

Composition of functions

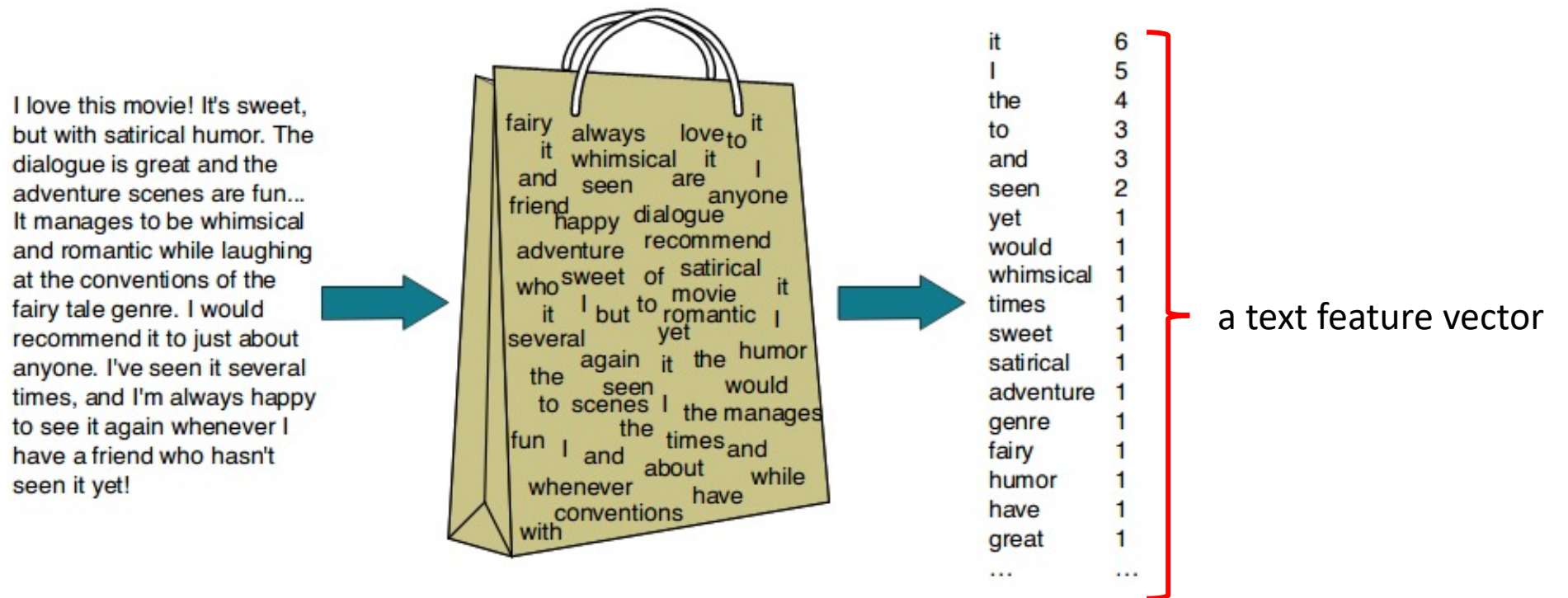
CPT_S 434/534 Neural network design and application

In last class

- Bag-of-words features (hand-crafted)
- History of convolutional neural networks
- Feedforward networks

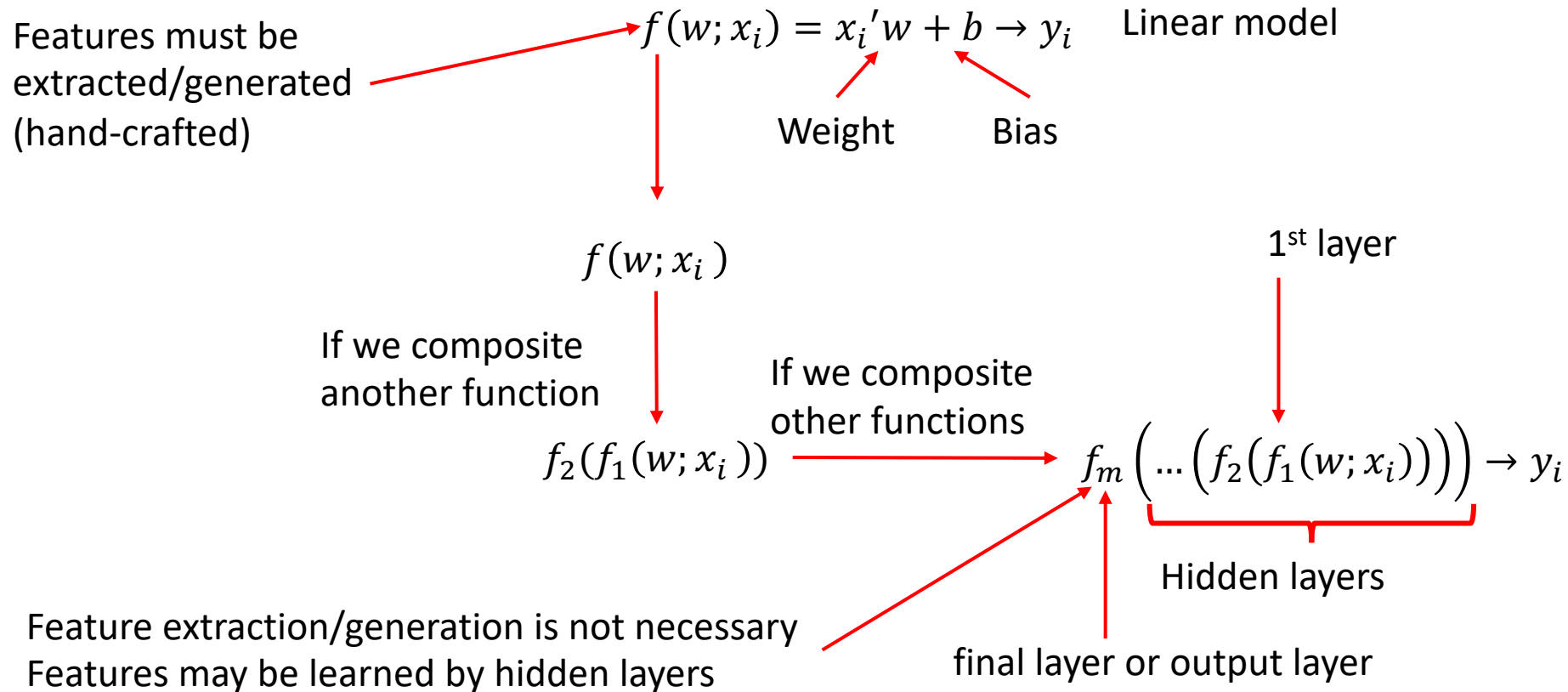
Bag-of-words features

- TF-IDF (term frequency–inverse document frequency)



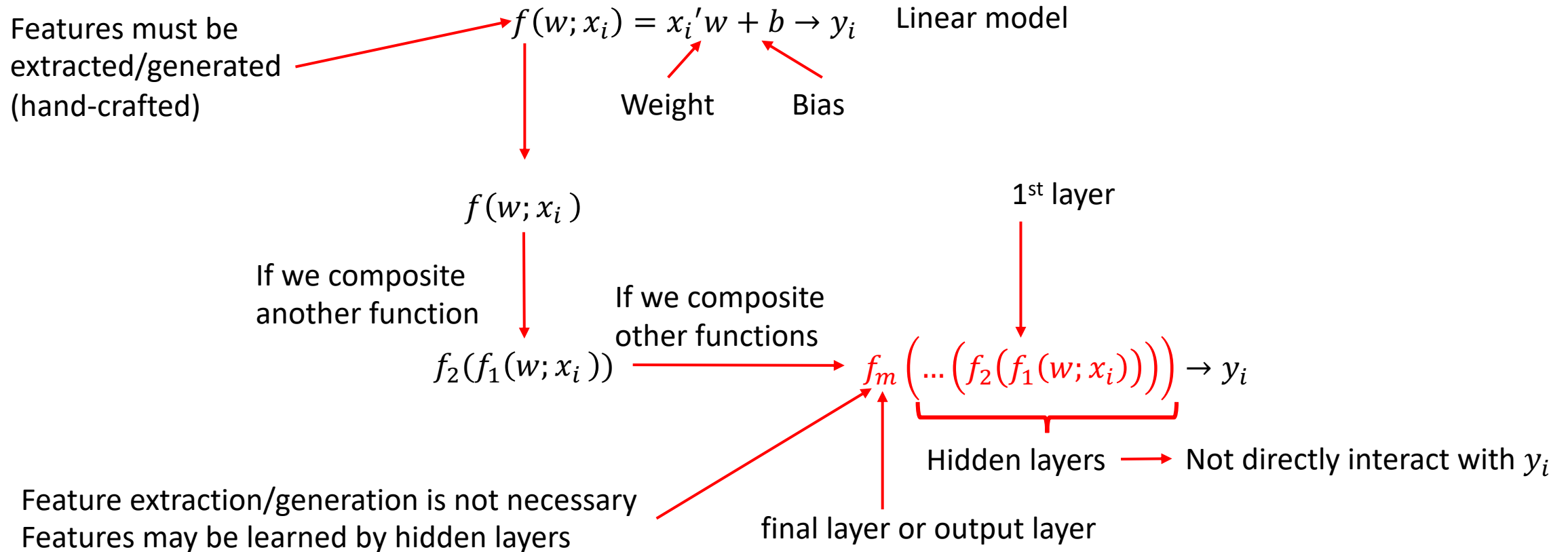
Feedforward networks

- Or multilayer perceptrons (MLPs)



Feedforward networks

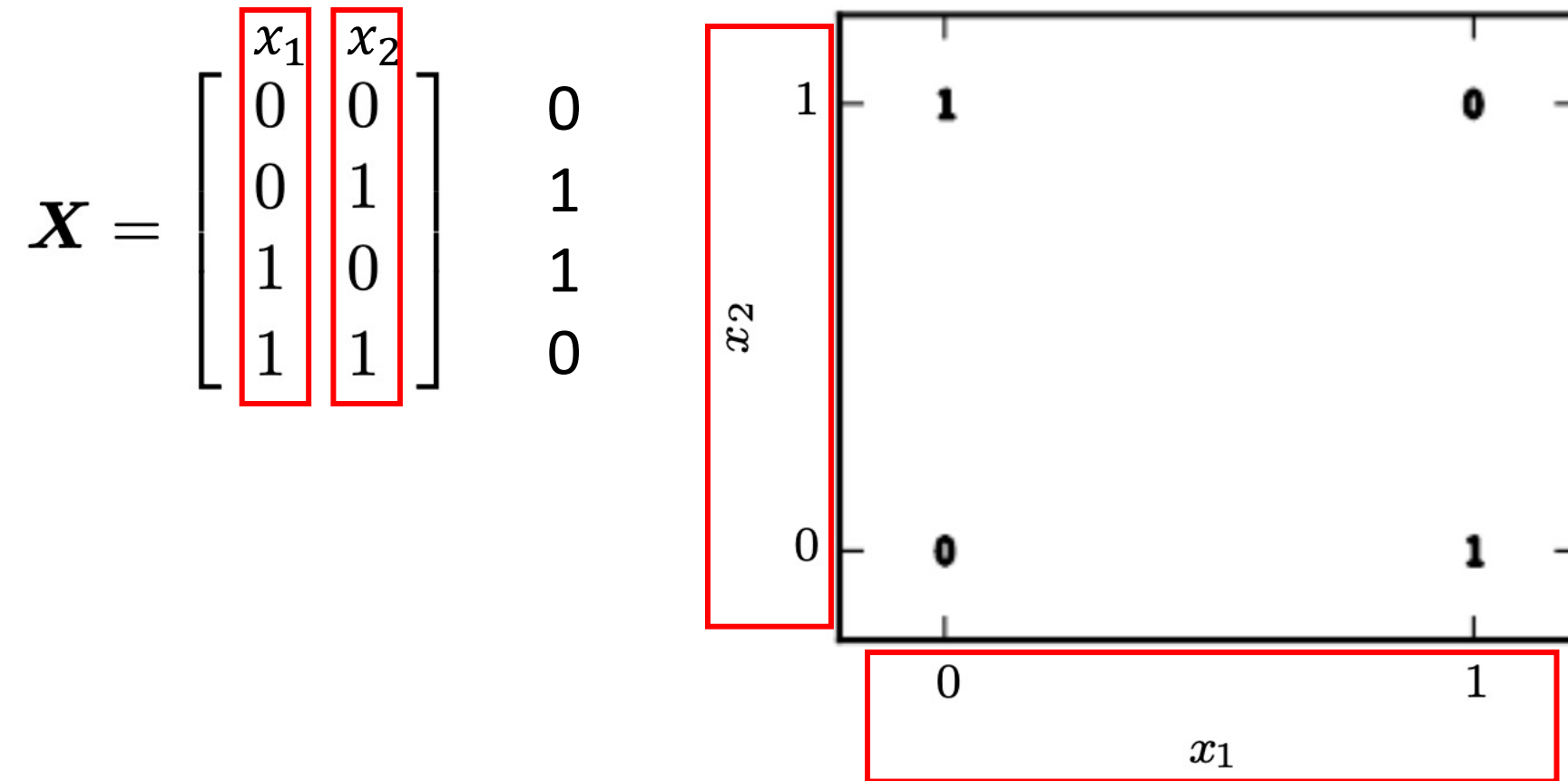
- Or multilayer perceptrons (MLPs)



