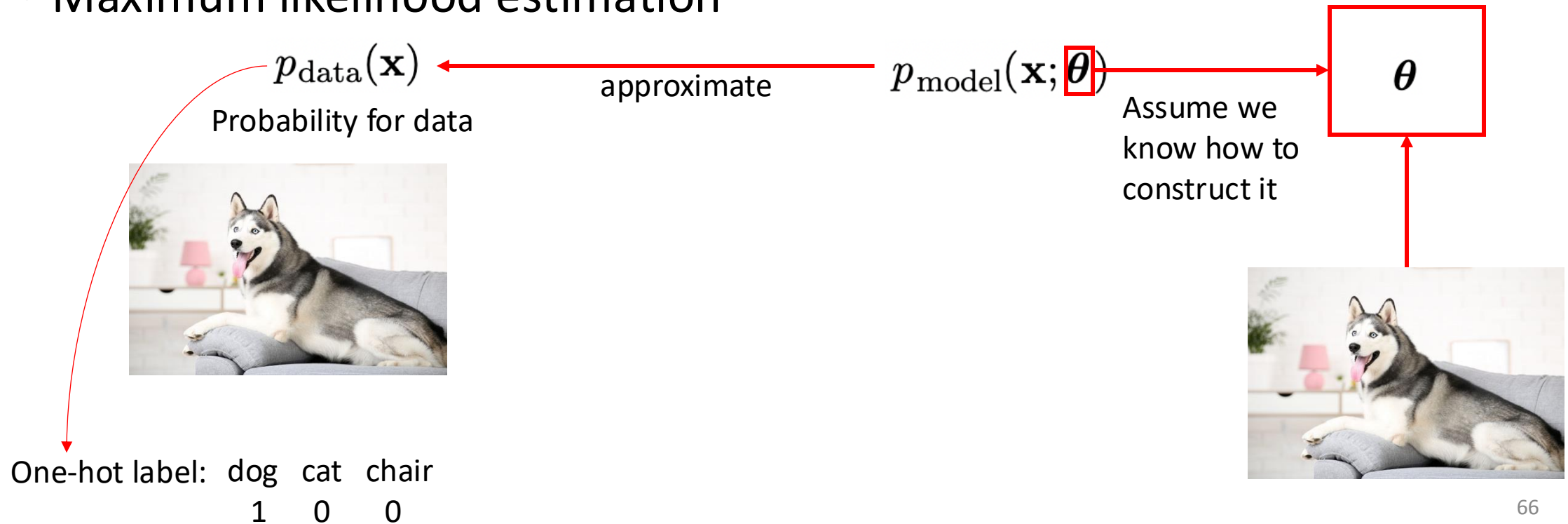
- How to interact with **groundtruth labels**?
- Maximum likelihood estimation



# Cost function and output units

$$f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$$

- How to interact with **groundtruth labels**?
- Maximum likelihood estimation



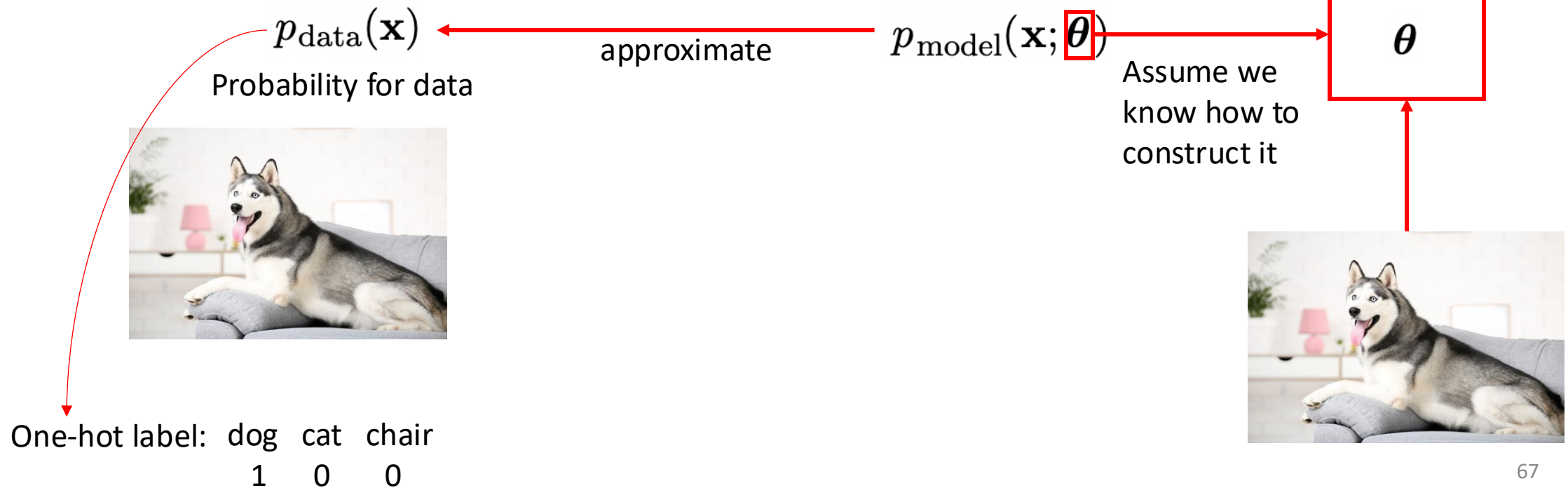
# Cost function and output units

$$f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$$

- How to interact with **groundtruth labels**?
- Maximum likelihood estimation

Probability prediction:

