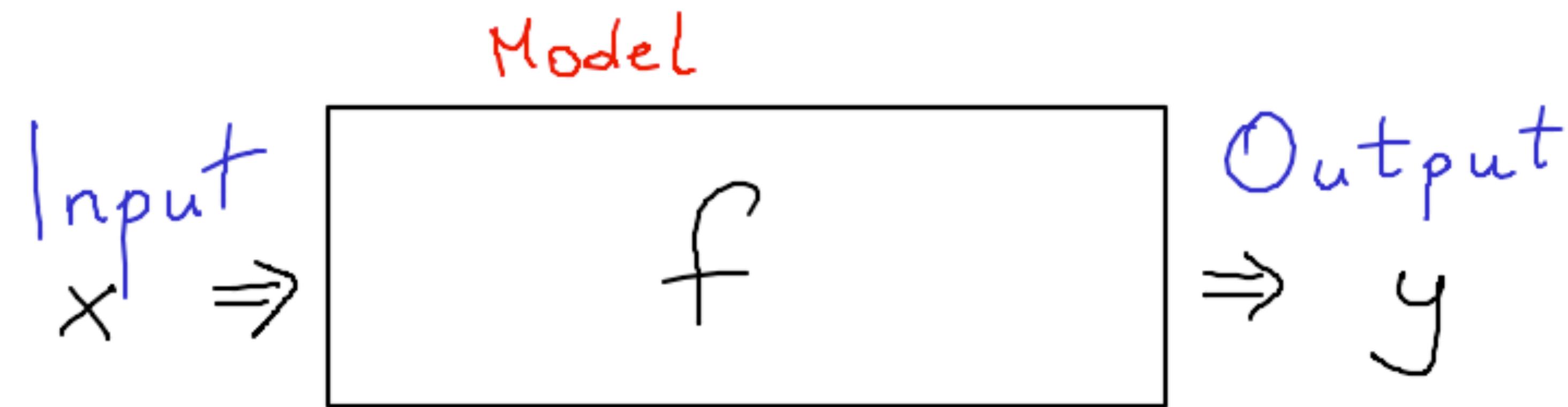


# **Neural Networks**

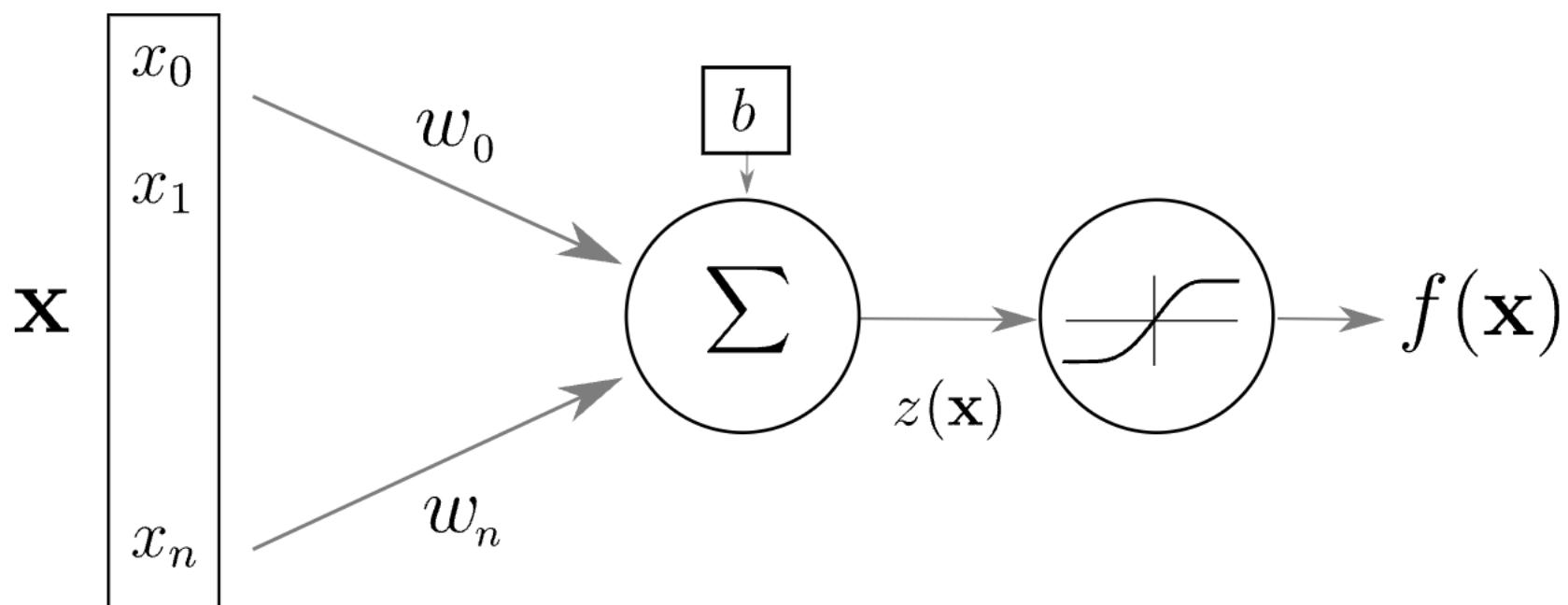
## **Calculus**

# Computation graph (recap)

- neural network = parametrized, non-linear function



# Artificial neuron



$$z(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