Today's class

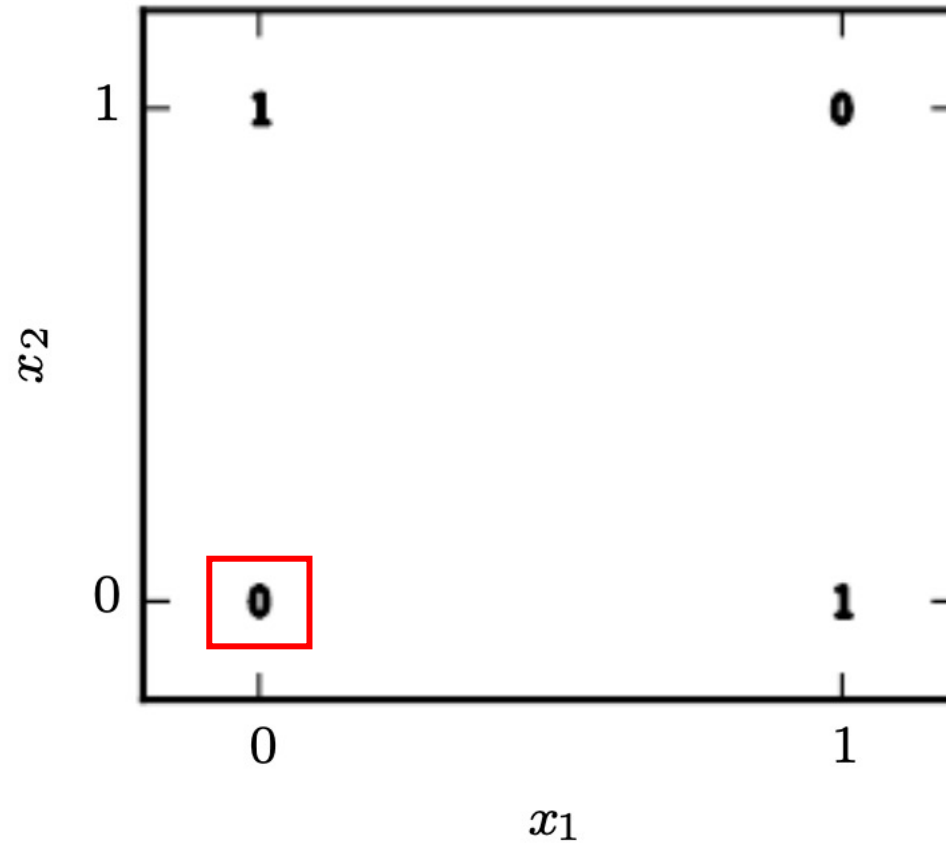
- What makes feedforward network different from linear model
 - Understanding its structure

Learning XOR function



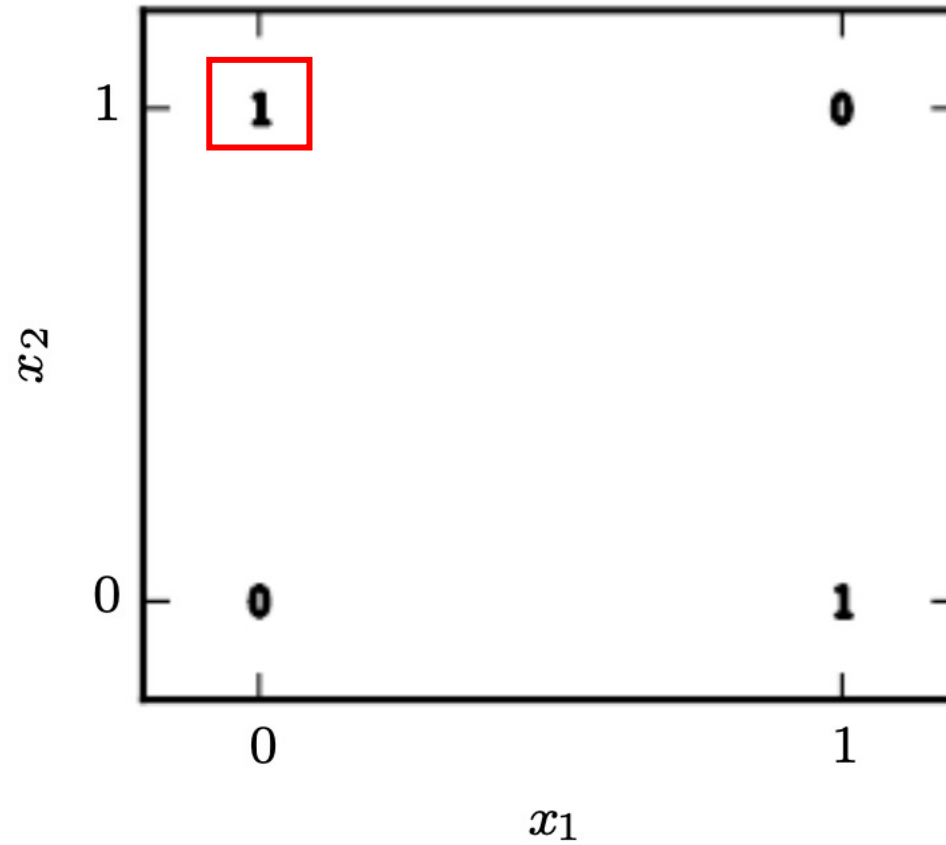
Learning XOR function

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



Learning XOR function

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

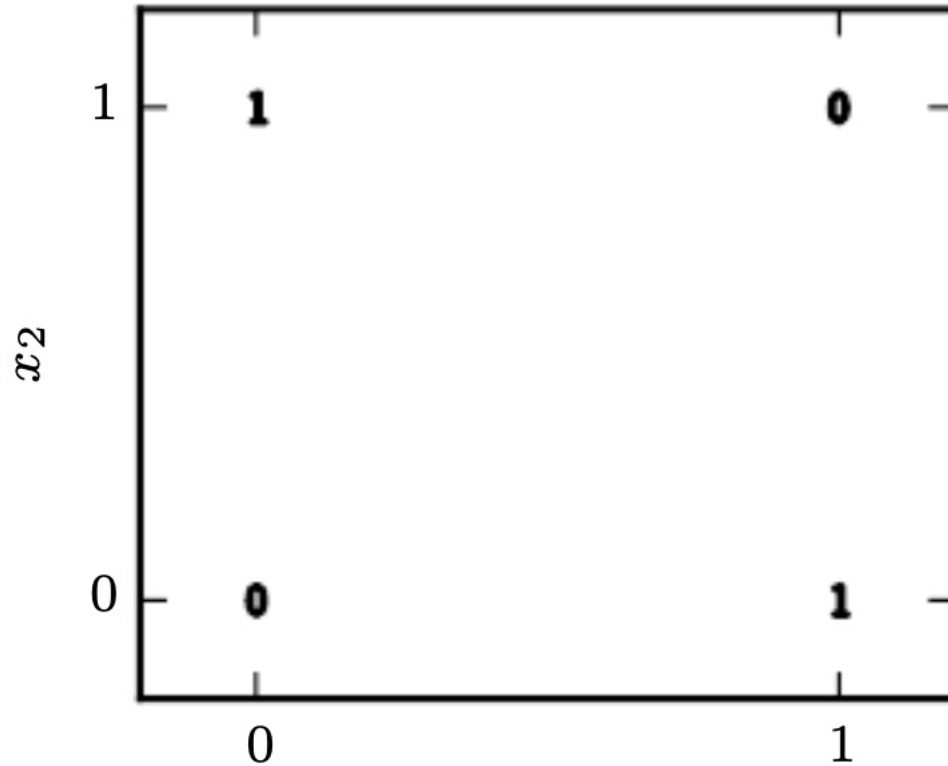


Learning XOR function

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

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

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

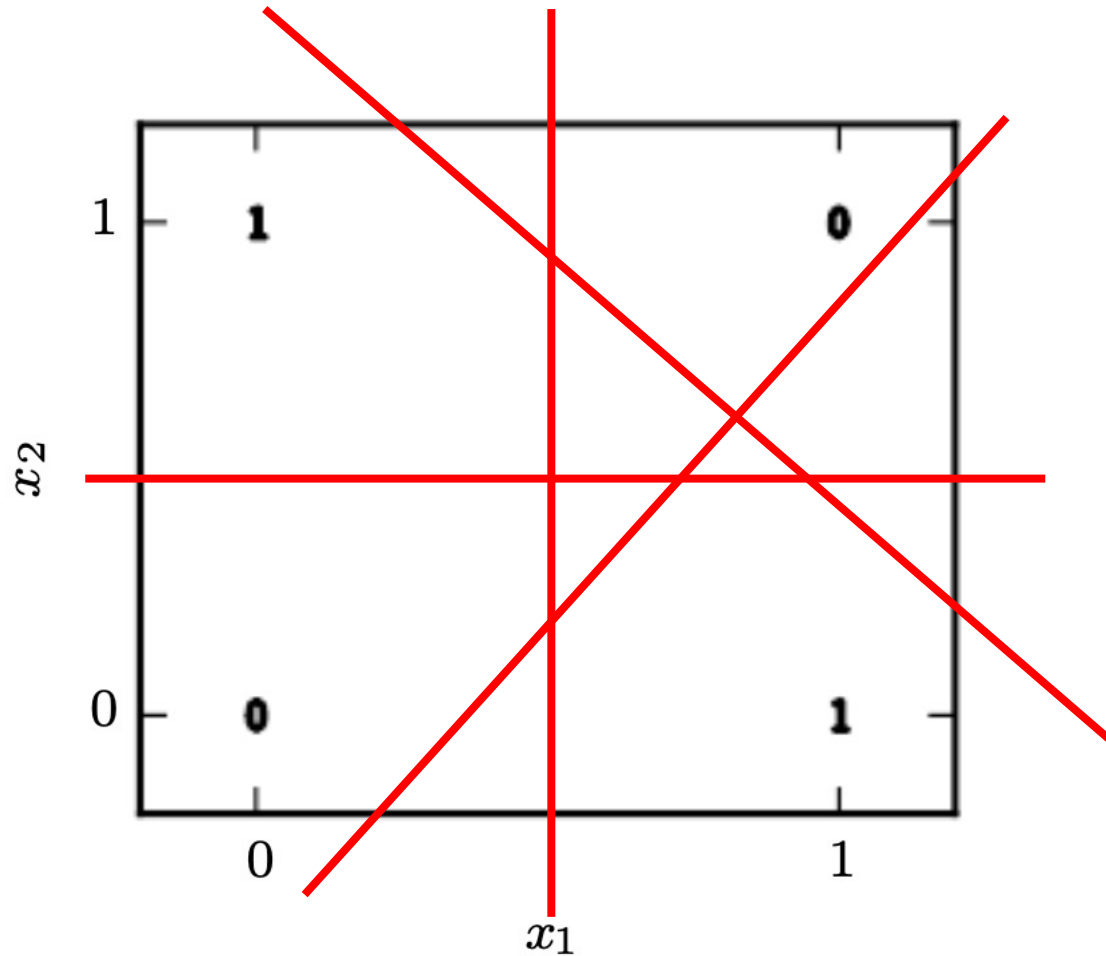


Learning XOR function

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

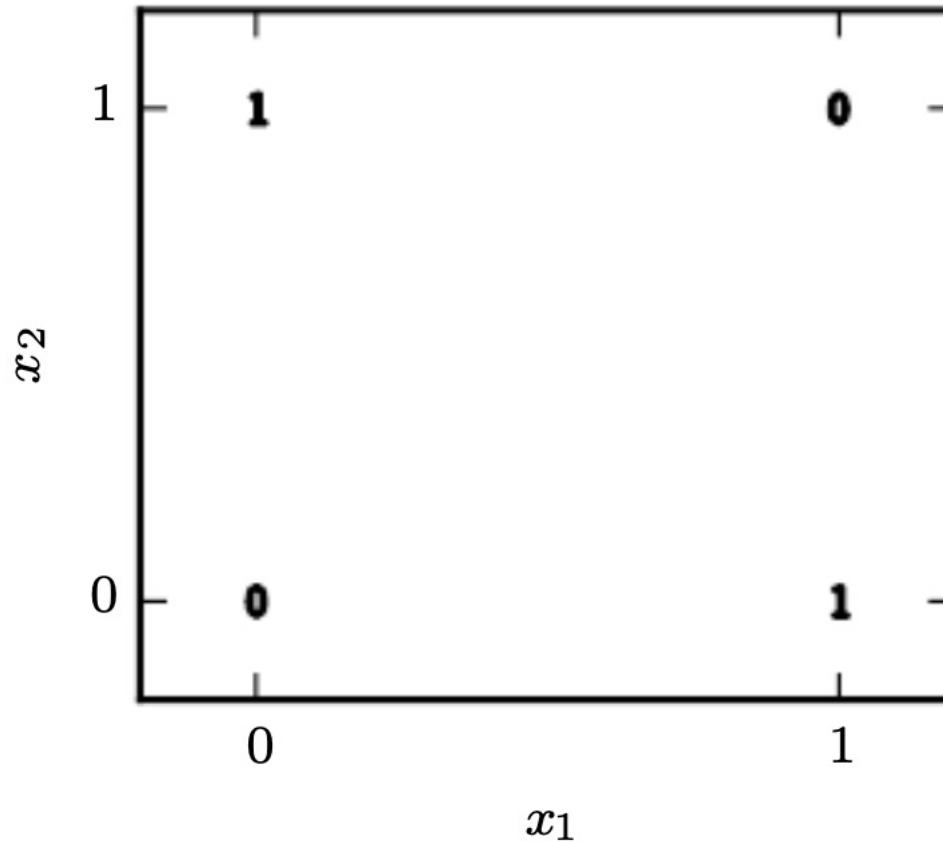
Q: Can we use a linear model to separate 0/1 classes?