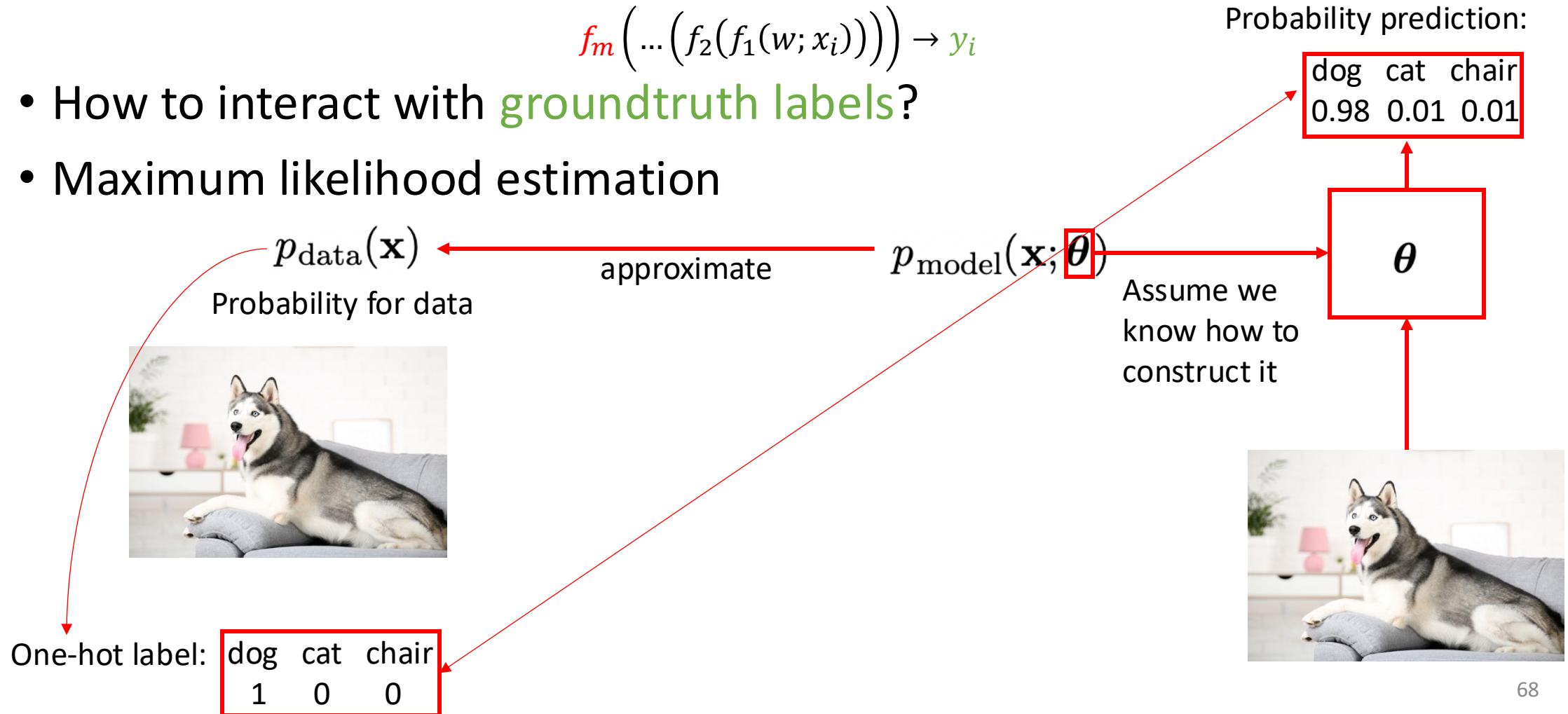
dog	cat	chair
0.98	0.01	0.01



# Cost function and output units

$$f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$$

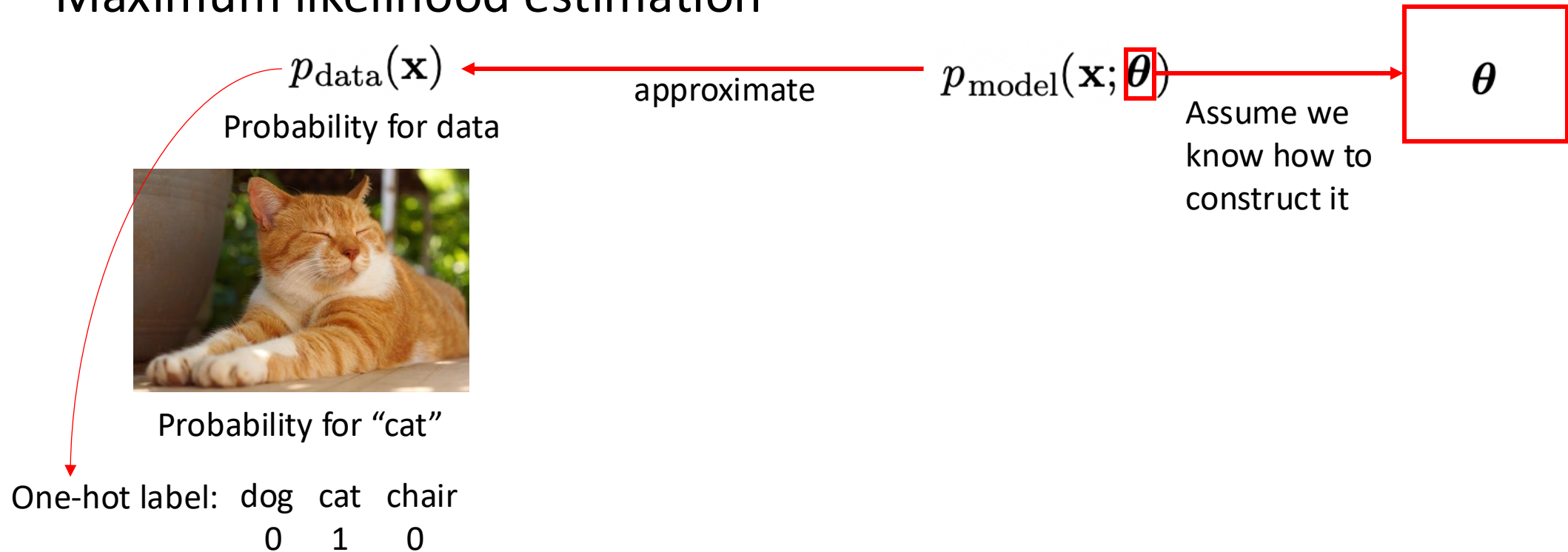
- How to interact with **groundtruth labels**?
- Maximum likelihood estimation



# Cost function and output units

$$f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$$

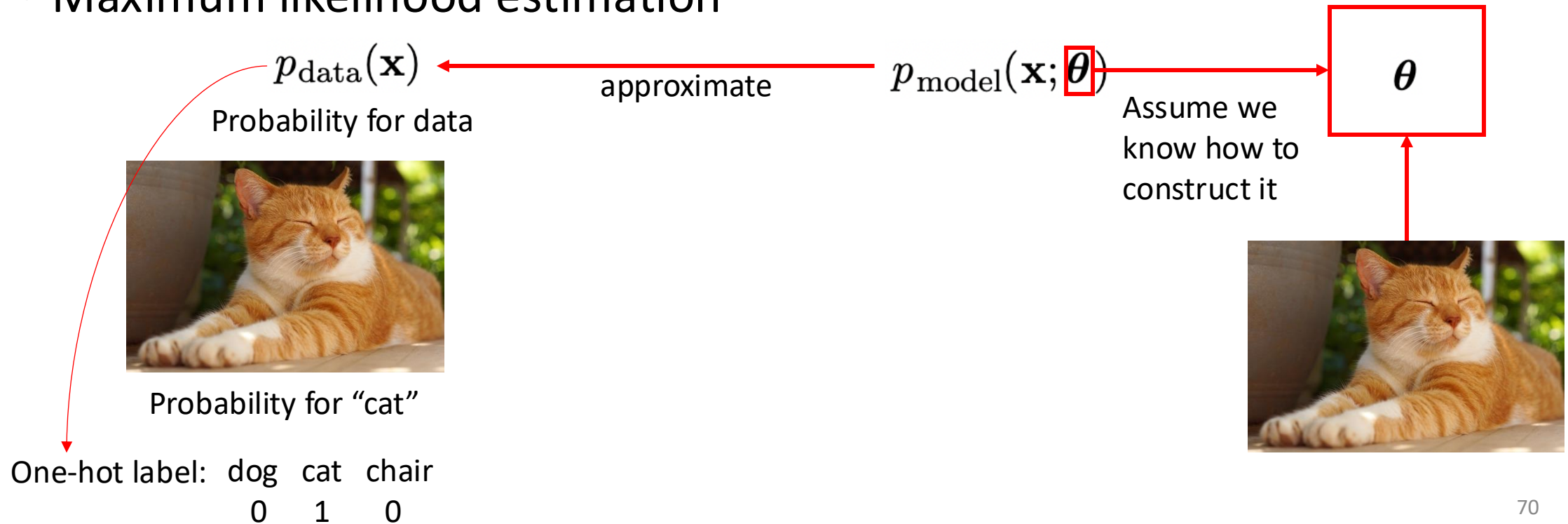
- How to interact with **groundtruth labels**?
- Maximum likelihood estimation