$$f(\mathbf{x}) = g(\mathbf{w}^T \mathbf{x} + b)$$

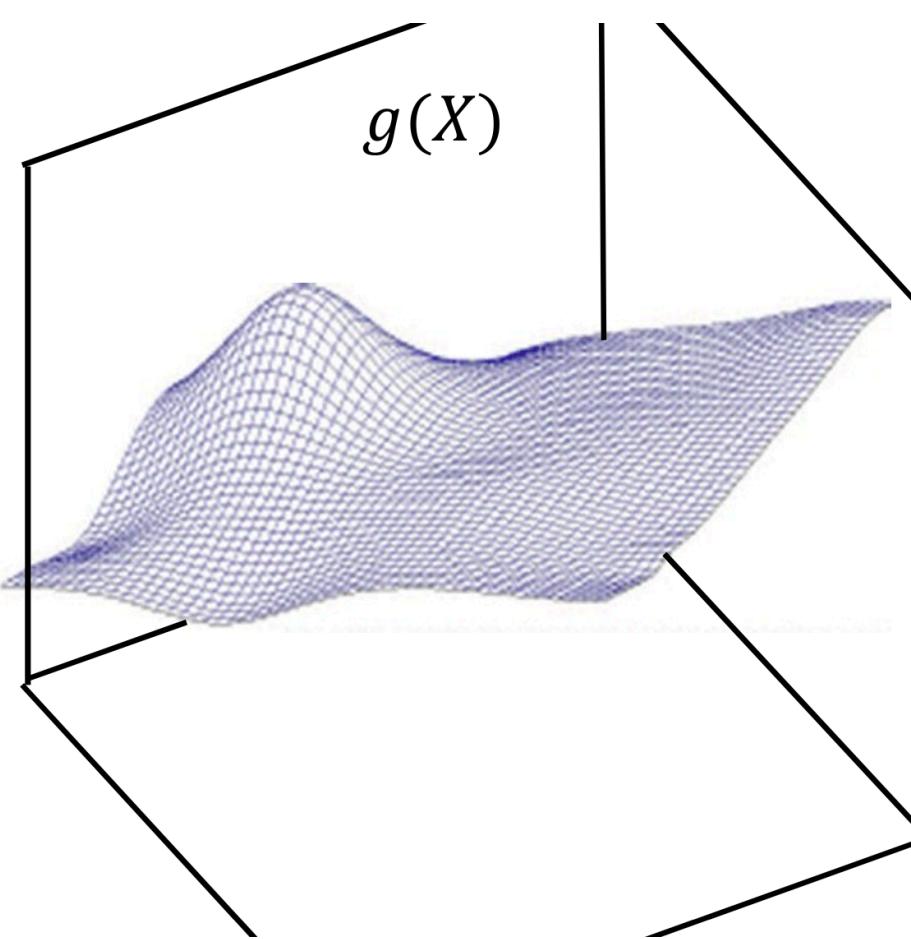
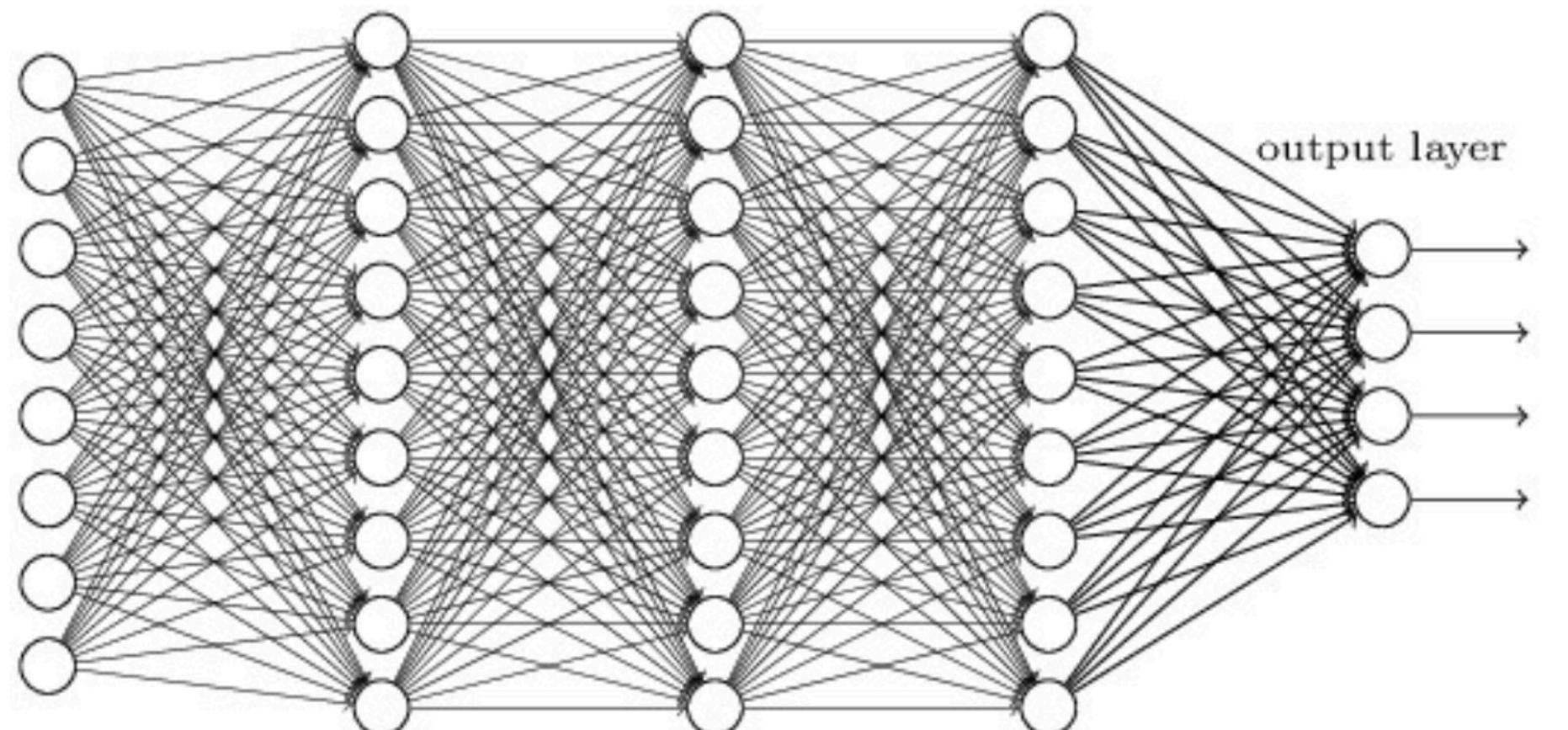
- $\mathbf{x}, f(\mathbf{x})$  input and output
- $z(\mathbf{x})$  pre-activation
- $\mathbf{w}, b$  weights and bias
- $g$  activation function

# Supervised learning

- regression
- classification
- ...