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



Learning XOR function

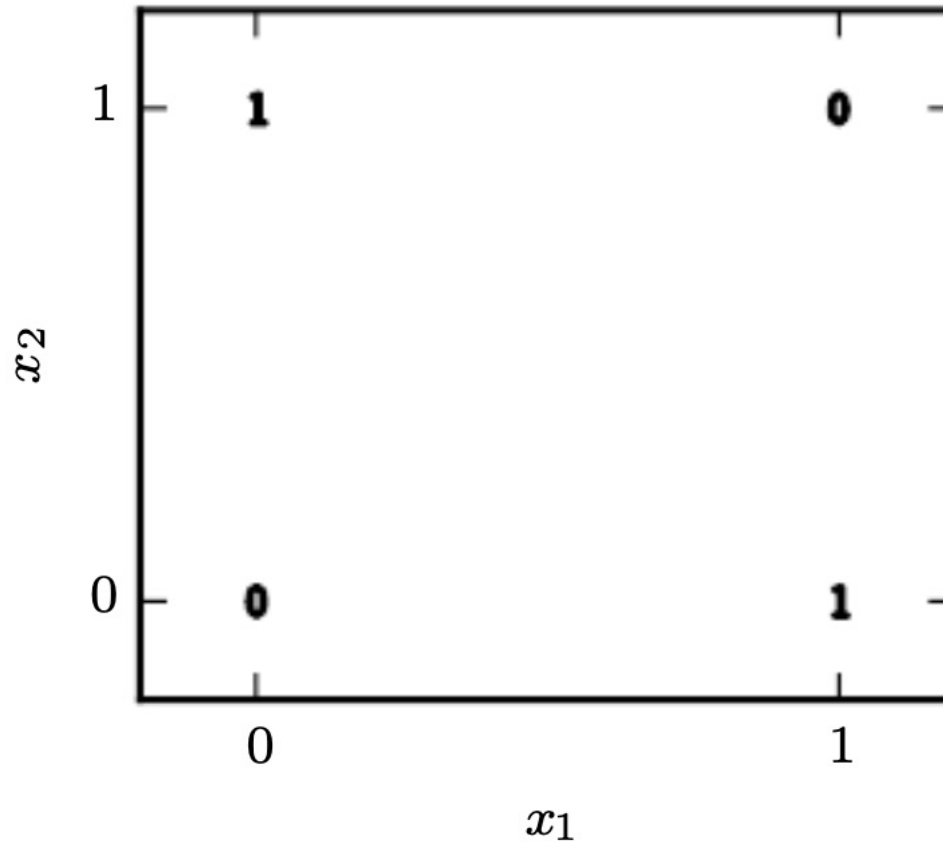
$$\mathbf{X} = \begin{array}{cc|c} & x_1 & x_2 & \\ \hline & 0 & 0 & 0 \\ & 0 & 1 & 1 \\ & 1 & 0 & 1 \\ & 1 & 1 & 0 \end{array}$$



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

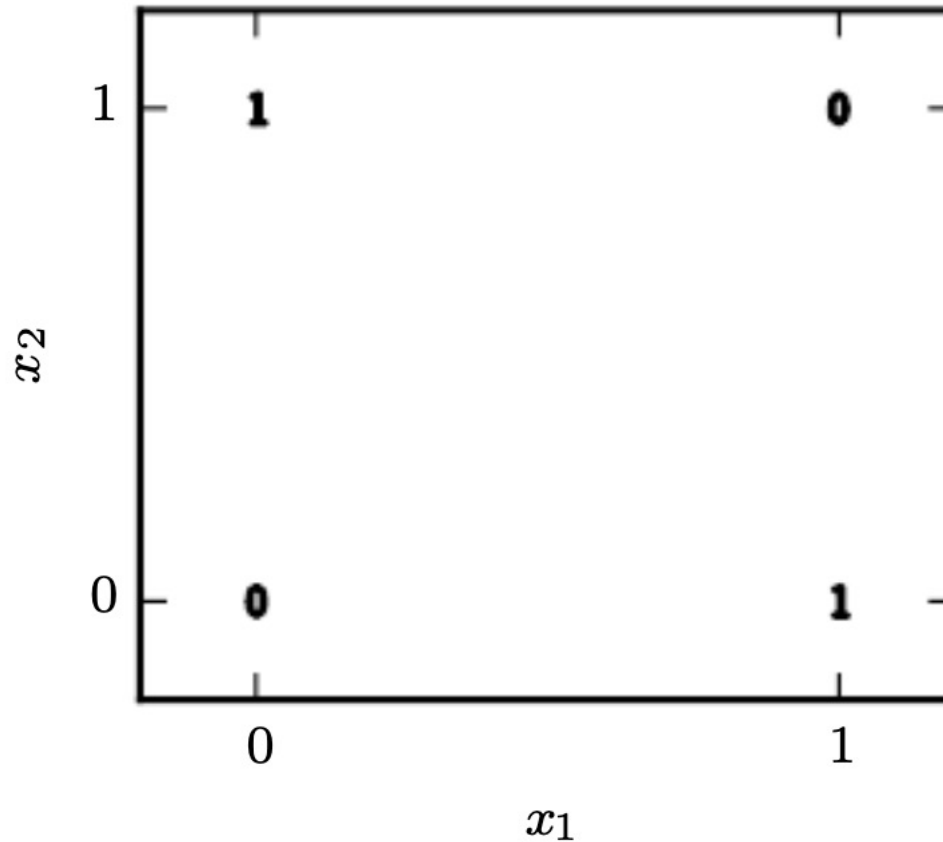
$$\mathbf{X} = \begin{array}{cc|c} & x_1 & x_2 & \\ \hline & 0 & 0 & 0 \\ & 0 & 1 & 1 \\ & 1 & 0 & 1 \\ & 1 & 1 & 0 \end{array}$$



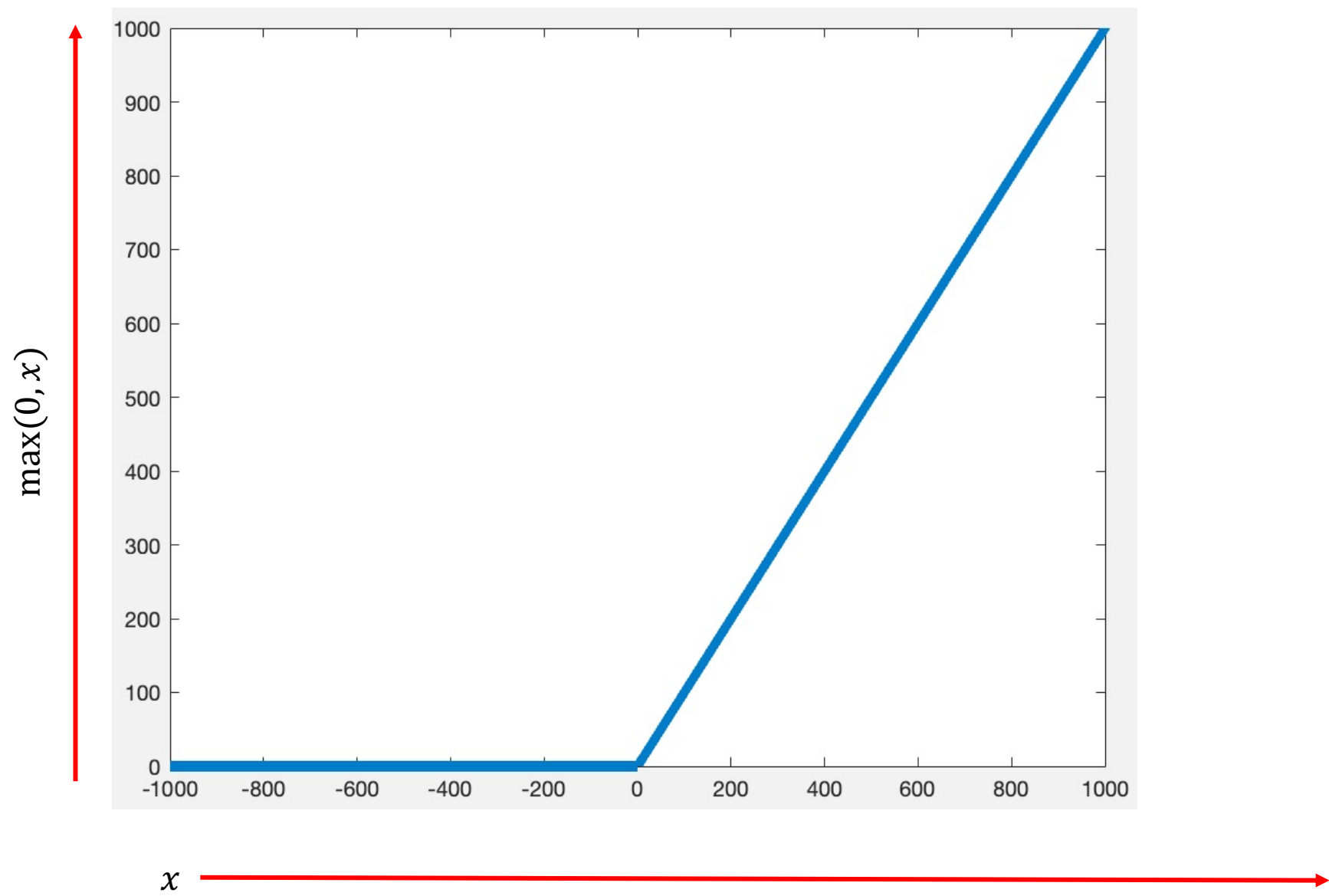
What if we use a nonlinear function as: $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + b$.

Learning XOR function

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

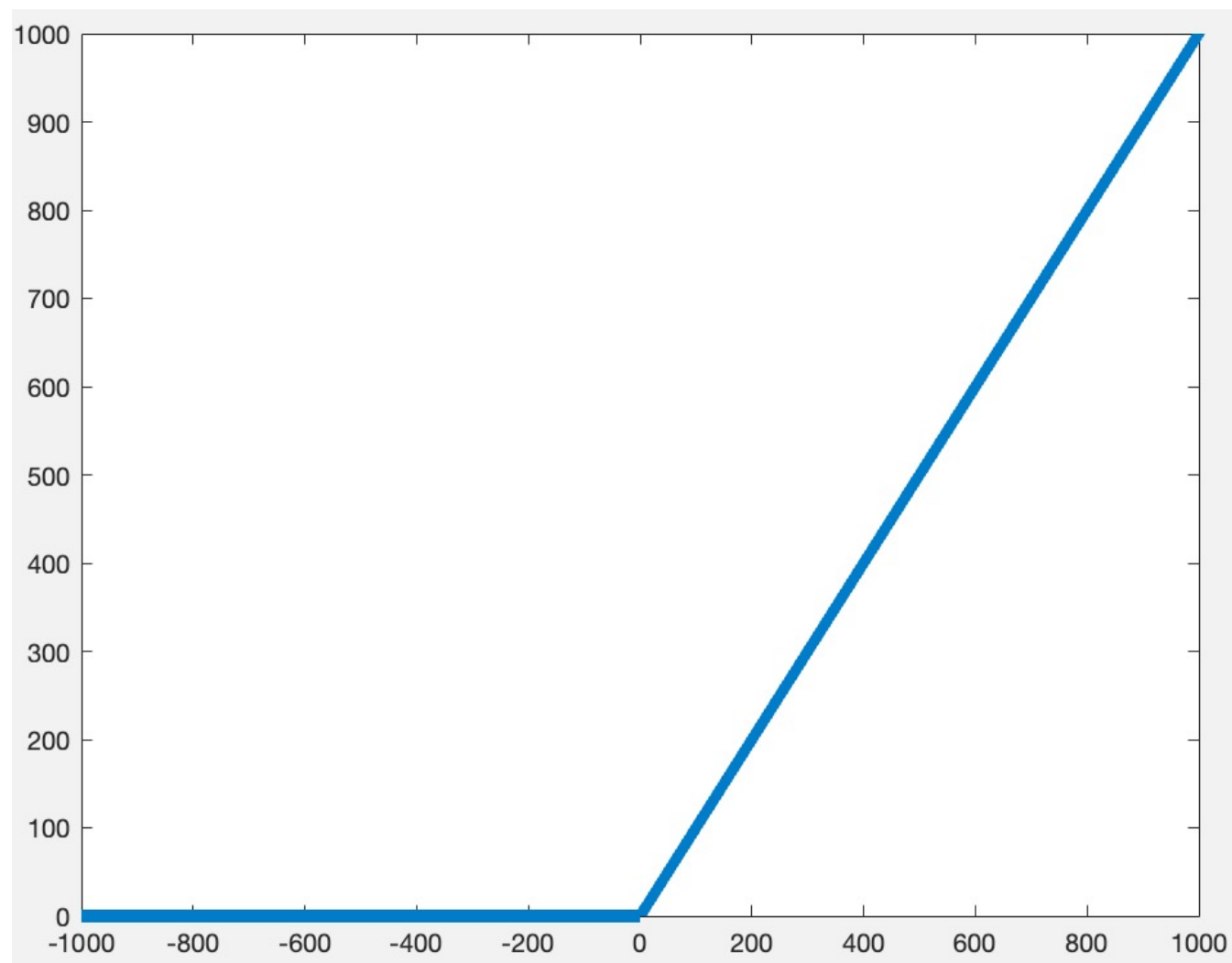


What if we use a nonlinear function as: $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + b$.