# Cost function and output units

$$f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$$

- How to interact with **groundtruth labels**?
- Maximum likelihood estimation



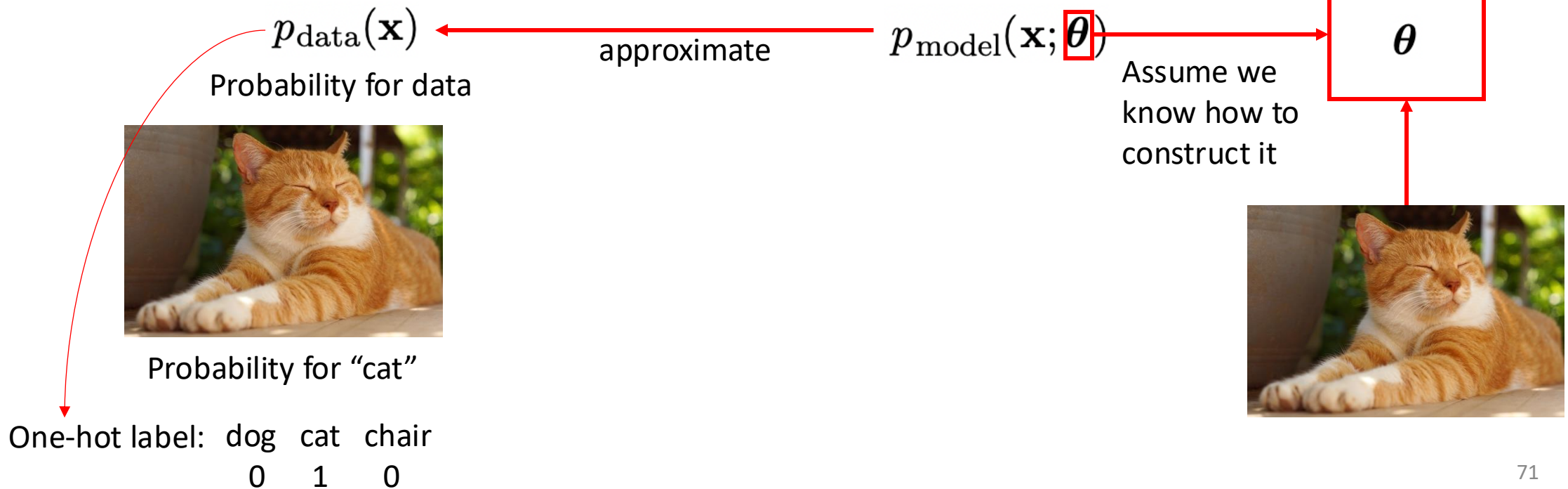
# Cost function and output units

$$f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$$

- How to interact with **groundtruth labels**?
- Maximum likelihood estimation

Probability prediction:

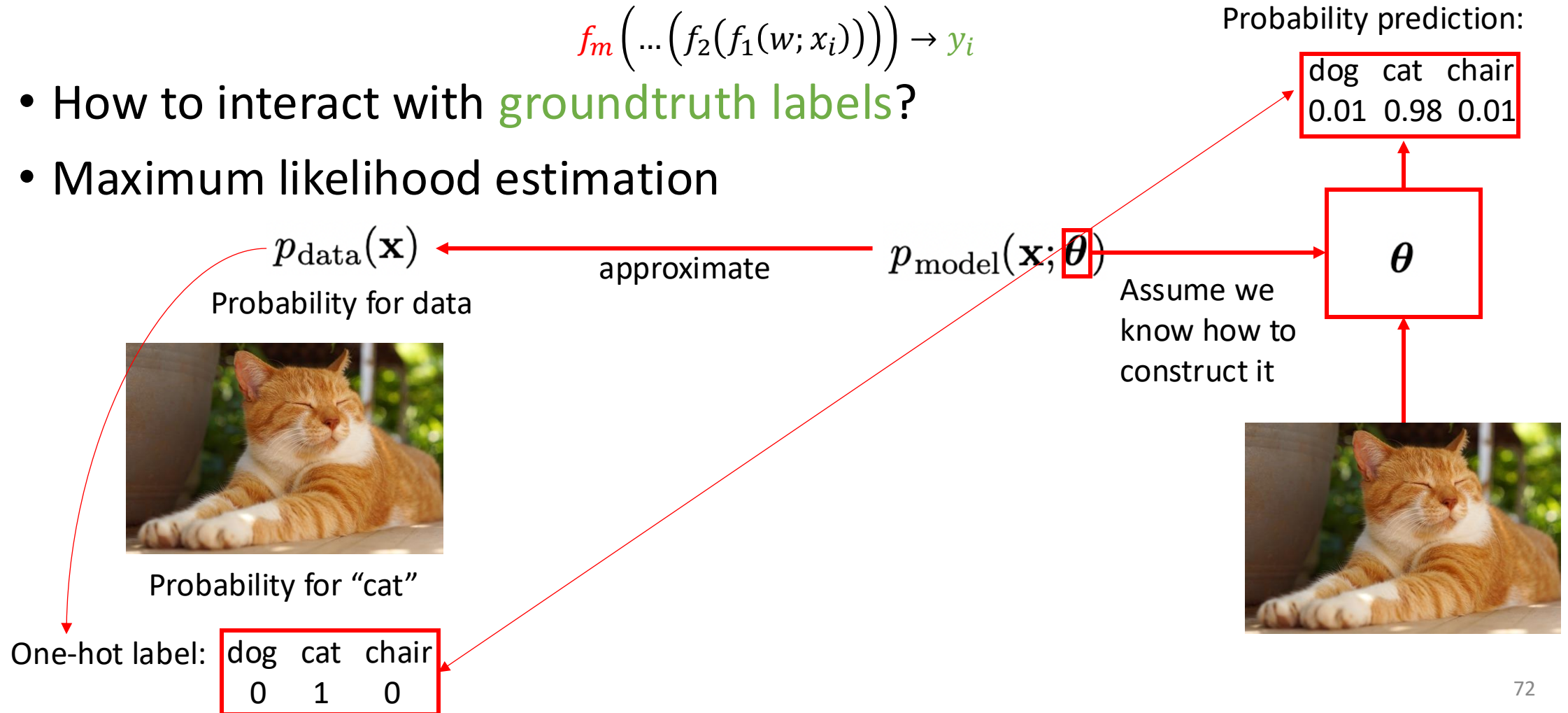
dog	cat	chair
0.01	0.98	0.01



# Cost function and output units

$$f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$$

- How to interact with **groundtruth labels**?
- Maximum likelihood estimation

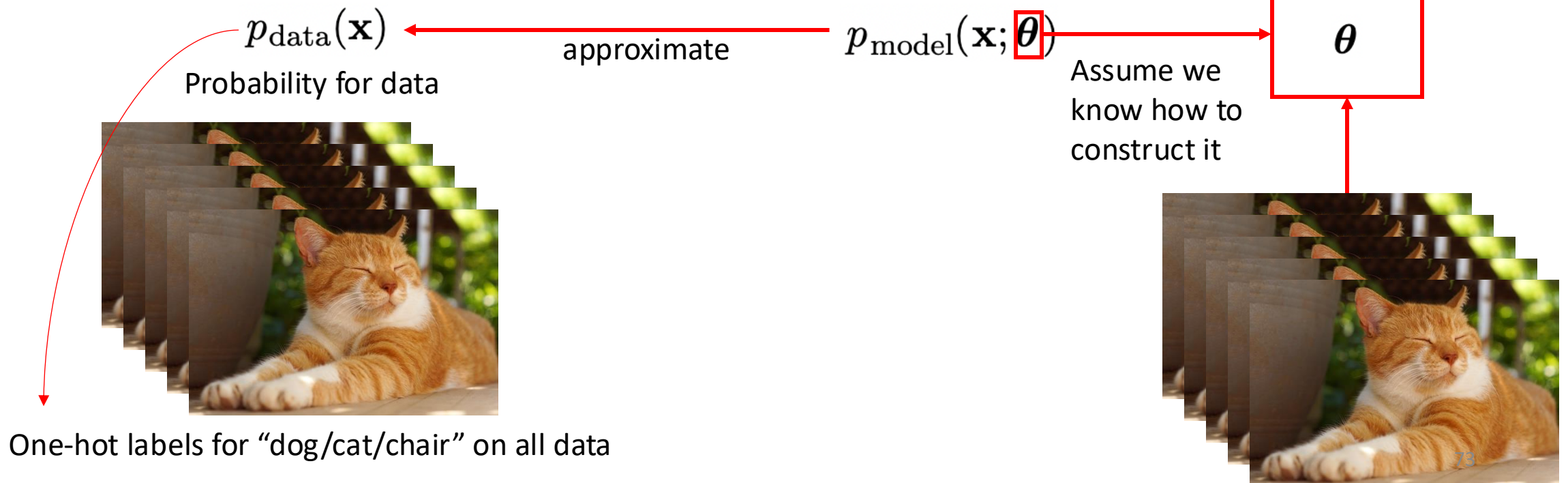




# Cost function and output units

$$f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$$

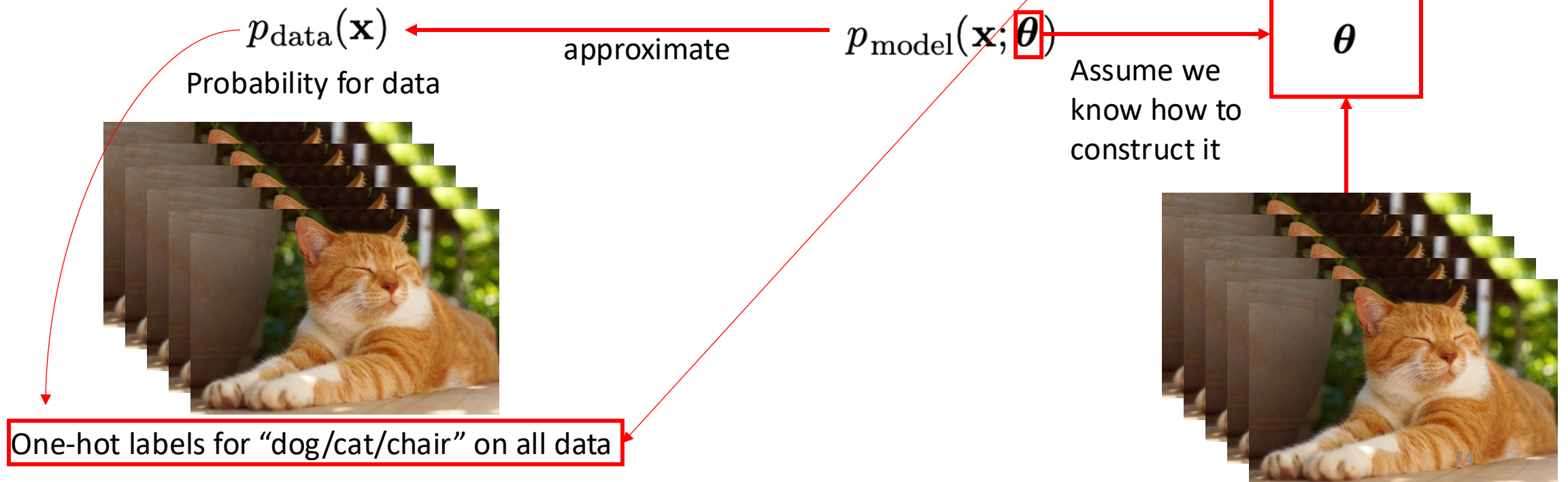
- How to interact with **groundtruth labels**?
- Maximum likelihood estimation



# Cost function and output units

$$f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$$

- How to interact with **groundtruth labels**?
- Maximum likelihood estimation



# Cost function and output units

- MLE and KL divergence

$$\min_{p_{model}} D_{KL}(\hat{p}_{data} || p_{model}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{data}} [\log \hat{p}_{data}(\mathbf{x}) - \log p_{model}(\mathbf{x})]$$

# Cost function and output units

- MLE and KL divergence

Empirical distribution (from training set): cannot scan all possible data  $p_{\text{data}}(\mathbf{x})$

$$\min_{p_{\text{model}}} D_{\text{KL}}(\hat{p}_{\text{data}} \| p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log \hat{p}_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x})]$$

# Cost function and output units

- MLE and KL divergence

Empirical distribution (from training set): cannot scan all possible data  $p_{\text{data}}(\mathbf{x})$

$$\min_{p_{\text{model}}} D_{\text{KL}}(\hat{p}_{\text{data}} || p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log \hat{p}_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x})]$$

If we minimize the KL divergence between the two distributions:

# Cost function and output units

- MLE and KL divergence

Empirical distribution (from training set): cannot scan all possible data  $p_{\text{data}}(\mathbf{x})$

$$\min_{p_{\text{model}}} D_{\text{KL}}(\hat{p}_{\text{data}} || p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log \hat{p}_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x})]$$

If we minimize the KL divergence between the two distributions:

Equivalent to

$$\min_{p_{\text{model}}} - \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x})]$$

# Cost function and output units

- MLE and KL divergence

Empirical distribution (from training set): cannot scan all possible data  $p_{\text{data}}(\mathbf{x})$

$$\min_{p_{\text{model}}} D_{\text{KL}}(\hat{p}_{\text{data}} || p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log \hat{p}_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x})]$$

If we minimize the KL divergence between the two distributions:

Equivalent to

$$\min_{p_{\text{model}}} - \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x})]$$

$$\theta_{\text{ML}} = \arg \max_{\theta} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \theta)$$

# Cost function and output units

- MLE and KL divergence

Empirical distribution (from training set): cannot scan all possible data  $p_{\text{data}}(\mathbf{x})$

$$\min_{p_{\text{model}}} D_{\text{KL}}(\hat{p}_{\text{data}} || p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log \hat{p}_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x})]$$

If we minimize the KL divergence between the two distributions:

Equivalent to

$$\min_{p_{\text{model}}} - \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x})]$$

$$\theta_{\text{ML}} = \arg \max_{\theta} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \theta)$$

$$\Rightarrow \arg \max_{\theta} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{x}; \theta)$$



# Cost function and output units

- MLE and KL divergence

Empirical distribution (from training set): cannot scan all possible data  $p_{\text{data}}(\mathbf{x})$

$$\min_{p_{\text{model}}} D_{\text{KL}}(\hat{p}_{\text{data}} || p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log \hat{p}_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x})]$$

If we minimize the KL divergence between the two distributions:

Equivalent to

$$\min_{p_{\text{model}}} - \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x})]$$

$$\theta_{\text{ML}} = \arg \max_{\theta} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \theta)$$

Minimizing KL divergence

=  
MLE

$$\Rightarrow \arg \max_{\theta} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{x}; \theta)$$

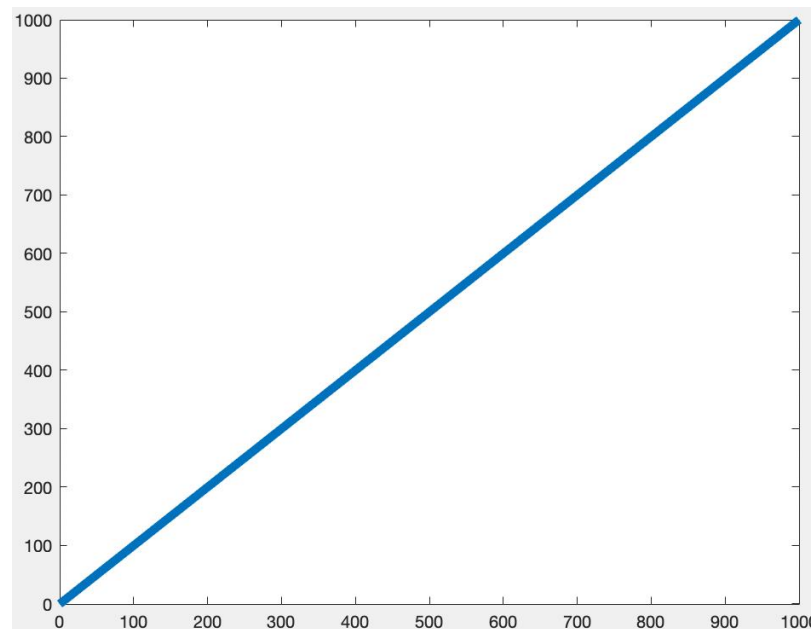
# What are hidden units?

$$f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$$



# What are hidden units?

$$f_m \left( \dots \left( f_2 \left( f_1(w; x_i) \right) \right) \right) \rightarrow y_i$$

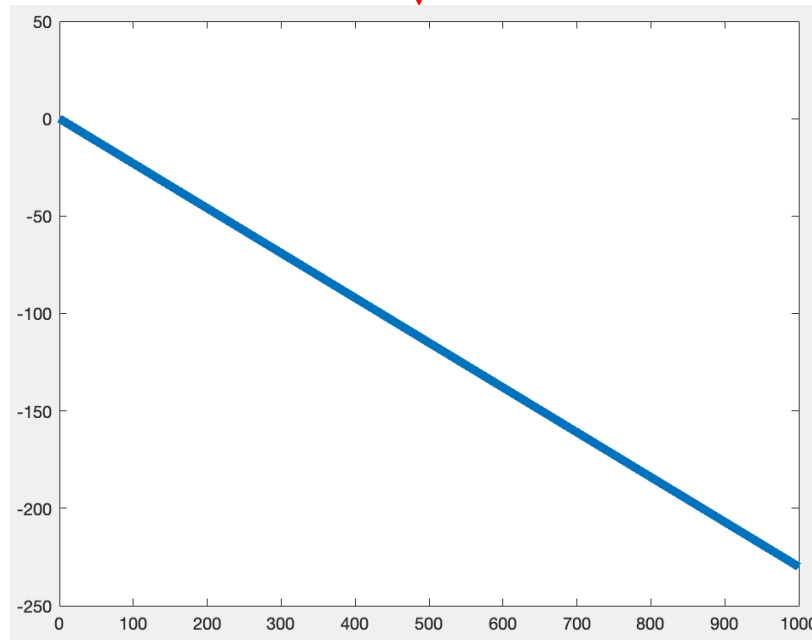


Q: what if for all layers:  
 $f(x) = a * x$ ?

$$f(x) = x$$

# What are hidden units?

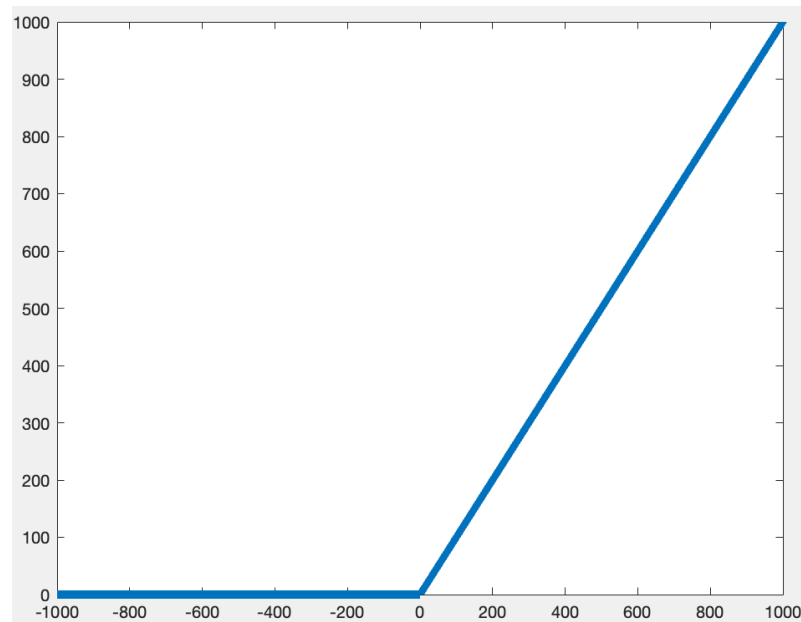
$$f(x) = a * b * c * d * x$$



Combination of **all linear** layers is still linear  
We are interested in nonlinear layers