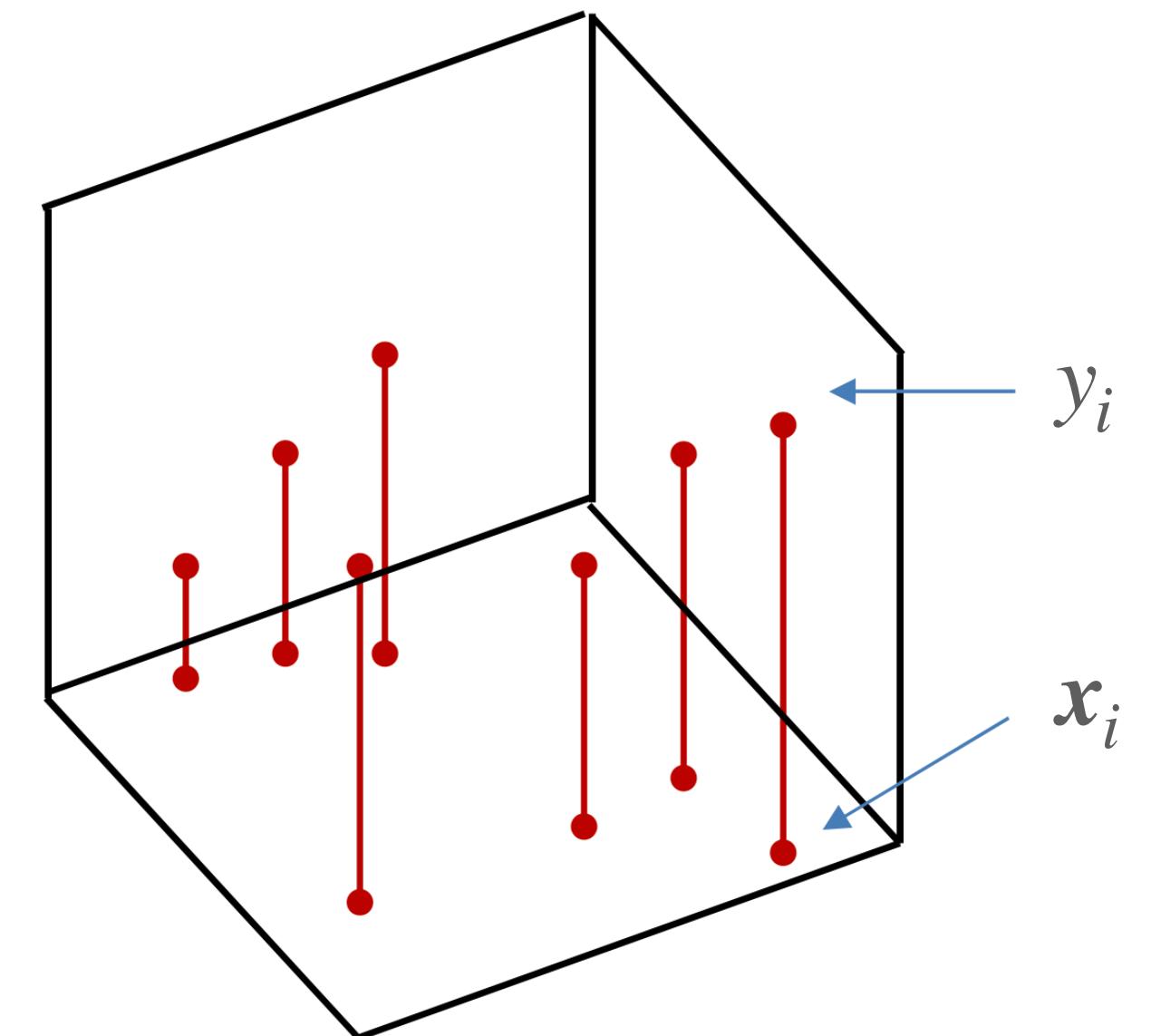
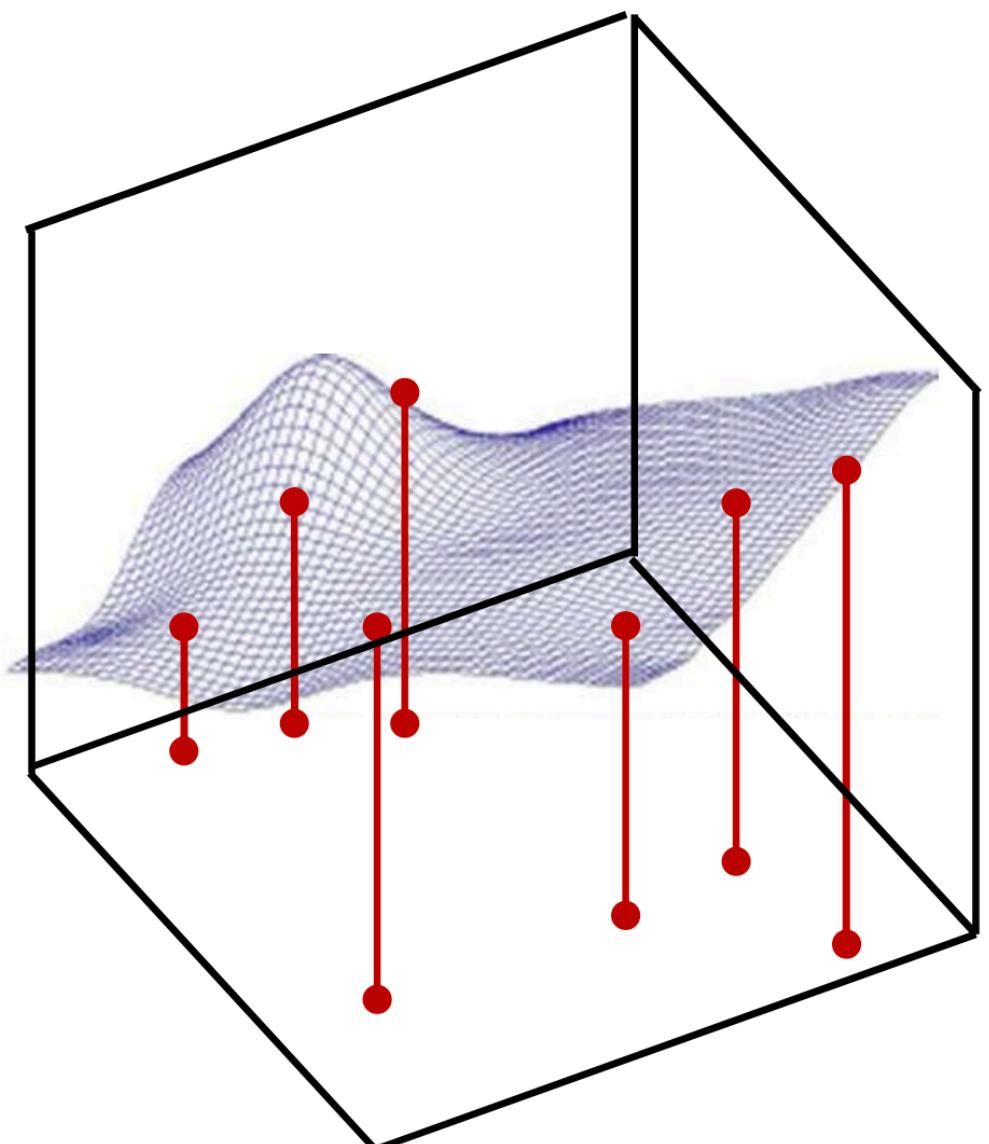
# What we learn?