Piecewise linear

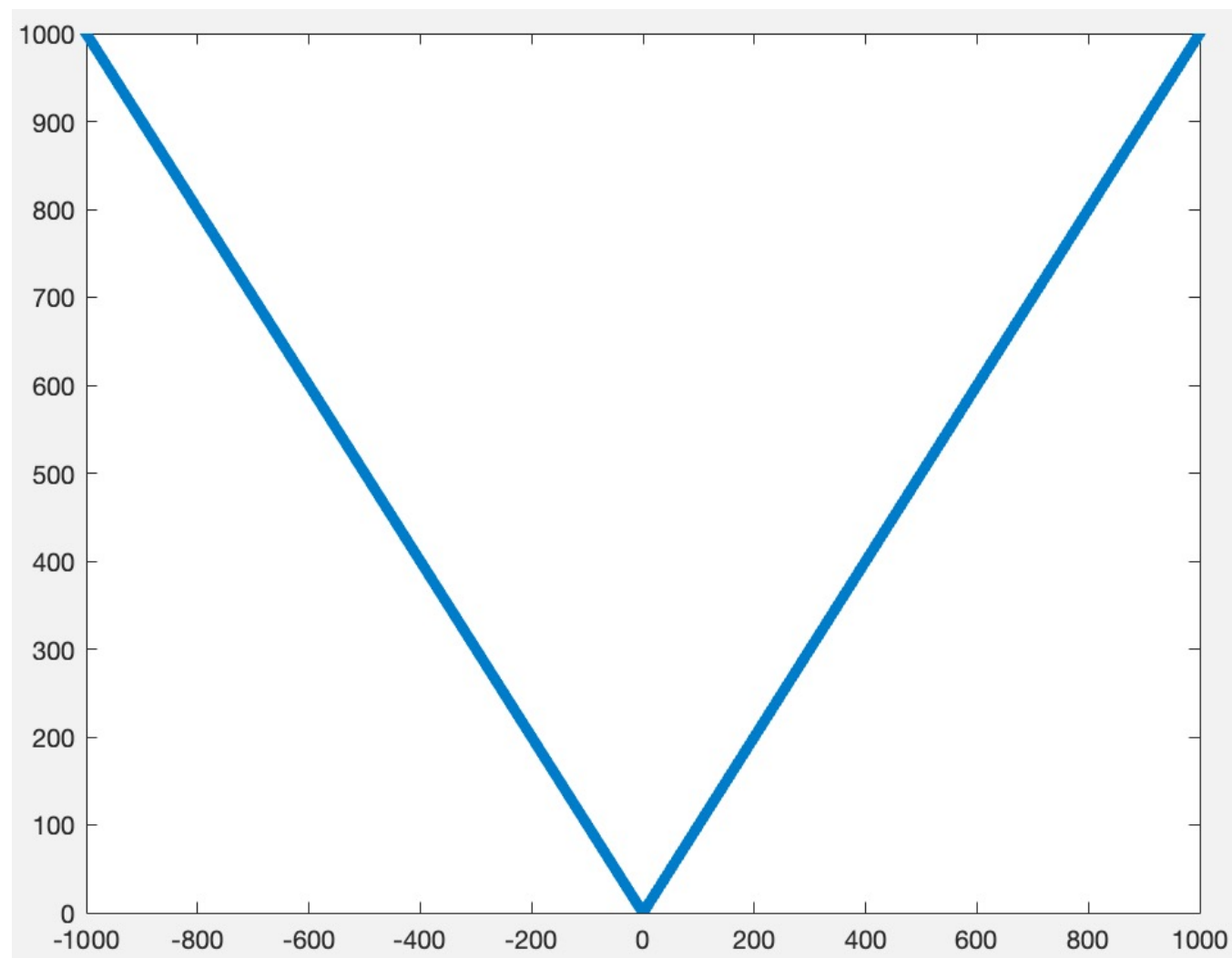
$\max(0, x)$



x

Piecewise linear

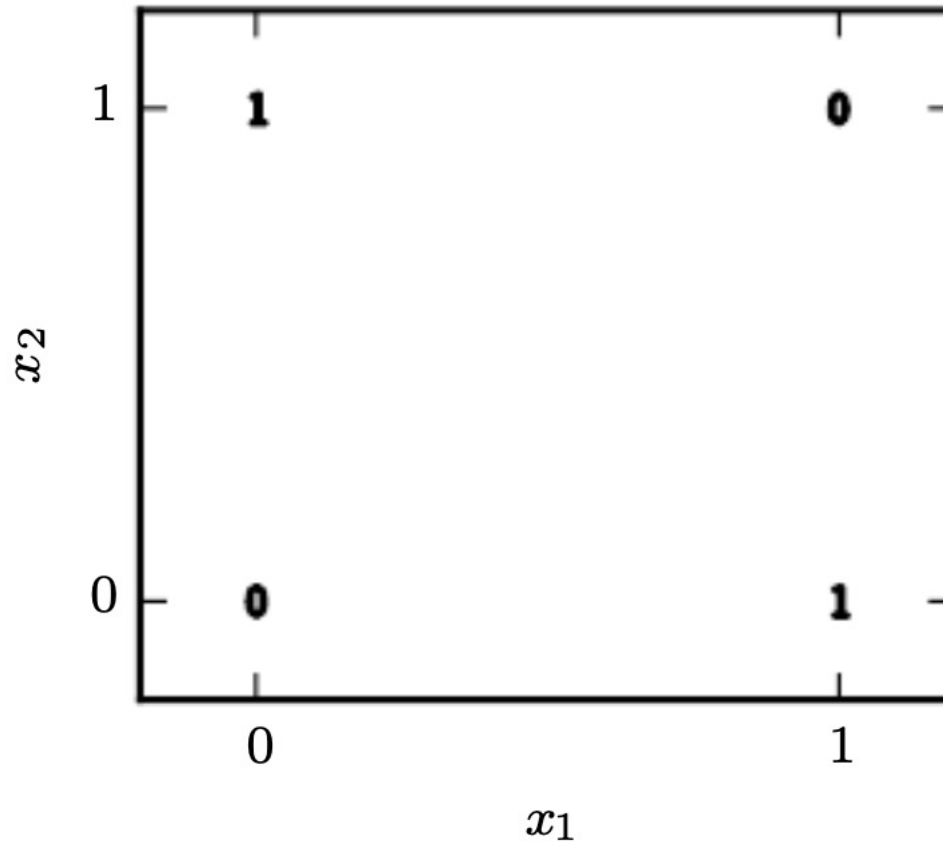
$\text{abs}(x)$



x

Learning XOR function

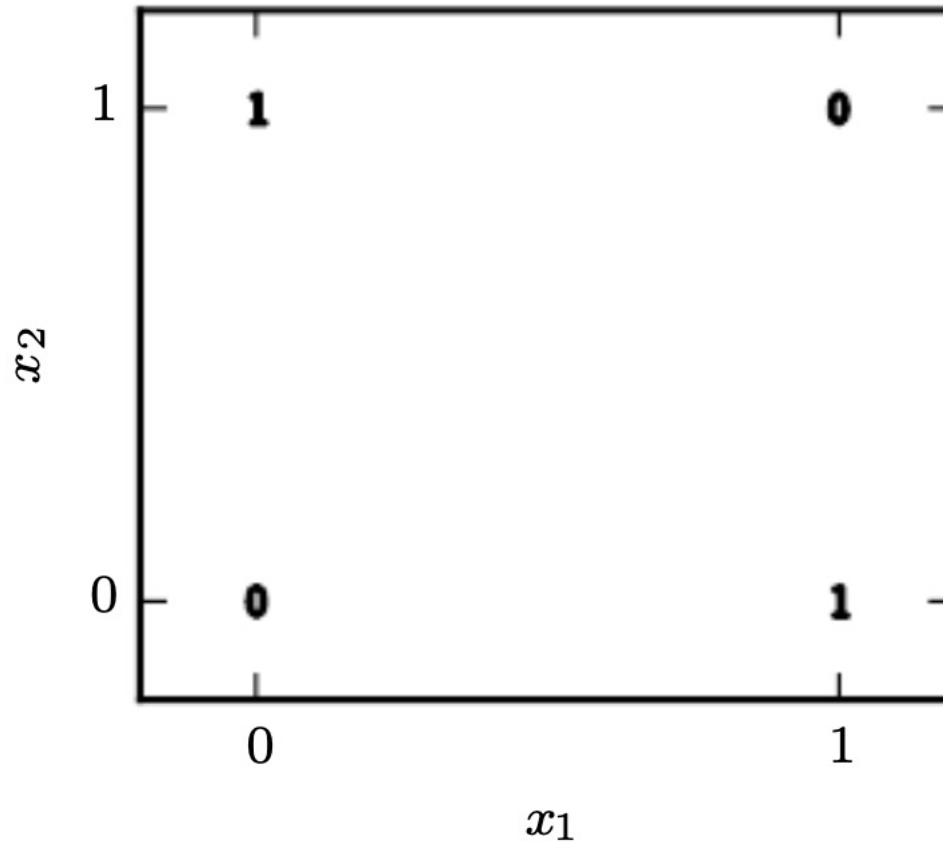
$$\mathbf{X} = \begin{array}{cc|c} & x_1 & x_2 & \\ \hline & 0 & 0 & 0 \\ & 0 & 1 & 1 \\ & 1 & 0 & 1 \\ & 1 & 1 & 0 \end{array}$$



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

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

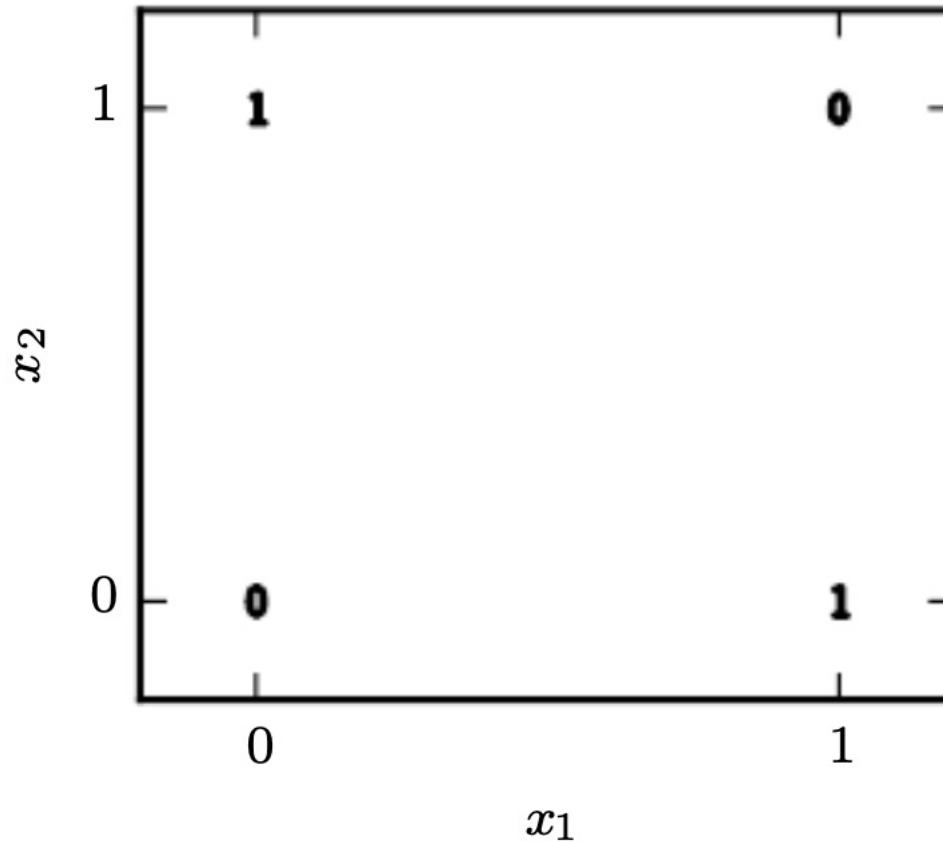


This output \rightarrow input of a linear function

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

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



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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

$$\mathbf{X} = \begin{matrix} & \begin{matrix} x_1 & x_2 \end{matrix} \\ \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} & \begin{matrix} 0 \\ 1 \\ 1 \\ 0 \end{matrix} \end{matrix} \xrightarrow{\text{red arrow}} \mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} \xrightarrow{\text{red arrow}} + \mathbf{c}$$
$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix},$$
$$\mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix},$$
$$\mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix},$$

What if we use a nonlinear function as: $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + \text{no}$

Learning XOR function

$$\mathbf{X} = \begin{matrix} & \begin{matrix} x_1 & x_2 \end{matrix} \\ \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} & \begin{matrix} 0 \\ 1 \\ 1 \\ 0 \end{matrix} \end{matrix} \xrightarrow{\text{red arrow}} \mathbf{XW} = \begin{bmatrix} \boxed{0} & \boxed{0} \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} \xrightarrow{\text{red arrow}} + \mathbf{c}$$

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

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

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

What if we use a nonlinear function as: $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \boxed{\mathbf{W}^\top \mathbf{x} + \mathbf{c}}\} + \text{no}$

Learning XOR function

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

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

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

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

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

$$\mathbf{X} = \begin{matrix} & \begin{matrix} x_1 & x_2 \end{matrix} \\ \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} & \begin{matrix} 0 \\ 1 \\ 1 \\ 0 \end{matrix} \end{matrix} \xrightarrow{\text{red arrow}} \mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} \xrightarrow{\text{red arrow}} \boxed{\begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}} \quad \mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \\
 \mathbf{c} = \boxed{\begin{bmatrix} 0 \\ -1 \end{bmatrix}}, \quad \mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix},$$

What if we use a nonlinear function?