# ReLU (Rectified Linear Unit)

Activation function

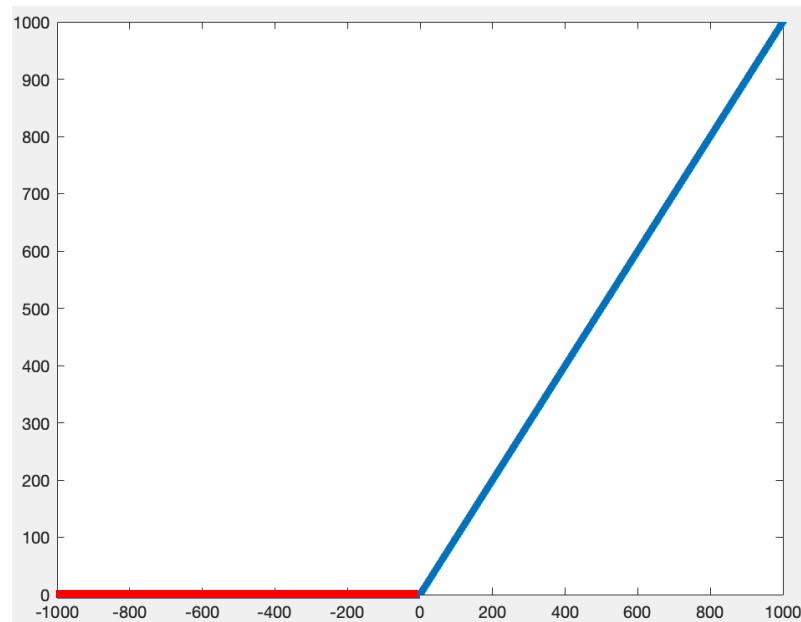


$$f(x) = \max(0, x)$$

# ReLU (Rectified Linear Unit)

- Dying ReLU issue

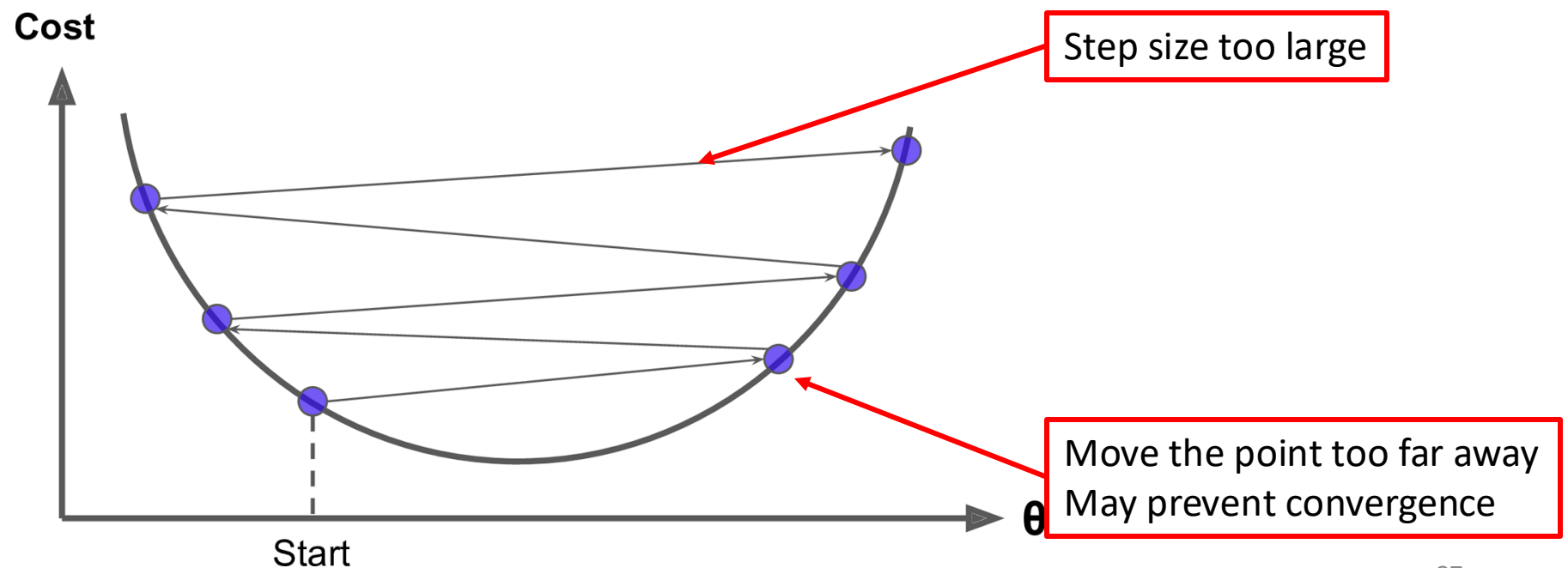
Activation function



$$f(x) = \max(0, x)$$

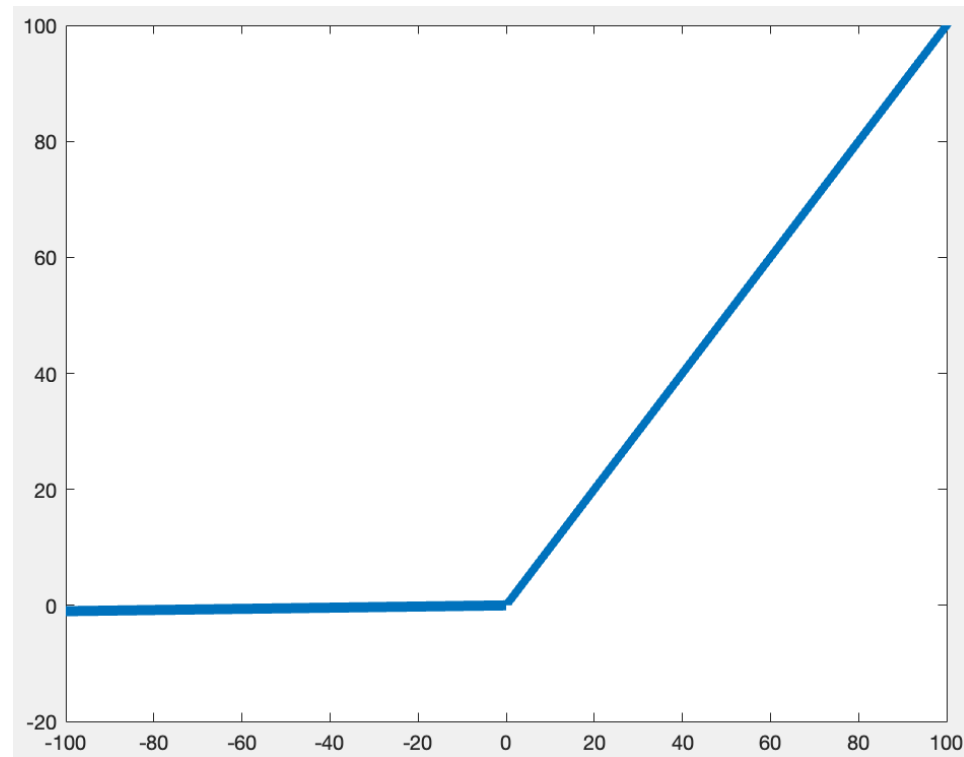
# Determining model parameters

- When to terminate GD (determining T)?
  - Main factors influencing convergence rate?
    - Step size (learning rate)