- network parameters  $W$
- $f(X, W) \approx g(X)$



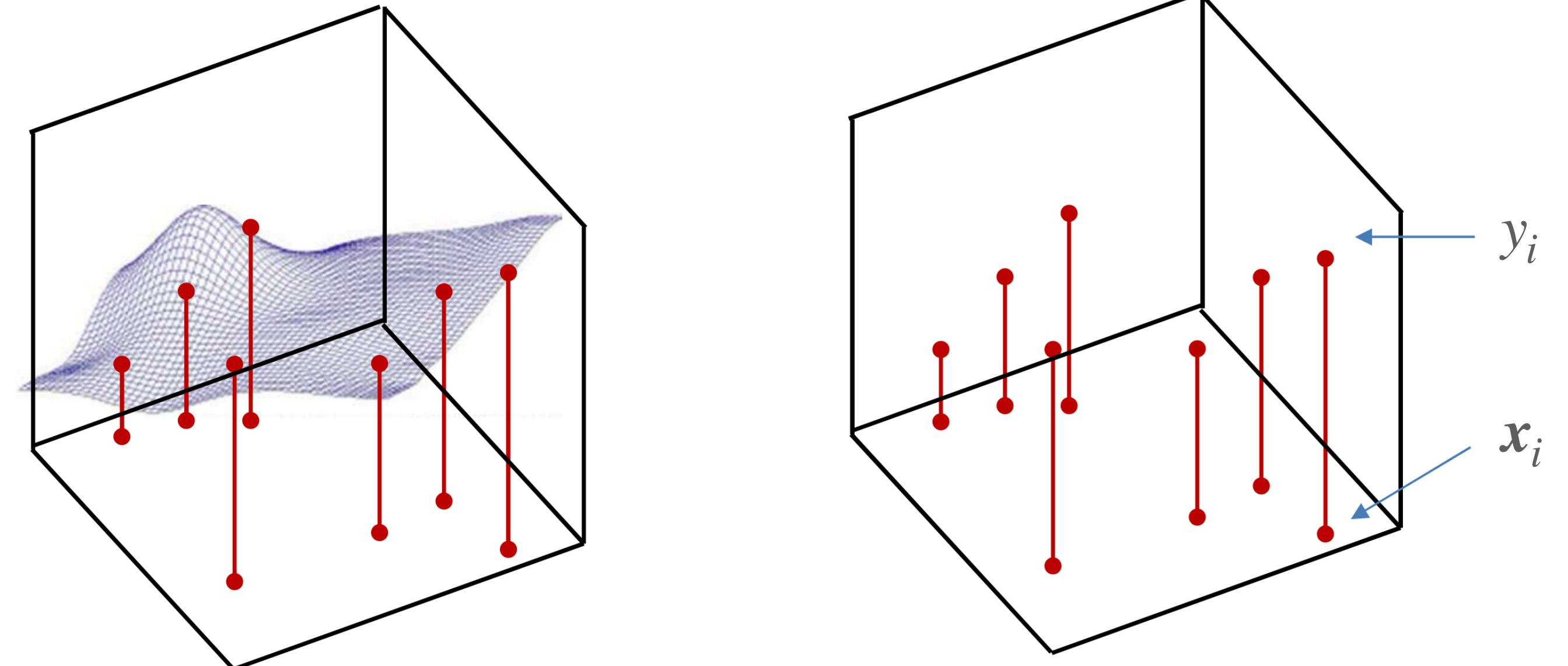
# Learning by sampling

- training set  $X_{\text{train}} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T \in \mathbb{R}^{N \times n}$
- $f(X, \mathbf{W}) \approx g(X)$
- Questions:
  - good sampling?
- Examples:
  - images + labels
  - speech recordings and their transcription



# Learning the function

- Estimate the training **parameters** to *fit* the training data points
- Questions:
  - network architecture?  
basis functions?
  - approximation error?

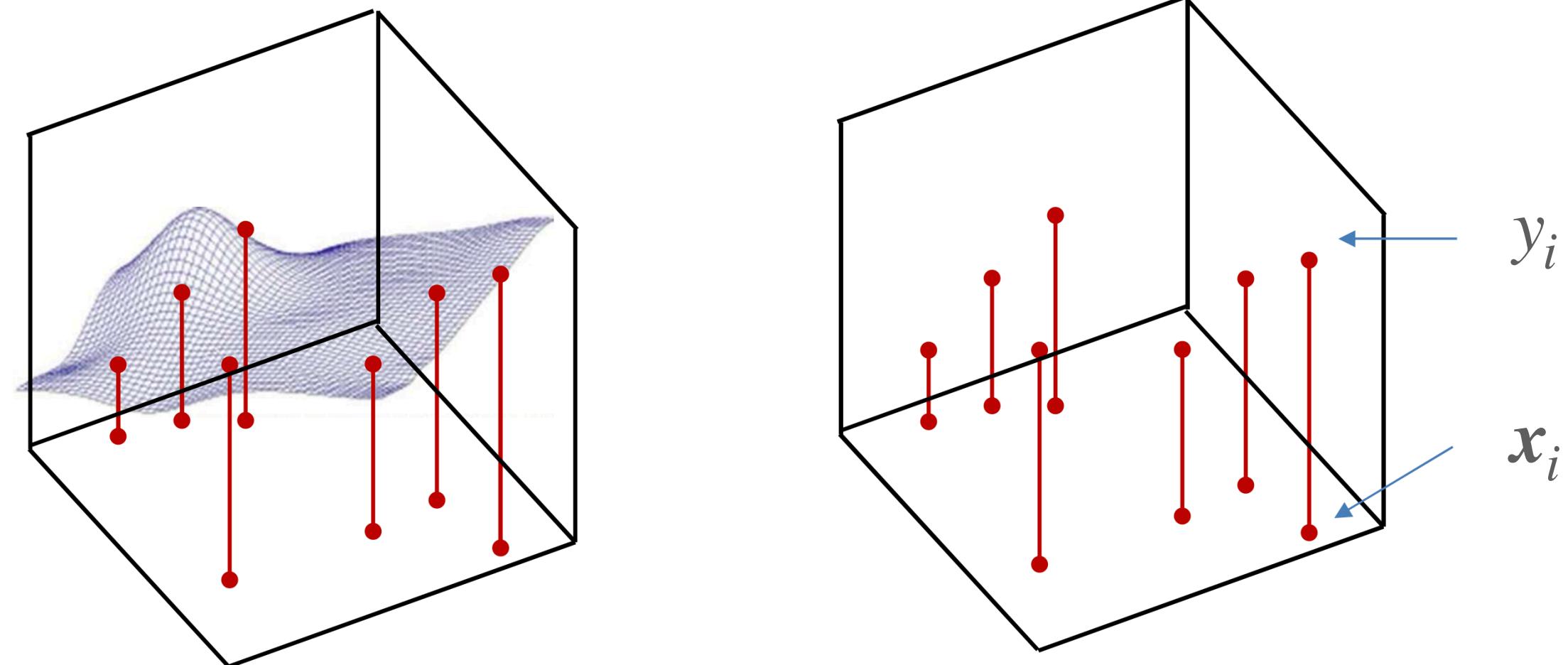


# Learning the function

## Error function (loss function)

- Given the training set  $X_{\text{train}} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T \in \mathbb{R}^{N \times n}$

$$\text{Err}(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^N \text{div}(f(\mathbf{x}_i; \mathbf{W}), y_i)$$