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

Learning XOR function

$$\mathbf{X} = \begin{matrix} & \begin{matrix} x_1 & x_2 \end{matrix} \\ \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} & \begin{matrix} 0 \\ 1 \\ 1 \\ 0 \end{matrix} \end{matrix} \xrightarrow{\text{red arrow}} \mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} \xrightarrow{\text{red arrow}} \boxed{\begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}} + \mathbf{c} \quad \mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \\
 \mathbf{c} = \boxed{\begin{bmatrix} 0 \\ -1 \end{bmatrix}}, \\
 \mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix},$$

The inner linear model

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

Learning XOR function

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

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

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

$\xrightarrow{\text{red arrow, max}(0, \cdot)}$

$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$

$\xleftarrow{\text{red arrow, } \mathbf{w}^\top}$

$\begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$

What if we use a nonlinear function as:

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

Learning XOR function

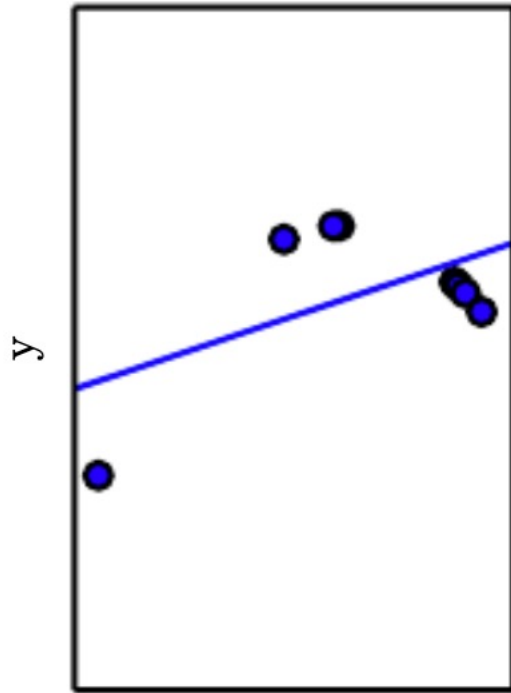
$$\begin{array}{c}
 \mathbf{X} = \begin{bmatrix} x_1 & x_2 \\ 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} \xrightarrow{\quad} \mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} \xrightarrow{+ \mathbf{c}} \begin{bmatrix}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

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

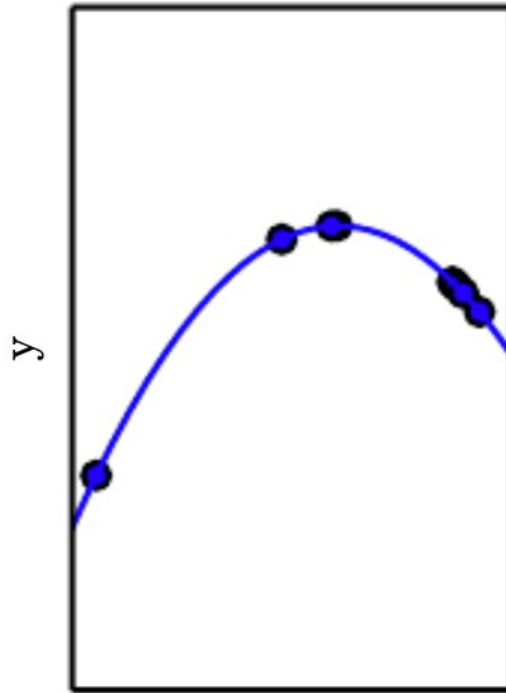
$$\mathbf{W}^\top \mathbf{x} + \mathbf{c} \quad \text{v.s.} \quad \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\}$$

Learning XOR function



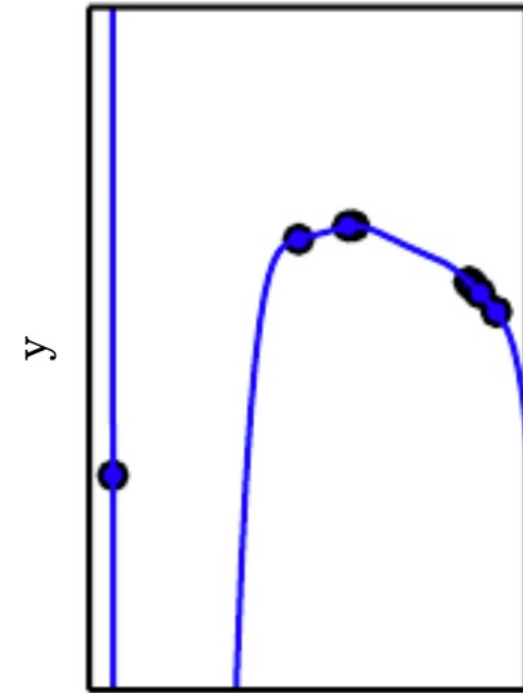
x_0

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



x_0

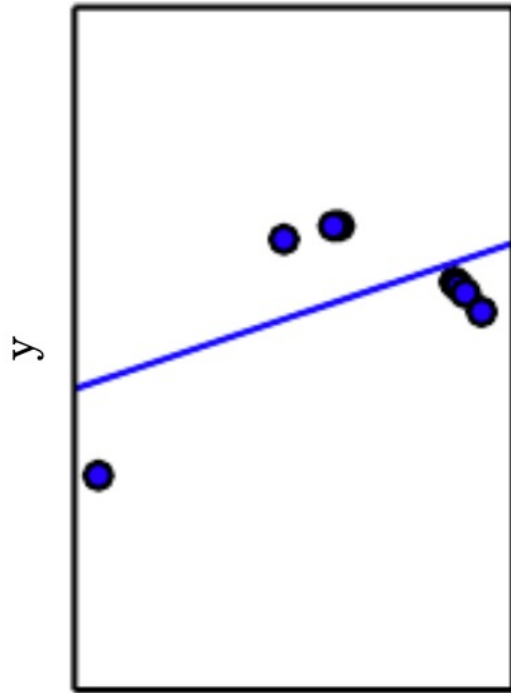
Quadratic model
 $f(w; x) = w_2 x^2 + w_1 x^1 + w_0$



x_0

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

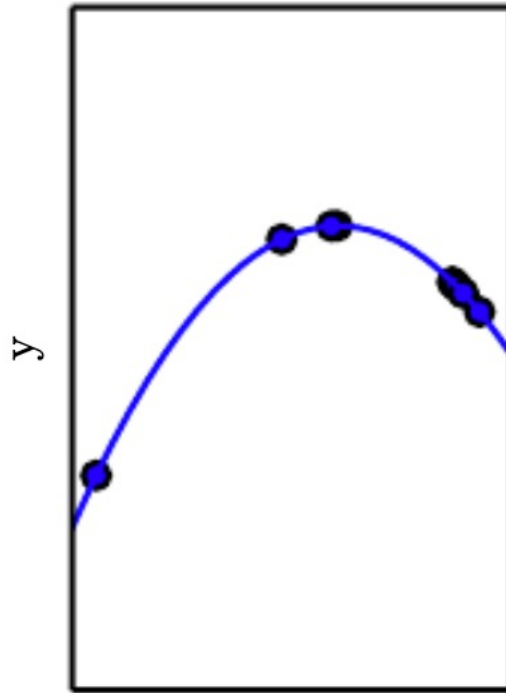
Learning XOR function



x_0

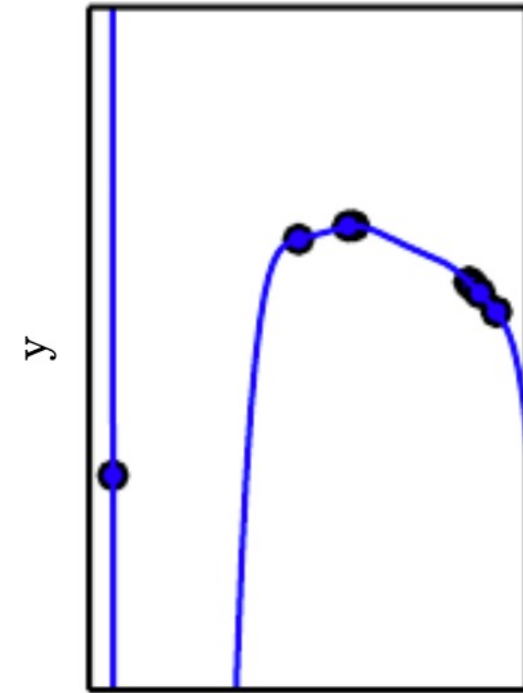
Linear model
 $f(w; x) = w_1 x^1 + w_0$

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



x_0

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

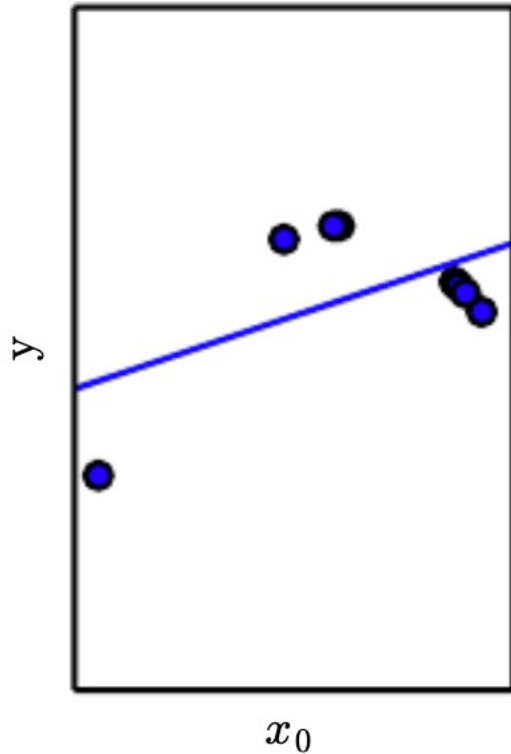


x_0

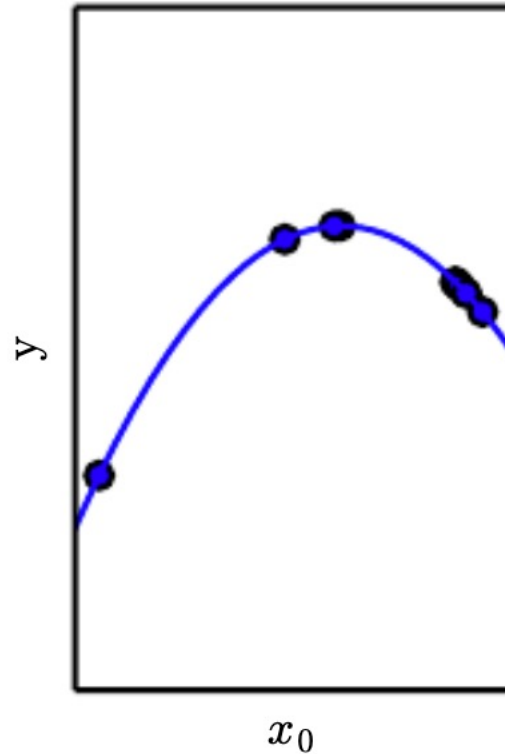
Polynomial model (9 degree)

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

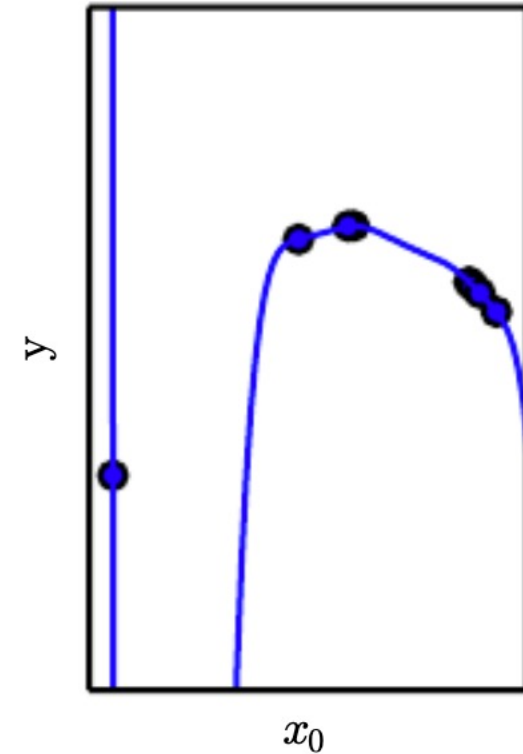
Learning XOR function



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

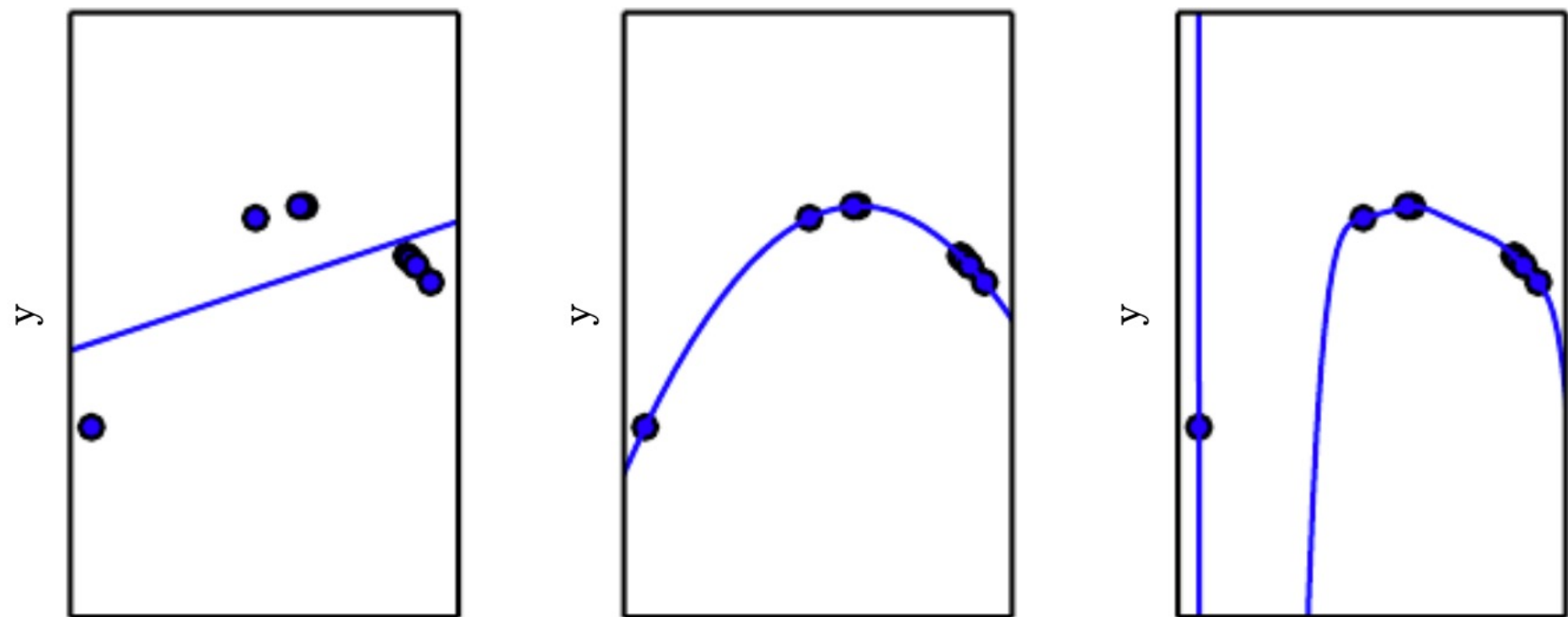


Quadratic model
 $f(w; x) = w_2 x^2 + w_1 x^1 + w_0$
$$= \sum_{i=3}^9 0 \cdot x^i + \sum_{i=1}^2 w_i x^i + w_0$$



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

$I=5$

$I=1$

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 \leftarrow 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 \leftarrow 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?