# Leaky ReLU

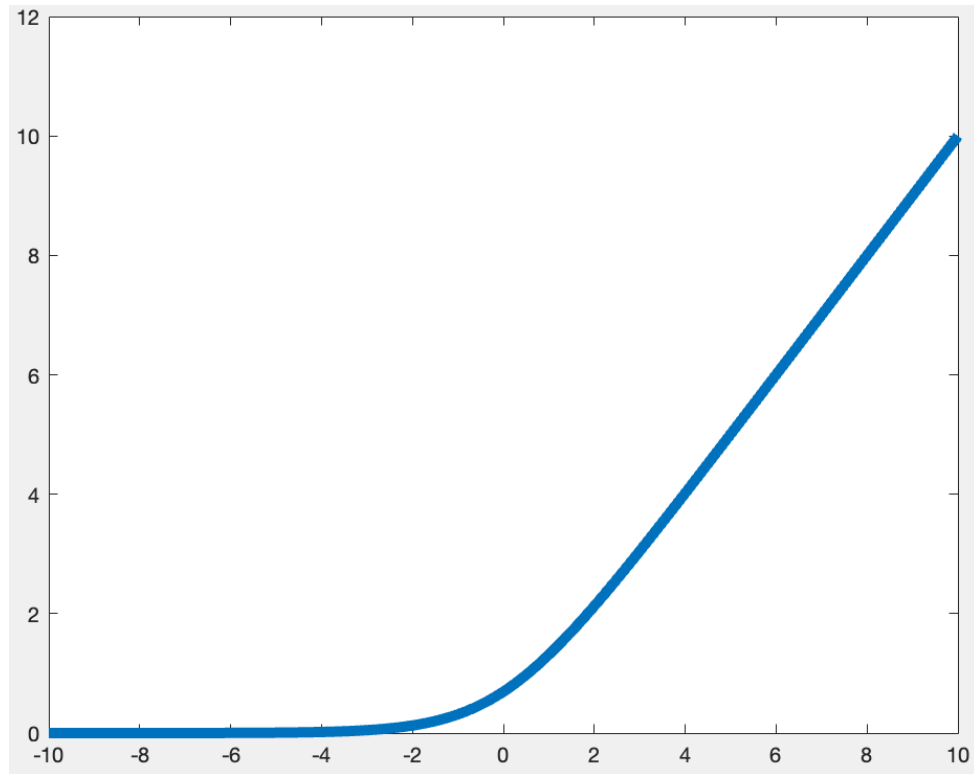
$$f(x_j^i) = \max(0.01x_j^i, x_j^i)$$





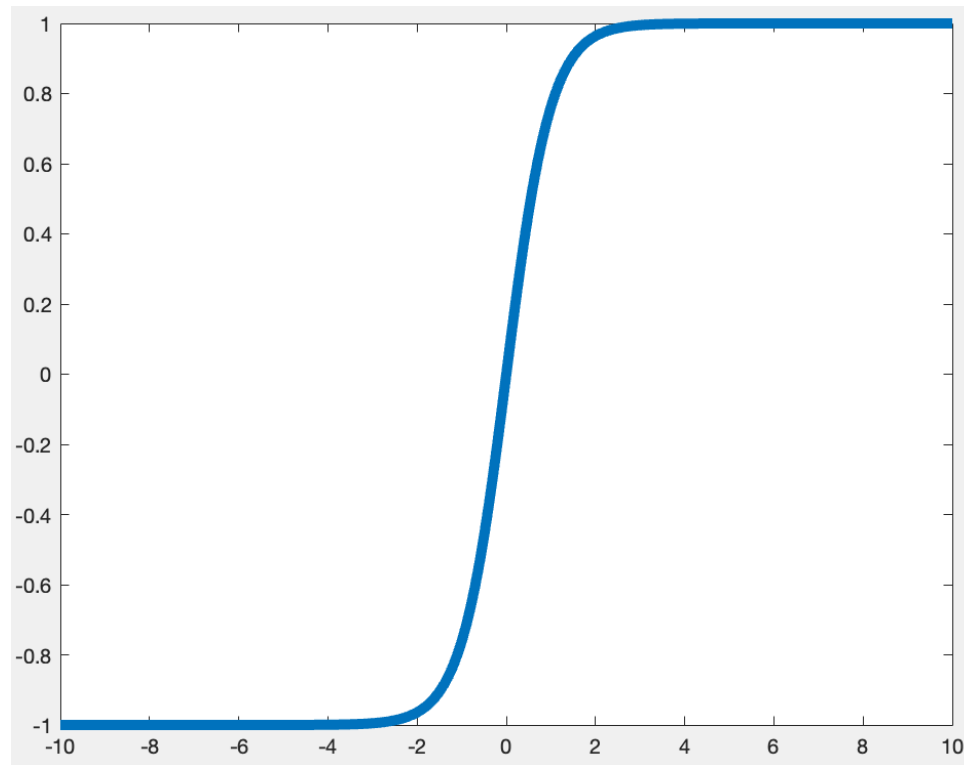
# Smooth ReLU/softplus

$$a_j^i = f(x_j^i) = \log(1 + \exp(x_j^i))$$



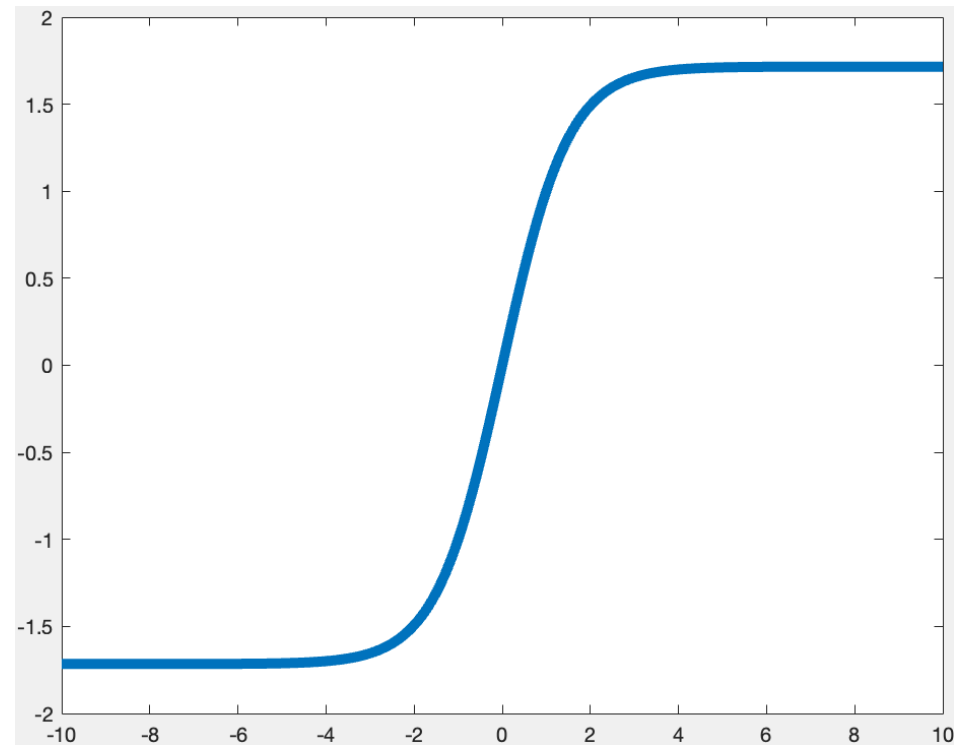
# Tanh

$$f(x_j^i) = \tanh(x_j^i)$$



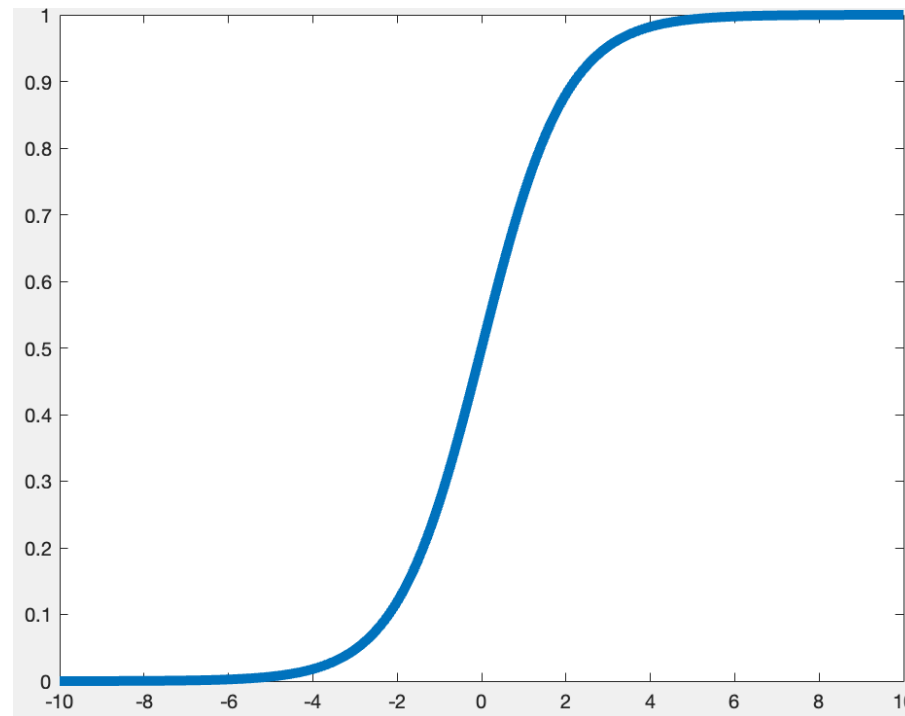
# LeCun's Tanh [1]

$$f(x_j^i) = 1.7159 \tanh\left(\frac{2}{3}x_j^i\right)$$