# Logistic regression

## Binary classification

- input data ( $i = 1, 2, \dots, N$ )

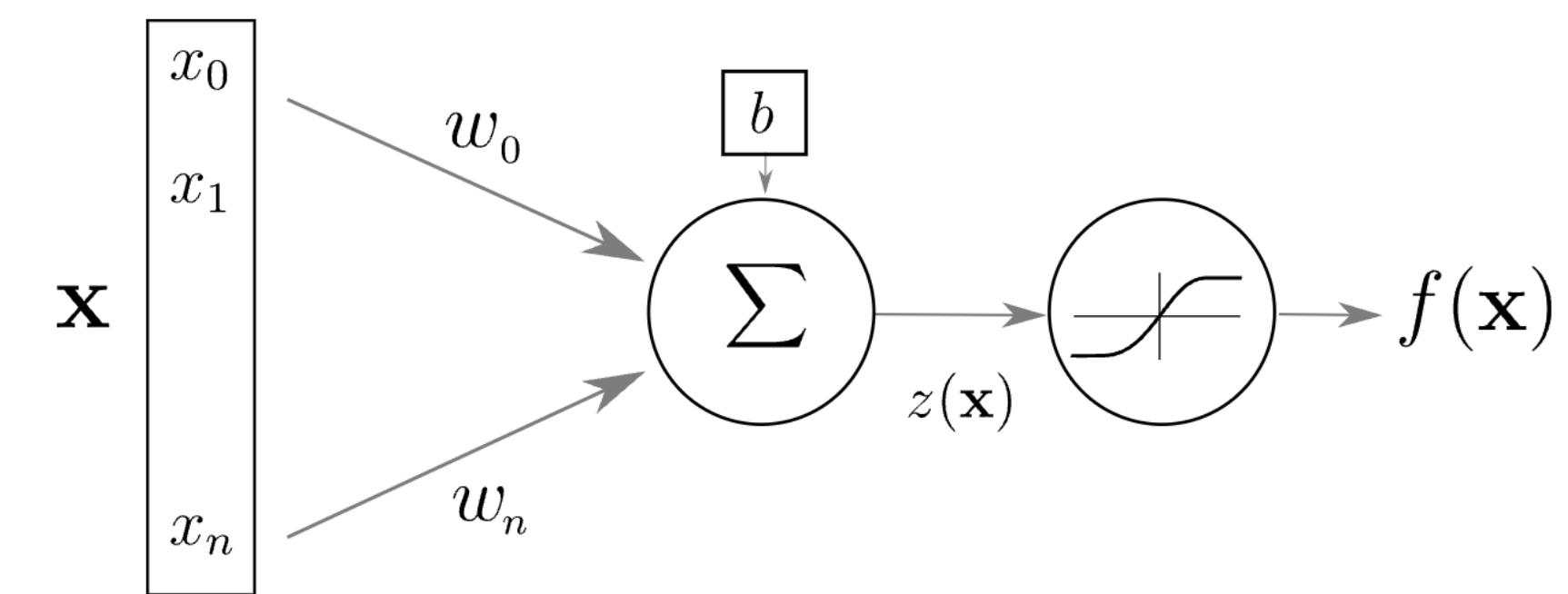
$$\mathbf{x}_i \in \mathbb{R}^n, \quad y_i \in \{0, 1\}$$

- parameters

$$\mathbf{w} \in \mathbb{R}^n, \quad b \in \mathbb{R}$$

- learning function (for given sample point  $(\mathbf{x}, y)$ )