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

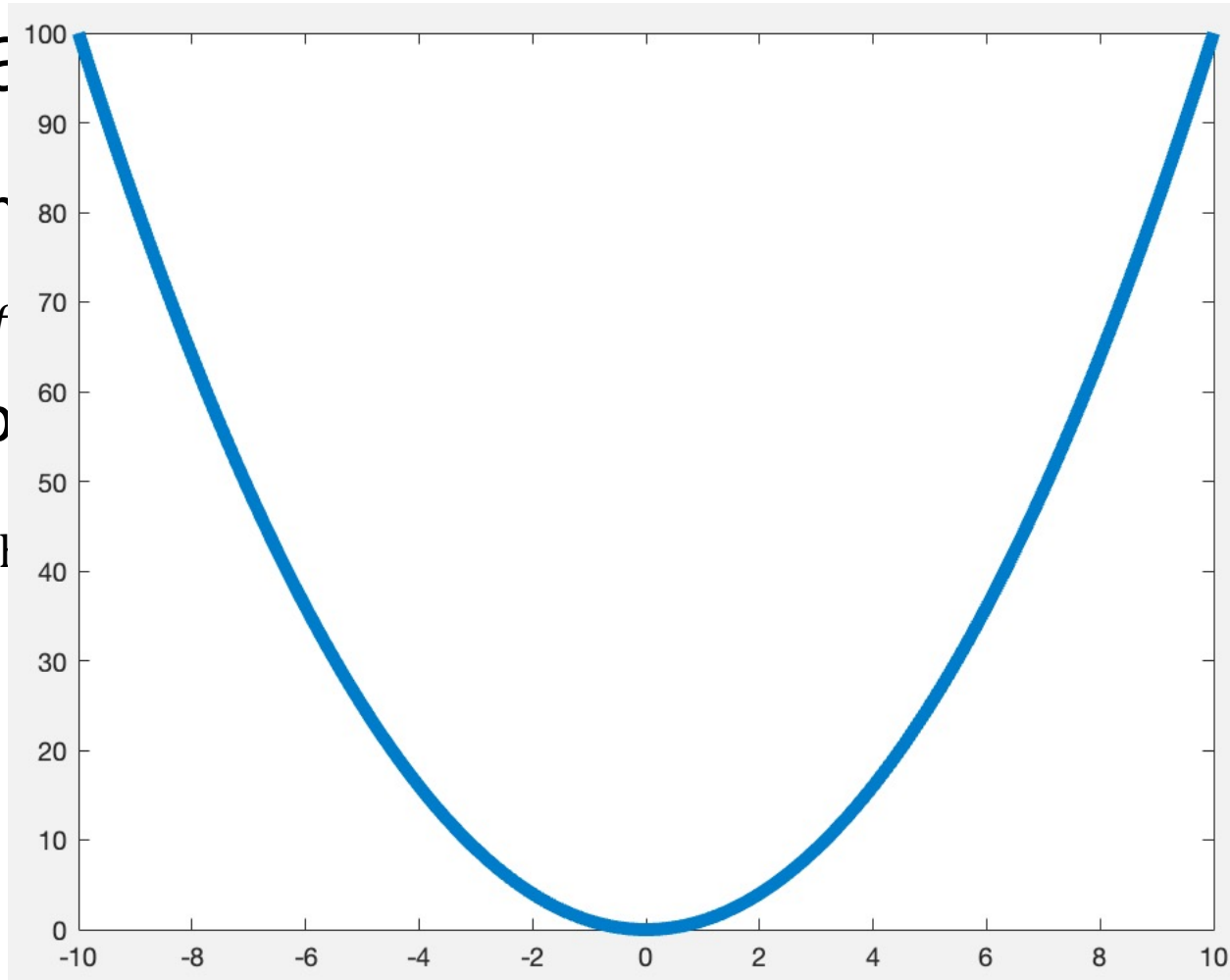
What makes feedforward network different from linear

- Nonlinear function

$$f_3(f$$

- Q: Why comp

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



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

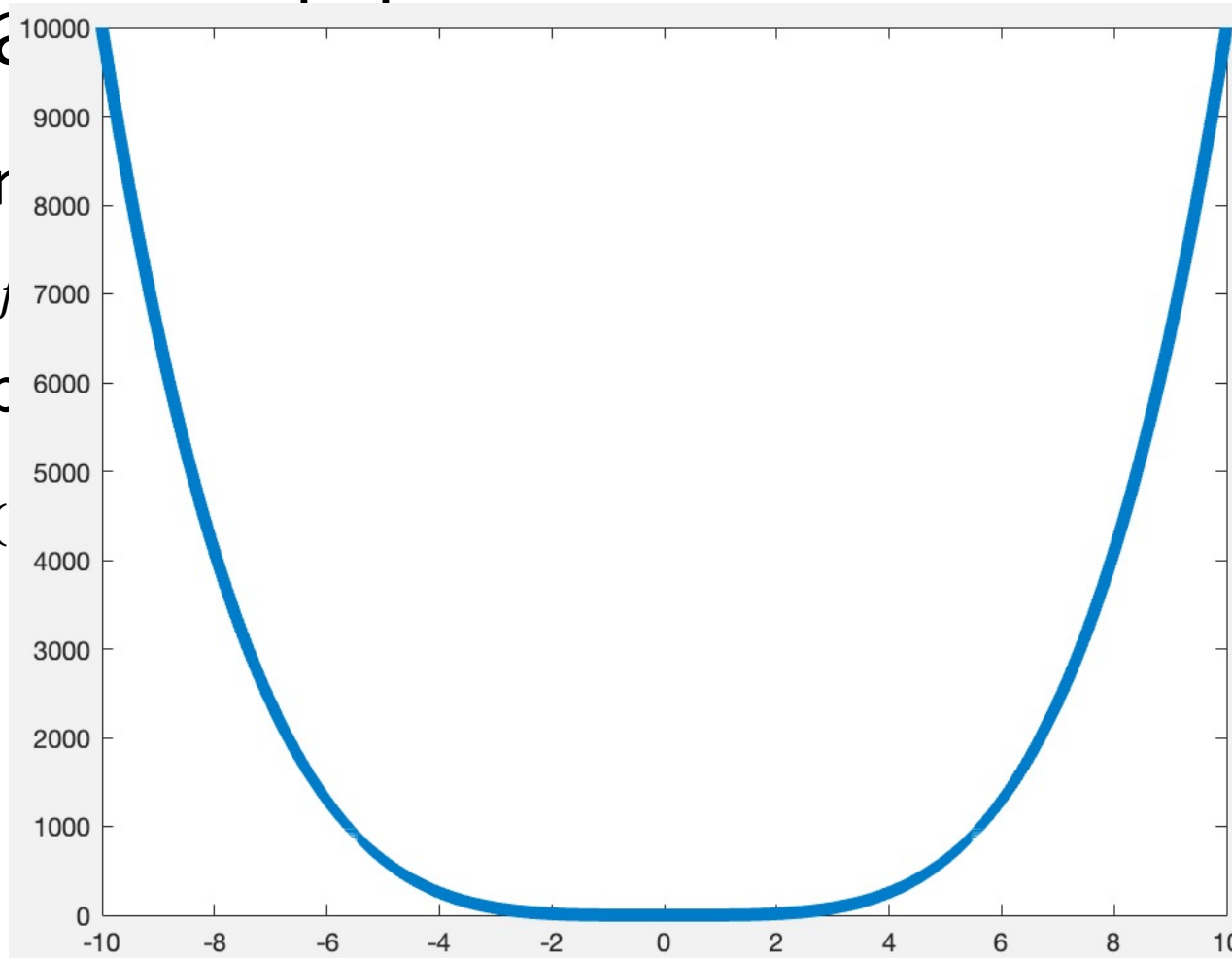
What makes feedforward network different from linear?

- Nonlinear function

$$f_3(j)$$

- Q: Why compute

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



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

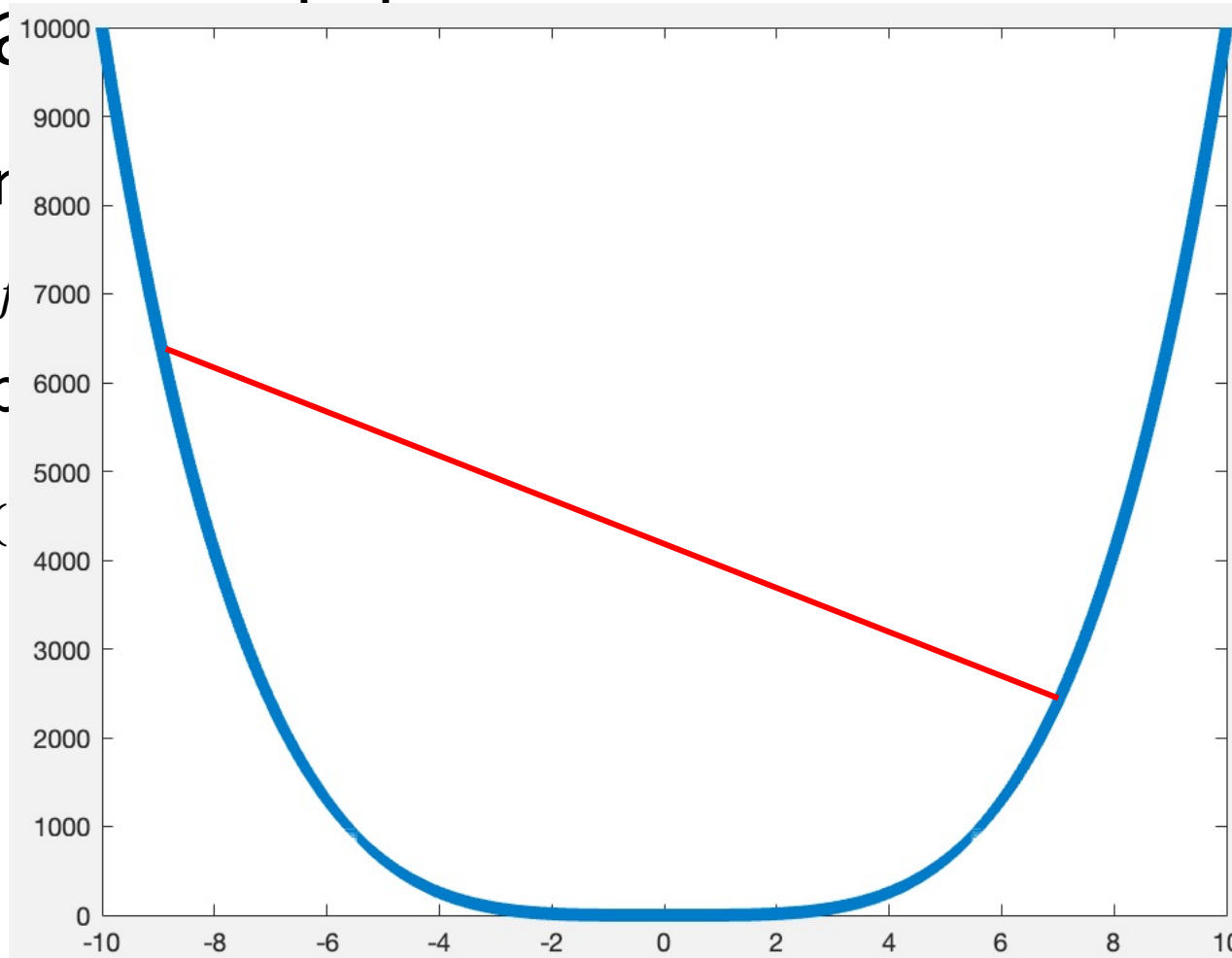
What makes feedforward network different from linear?