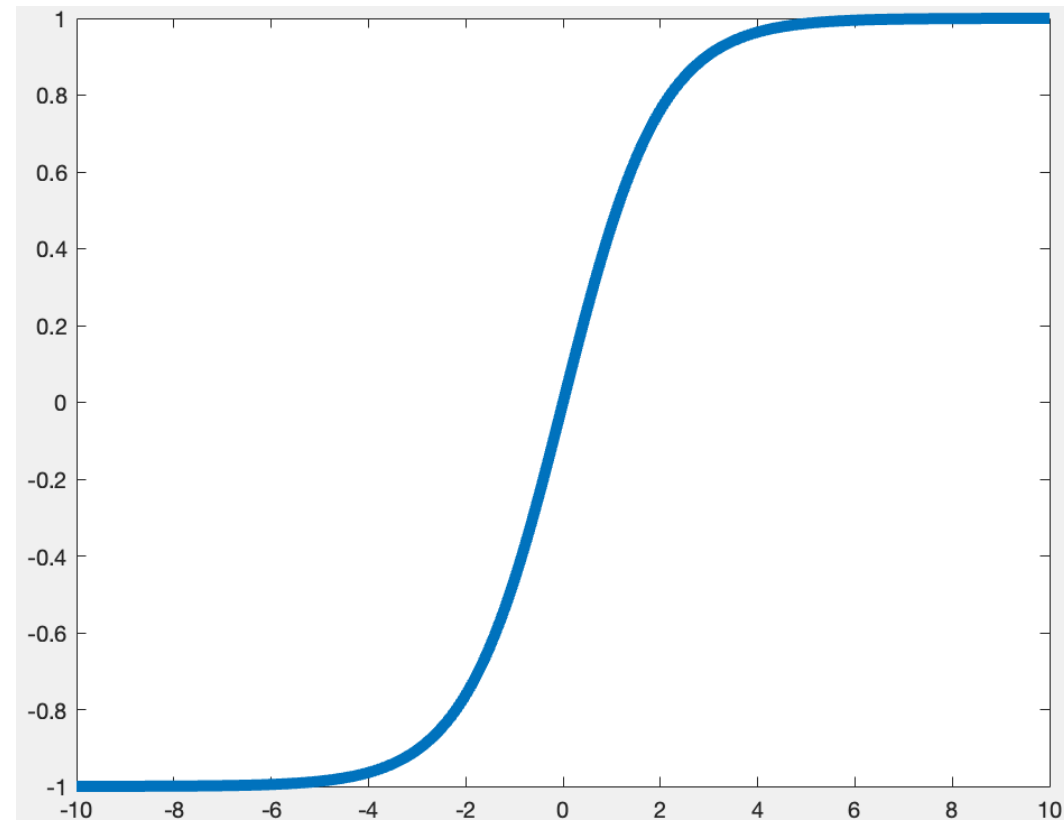
# Sigmoid

$$f(x_j^i) = \frac{1}{1 + \exp(-x_j^i)}$$



# Bipolar sigmoid

$$f(x_j^i) = \frac{1 - \exp(-x_j^i)}{1 + \exp(-x_j^i)}$$



# References

- [1] LeCun, Yann A., Léon Bottou, Genevieve B. Orr, and Klaus-Robert Müller. "Efficient backprop." In *Neural networks: Tricks of the trade*, pp. 9-48. Springer, Berlin, Heidelberg, 2012.

Pdf online: <http://yann.lecun.com/exdb/publis/pdf/lecun-98b.pdf>