$$f(\mathbf{x}; \mathbf{w}, b) = p(y = 1 | \mathbf{x})$$

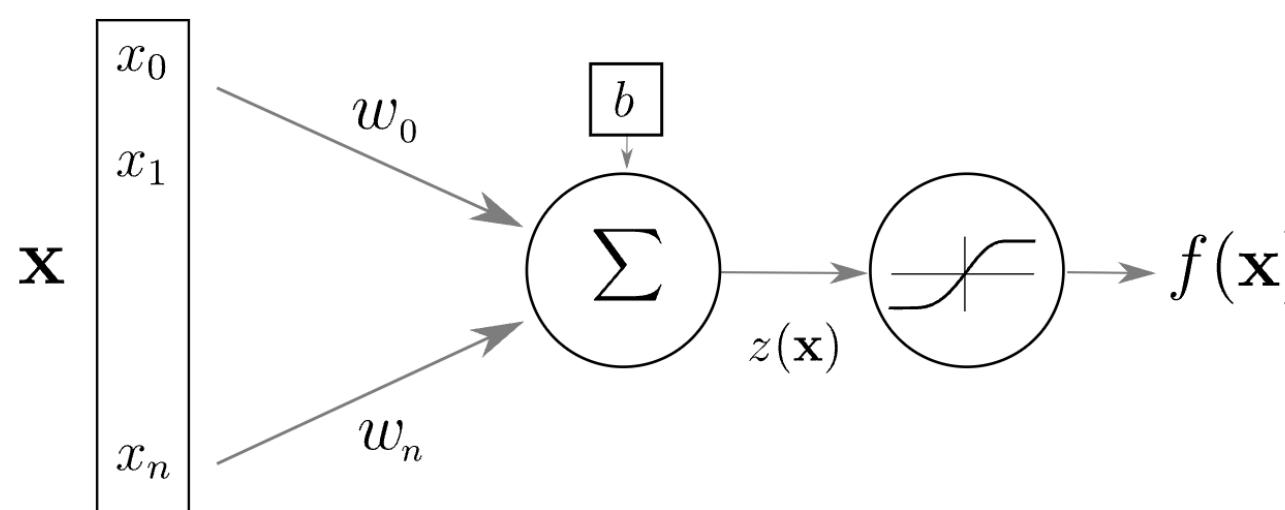


# Logistic regression

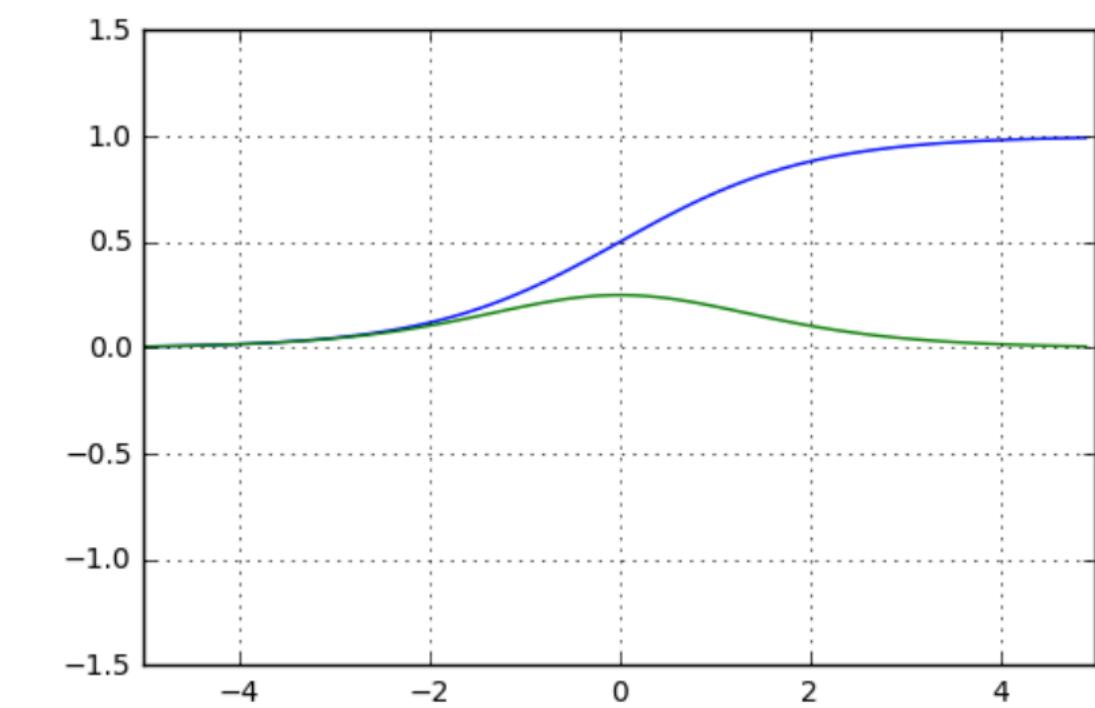
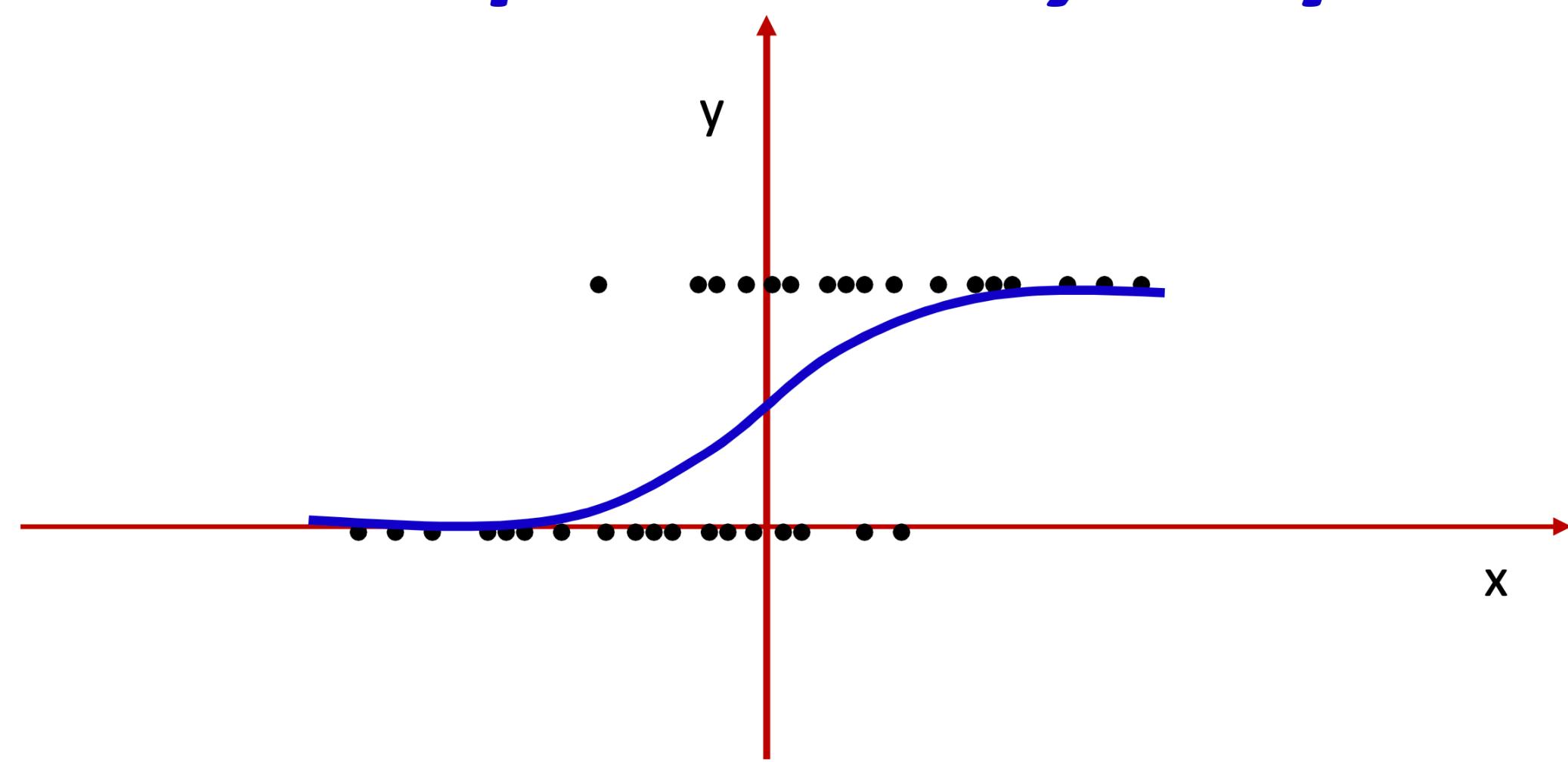
## Binary classification