- Nonlinear function

$$f_3(j)$$

- Q: Why compute

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



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

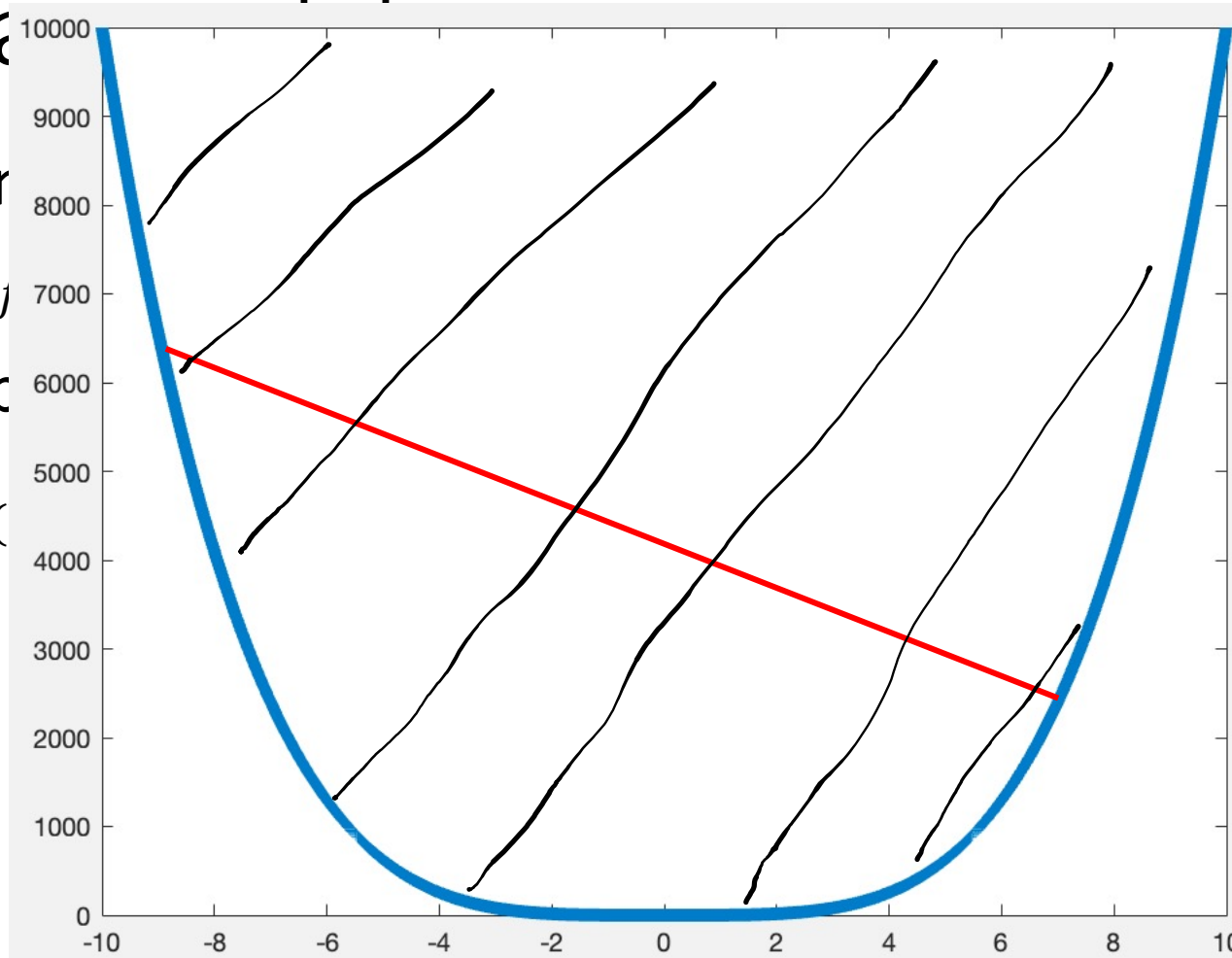
What makes feedforward network different from linear?

- Nonlinear function

$$f_3(x)$$

- 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 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

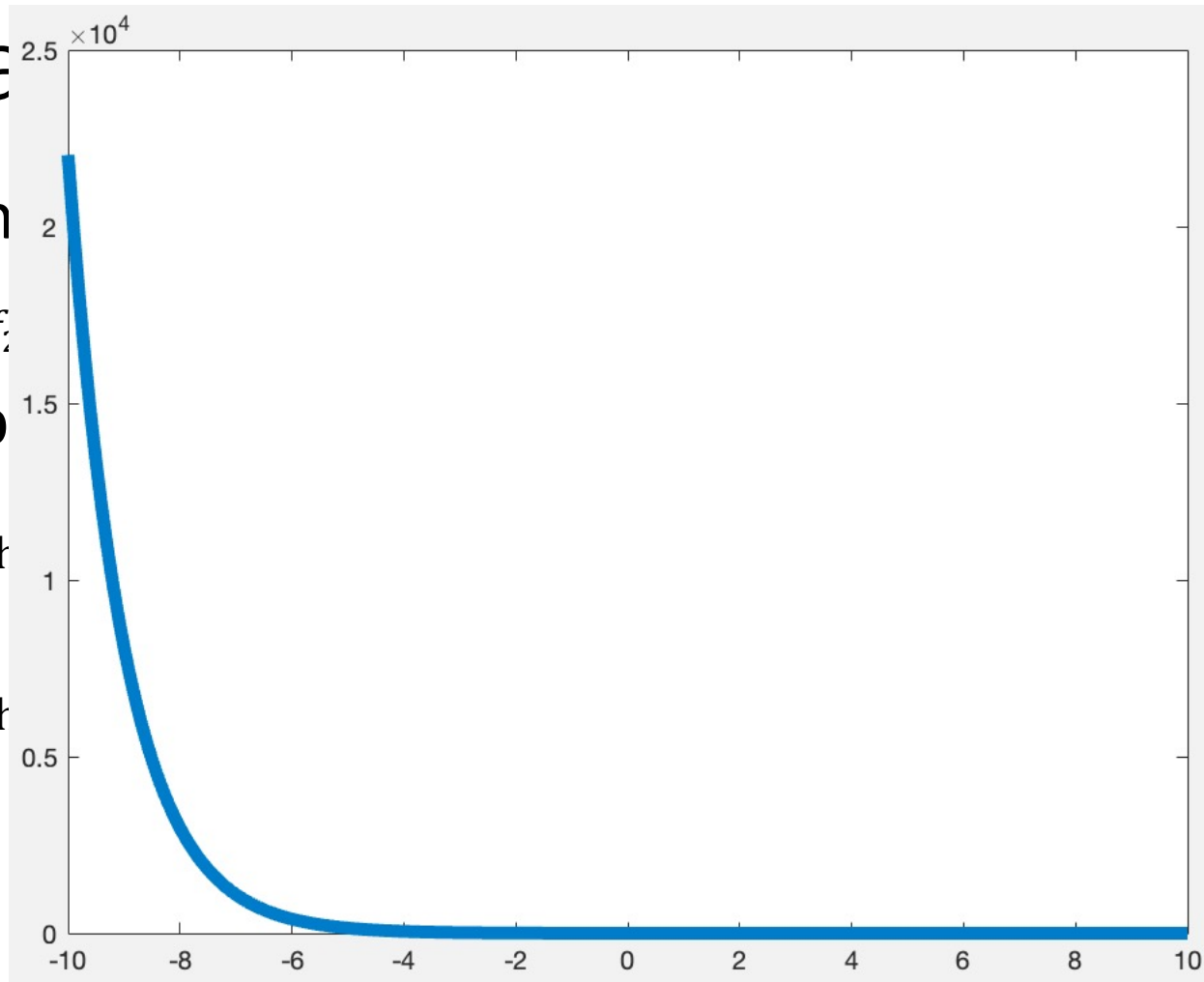
- Nonlinear function

$$f_3(f_2$$

- Q: Why comp

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

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



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

What makes feedforward network different from linear

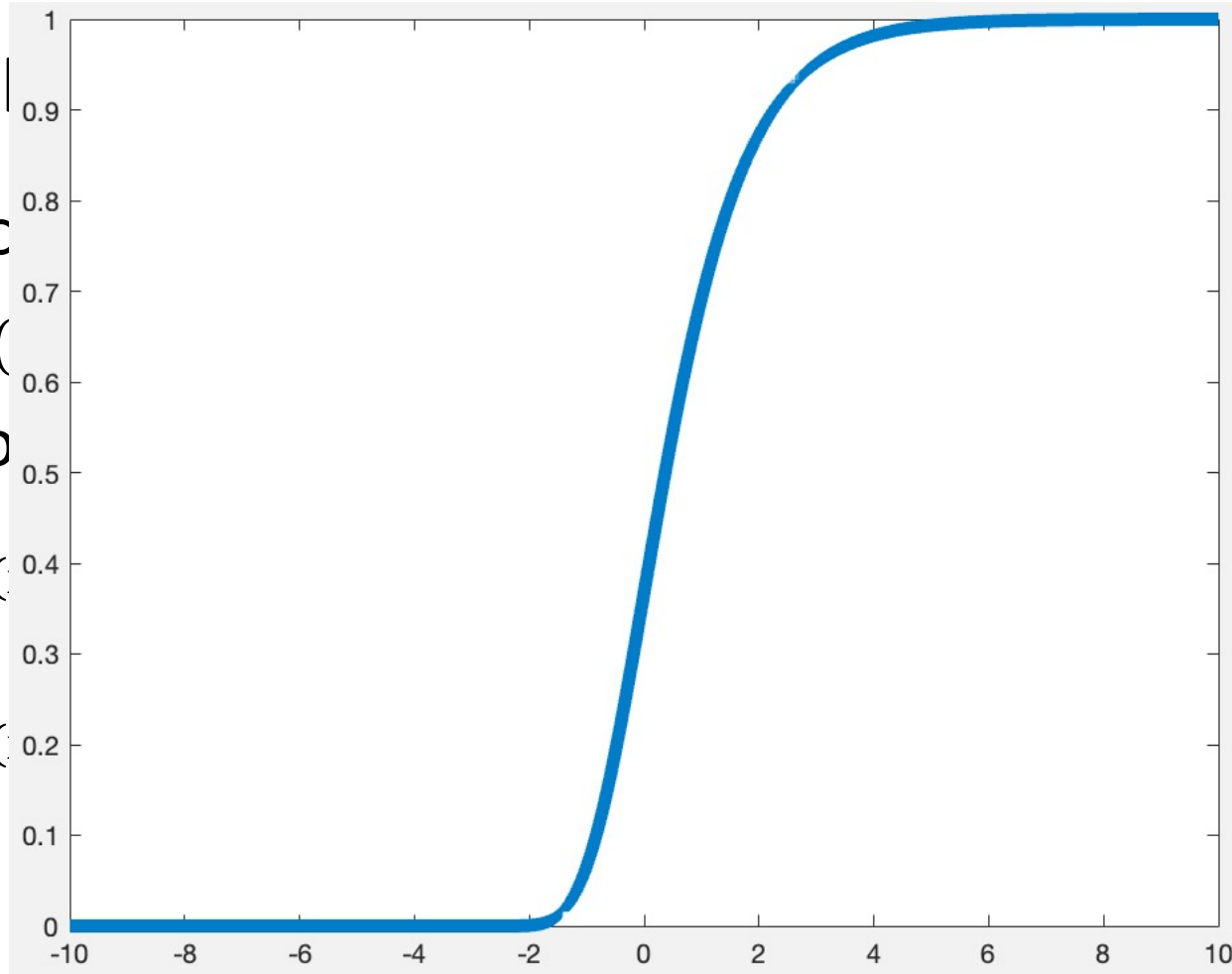
- Nonlinear func

$$f_3(f_2($$

- Q: Why compo

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

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



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

What makes feedforward network different from linear model

- Nonlinear functions in hidden layers

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

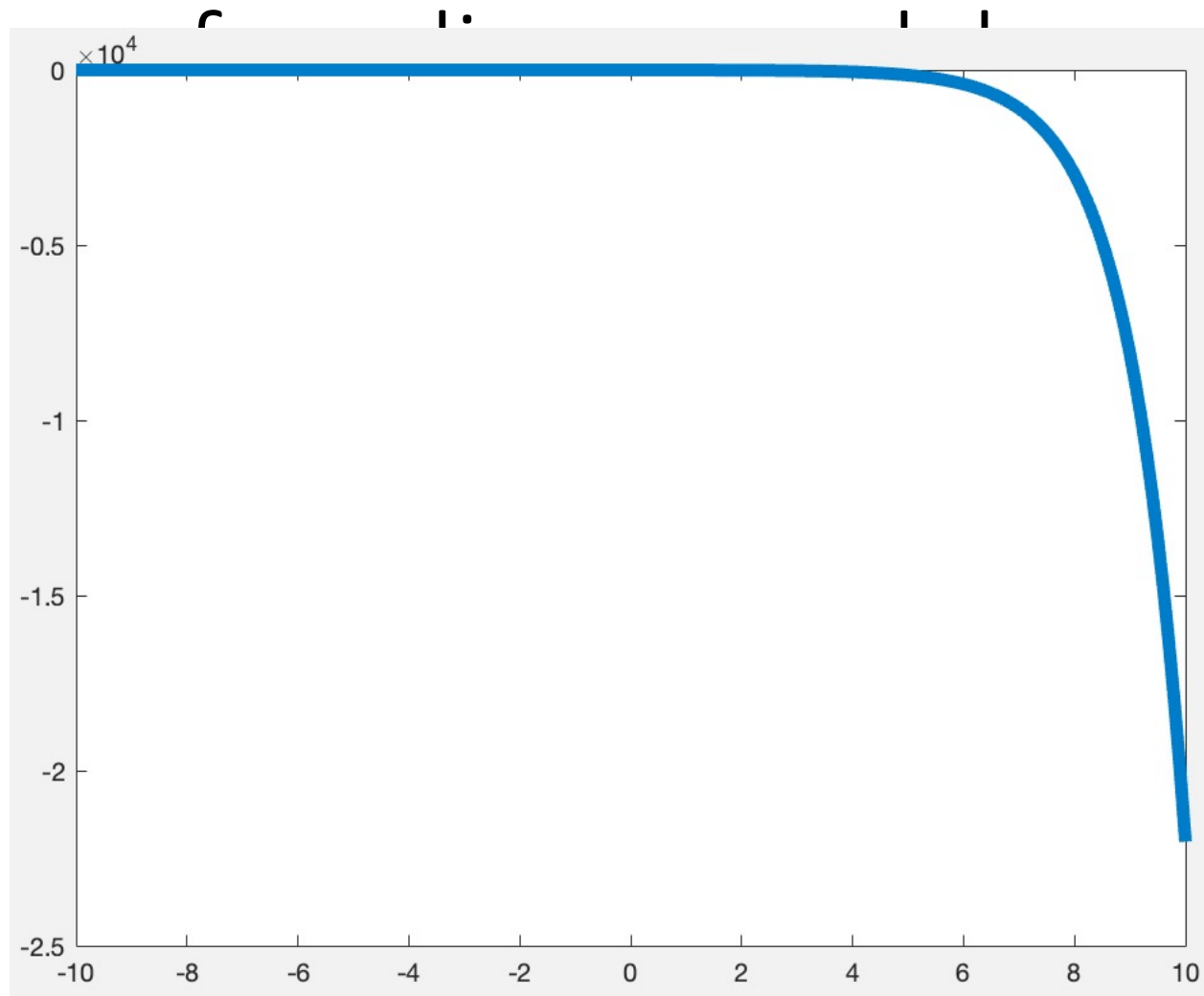
- Q: Why composition makes nonconvexity?