- sigmoid function

$$g(x) = \sigma(x)$$



The *probability* of  $y=1$



$$\text{sigm}(x) = \frac{1}{1 + e^{-x}}$$

$$\text{sigm}'(x) = \text{sigm}(x)(1 - \text{sigm}(x))$$

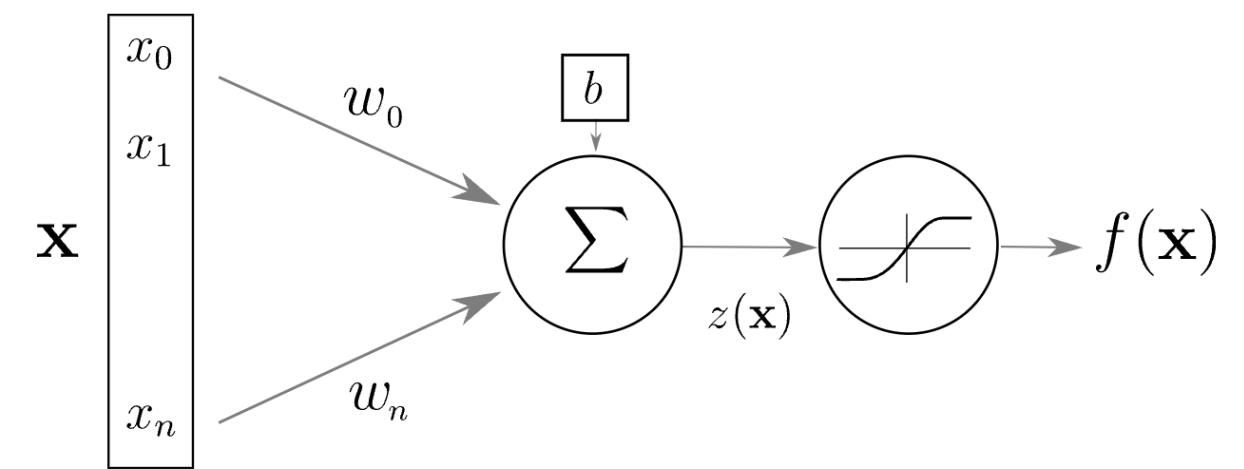
# Logistic regression

## Binary **classification** - loss function

- binary **cross entropy**

$$f(\mathbf{x}; \mathbf{w}, b) = \hat{y} = p(y = 1 | \mathbf{x}) \in (0, 1)$$

$$\text{div}(\hat{y}, y) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$



$$z(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

$$f(\mathbf{x}) = g(\mathbf{w}^T \mathbf{x} + b)$$

- $\mathbf{x}, f(\mathbf{x})$  input and output
- $z(\mathbf{x})$  pre-activation
- $\mathbf{w}, b$  weights and bias

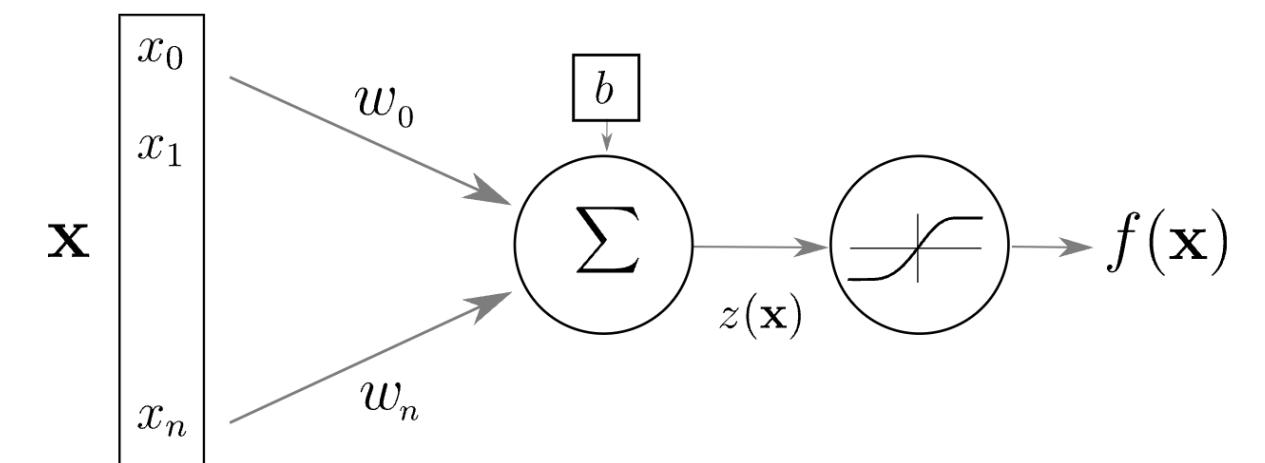
# Linear regression (recap)

## Loss function

- (scaled)  $L_2$  error (squared euclidean distance)

$$f(\mathbf{x}; \mathbf{w}, b) = \hat{y}$$

$$\text{div}(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$$



$$z(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

$$f(\mathbf{x}) = g(\mathbf{w}^T \mathbf{x} + b)$$

- $\mathbf{x}, f(\mathbf{x})$  input and output
- $z(\mathbf{x})$  pre-activation
- $\mathbf{w}, b$  weights and bias

# Logistic regression

## Derivative of loss function

$$z = \mathbf{w}^T \mathbf{x} + b, \quad \hat{y} = \sigma(z), \quad \frac{\partial \hat{y}}{\partial z} = \sigma(z)(1 - \sigma(z))$$

$$L(\mathbf{w}, b) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

$$\frac{\partial L}{\partial \hat{y}} = ? , \quad \frac{\partial L}{\partial z} = ? , \quad \frac{\partial L}{\partial \mathbf{w}} = ? , \quad \frac{\partial L}{\partial b} = ?$$

# Matrix calculus

## Notation

- $x \in \mathbb{R}$
- $\boldsymbol{x} = [x_1, x_2, \dots, x_n]^T \in \mathbb{R}^n$
- $f: \mathbb{R}^n \rightarrow \mathbb{R}, y = f(\boldsymbol{x})$
- $f: \mathbb{R}^n \rightarrow \mathbb{R}^m, \boldsymbol{y} = f(\boldsymbol{x})$

# Matrix calculus

## (Scalar) derivative

- $x \in \mathbb{R}$

$$f: \mathbb{R} \rightarrow \mathbb{R}$$

$$y = f(x), \quad \frac{\partial f}{\partial x} = f'(x)$$

# Matrix calculus

## Gradient (**horizontal vector**)

- $\mathbf{x} = [x_1, x_2, \dots, x_n]^T \in \mathbb{R}^n$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$y = f(\mathbf{x}), \quad \nabla_{\mathbf{x}} y = \nabla y = \frac{\partial y}{\partial \mathbf{x}} = \left[ \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right]$$