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

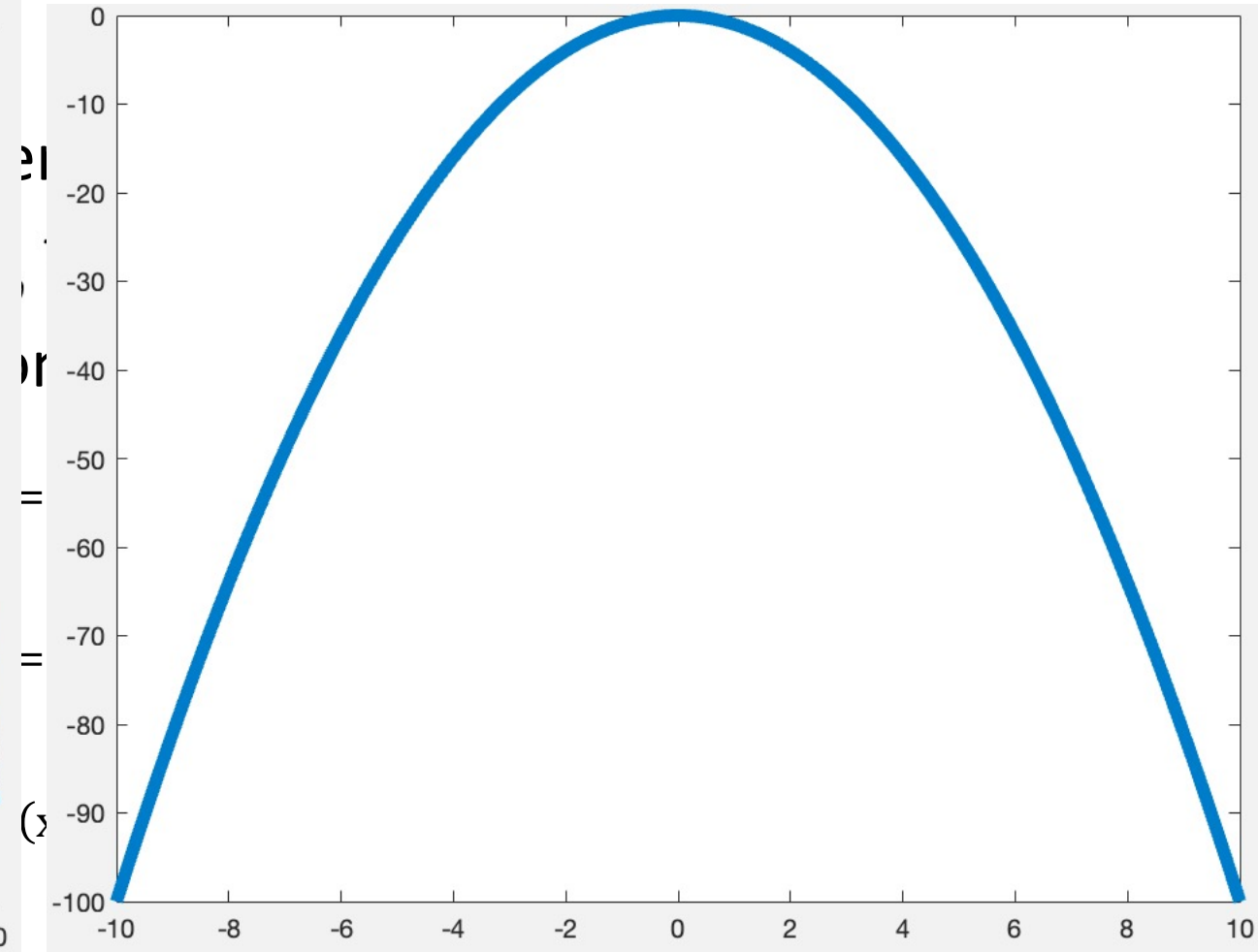
$$f(x) = g(h(x)) \quad \text{where } g(x) = h(x) = \exp(-x)$$

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

What makes feedforward network different



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



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

What makes feedforward network different from linear model

- Nonlinear function

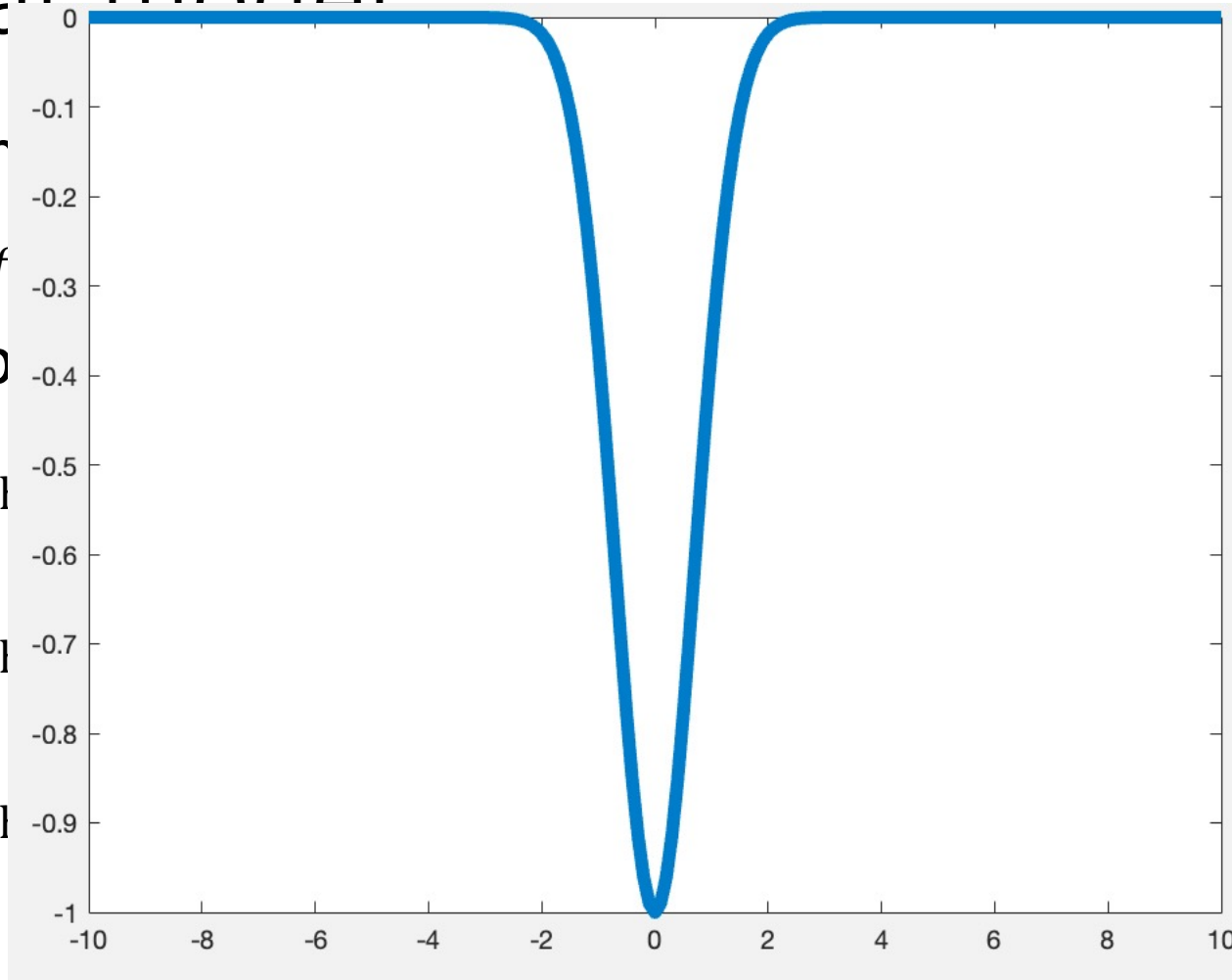
$$f_3(f_2(f_1(x)))$$

- 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

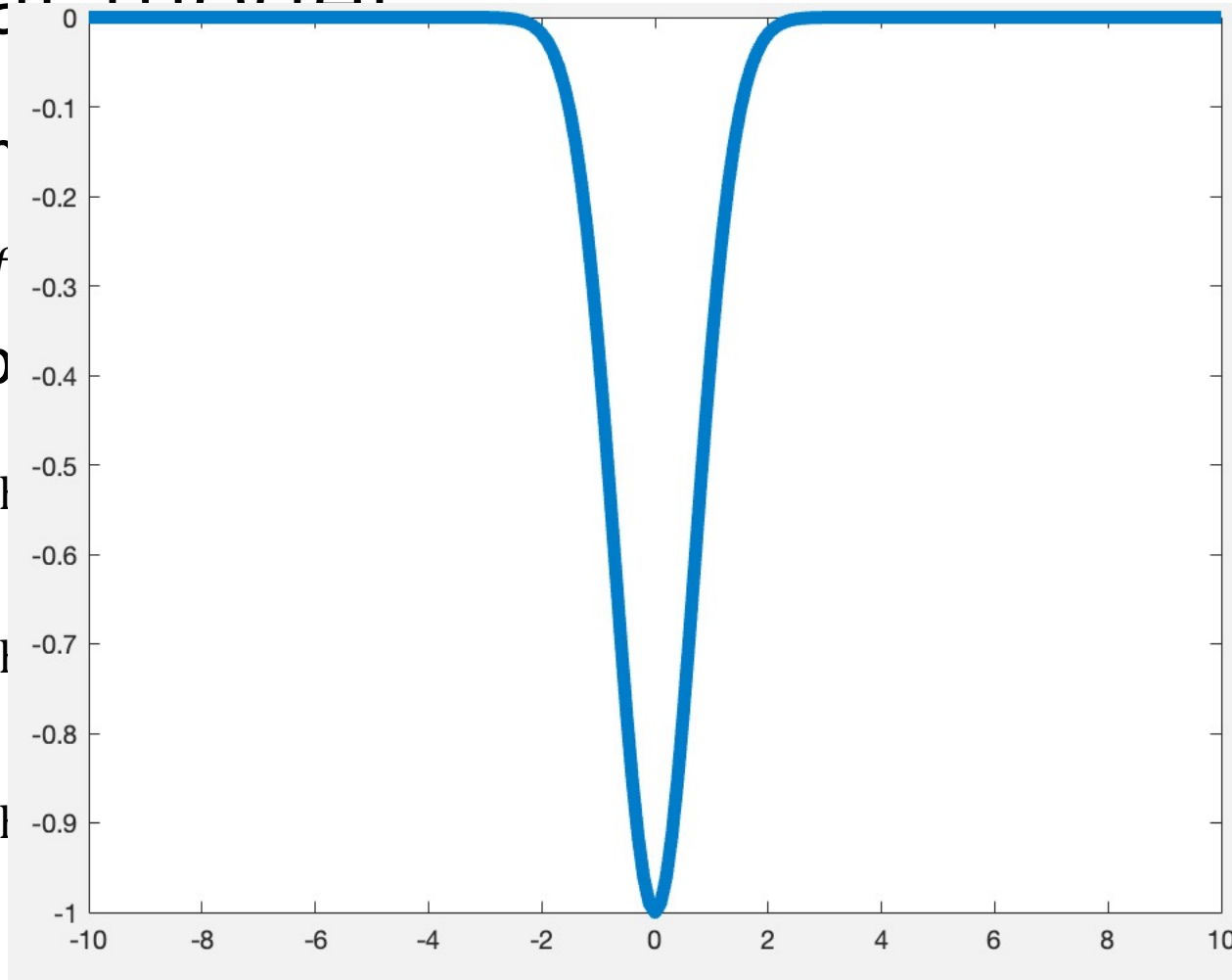
$$f_3(f_2(f_1(x)))$$

- 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

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