# Matrix calculus

## Gradient

- $\boldsymbol{x} = [x_1, x_2, \dots, x_n]^T \in \mathbb{R}^n$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$\boldsymbol{y} = f(\boldsymbol{x}) = \begin{bmatrix} f_1(\boldsymbol{x}) \\ f_2(\boldsymbol{x}) \\ \vdots \\ f_m(\boldsymbol{x}) \end{bmatrix},$$

$$\nabla_{\boldsymbol{x}} \boldsymbol{y} = \frac{\partial \boldsymbol{y}}{\partial \boldsymbol{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1}, \frac{\partial f_1}{\partial x_2}, \dots, \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1}, \frac{\partial f_2}{\partial x_2}, \dots, \frac{\partial f_2}{\partial x_n} \\ \vdots \\ \frac{\partial f_m}{\partial x_1}, \frac{\partial f_m}{\partial x_2}, \dots, \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

# Examples

- $\frac{\partial y}{\partial x}(c) = 0 \quad (c - \text{constant})$
- $\frac{\partial y}{\partial x}(x) = 1$
- $\frac{\partial y}{\partial x}(x) = ? \quad (x \in \mathbb{R}^n)$

# Examples

$$f : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

- $\frac{\partial}{\partial x}(x + w) = ?$

- $\frac{\partial}{\partial x}(x * w) = ?$

- $\frac{\partial}{\partial x}(x * x) = ?$

- $\frac{\partial}{\partial x}(Ax) = ?$

# Examples

$$f : \mathbb{R}^n \rightarrow \mathbb{R}$$

- $\frac{\partial}{\partial \mathbf{x}} \text{sum}(\mathbf{x}) = ?$
- $\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{w}) = ?$
- $\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{x}) = ?$
- $\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T A \mathbf{x}) = ?$

# Total derivative chain rule

$$\bullet \frac{\partial f(u)}{\partial x} = \frac{\partial f(u_1, u_2, \dots, u_n)}{\partial x} = \sum_i \frac{\partial f}{\partial u_i} \frac{\partial u_i}{\partial x} \quad (u_i = u_i(x), \quad x \in \mathbb{R})$$

# Chain rule

$$x \in \mathbb{R}$$

$$\begin{aligned} y = f(g(x)) = f \circ g(x) &= \begin{bmatrix} f_1(\mathbf{g}(x)) \\ f_2(\mathbf{g}(x)) \\ \vdots \\ f_n(\mathbf{g}(x)) \end{bmatrix} \\ \frac{\partial y}{\partial x} &= \begin{bmatrix} \frac{\partial f_1(\mathbf{g}(x))}{\partial x} \\ \frac{\partial f_2(\mathbf{g}(x))}{\partial x} \\ \vdots \\ \frac{\partial f_n(\mathbf{g}(x))}{\partial x} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial g_1} \frac{\partial g_1}{\partial x} + \frac{\partial f_1}{\partial g_2} \frac{\partial g_2}{\partial x} + \cdots + \frac{\partial f_1}{\partial g_n} \frac{\partial g_n}{\partial x} \\ \frac{\partial f_2}{\partial g_1} \frac{\partial g_1}{\partial x} + \frac{\partial f_2}{\partial g_2} \frac{\partial g_2}{\partial x} + \cdots + \frac{\partial f_2}{\partial g_n} \frac{\partial g_n}{\partial x} \\ \vdots \\ \frac{\partial f_n}{\partial g_1} \frac{\partial g_1}{\partial x} + \frac{\partial f_n}{\partial g_2} \frac{\partial g_2}{\partial x} + \cdots + \frac{\partial f_n}{\partial g_n} \frac{\partial g_n}{\partial x} \end{bmatrix} = \frac{\partial f}{\partial \mathbf{g}} \cdot \frac{\partial \mathbf{g}}{\partial x} \quad (\text{matrix} \cdot \text{vector}) \end{aligned}$$

# Chain rule

$$x \in \mathbb{R}^n$$

$$y = f(g(x)) = f \circ g(x) = \begin{bmatrix} f_1(\mathbf{g}(x)) \\ f_2(\mathbf{g}(x)) \\ \vdots \\ f_n(\mathbf{g}(x)) \end{bmatrix}$$
$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial f_1(\mathbf{g}(x))}{\partial x} \\ \frac{\partial f_2(\mathbf{g}(x))}{\partial x} \\ \vdots \\ \frac{\partial f_n(\mathbf{g}(x))}{\partial x} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial g_1} \frac{\partial g_1}{\partial x} + \frac{\partial f_1}{\partial g_2} \frac{\partial g_2}{\partial x} + \cdots + \frac{\partial f_1}{\partial g_n} \frac{\partial g_n}{\partial x} \\ \frac{\partial f_2}{\partial g_1} \frac{\partial g_1}{\partial x} + \frac{\partial f_2}{\partial g_2} \frac{\partial g_2}{\partial x} + \cdots + \frac{\partial f_2}{\partial g_n} \frac{\partial g_n}{\partial x} \\ \vdots \\ \frac{\partial f_n}{\partial g_1} \frac{\partial g_1}{\partial x} + \frac{\partial f_n}{\partial g_2} \frac{\partial g_2}{\partial x} + \cdots + \frac{\partial f_n}{\partial g_n} \frac{\partial g_n}{\partial x} \end{bmatrix} = \frac{\partial f}{\partial \mathbf{g}} \cdot \frac{\partial \mathbf{g}}{\partial x} \quad (\text{matrix} \cdot \text{matrix})$$

# Logistic regression

## Derivative of loss function

$$z = \mathbf{w}^T \mathbf{x} + b, \quad \hat{y} = \sigma(z), \quad \frac{\partial \hat{y}}{\partial z} = \sigma(z)(1 - \sigma'(z))$$

$$L(\mathbf{w}, b) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

$$\frac{\partial L}{\partial \hat{y}} = ? , \quad \frac{\partial L}{\partial z} = ? , \quad \frac{\partial L}{\partial \mathbf{w}} = ? , \quad \frac{\partial L}{\partial b} = ?$$

# Logistic regression

## Derivative of loss function

$$z = \mathbf{w}^T \mathbf{x} + b, \quad \hat{y} = \sigma(z), \quad \frac{\partial \hat{y}}{\partial z} = \sigma(z)(1 - \sigma'(z))$$

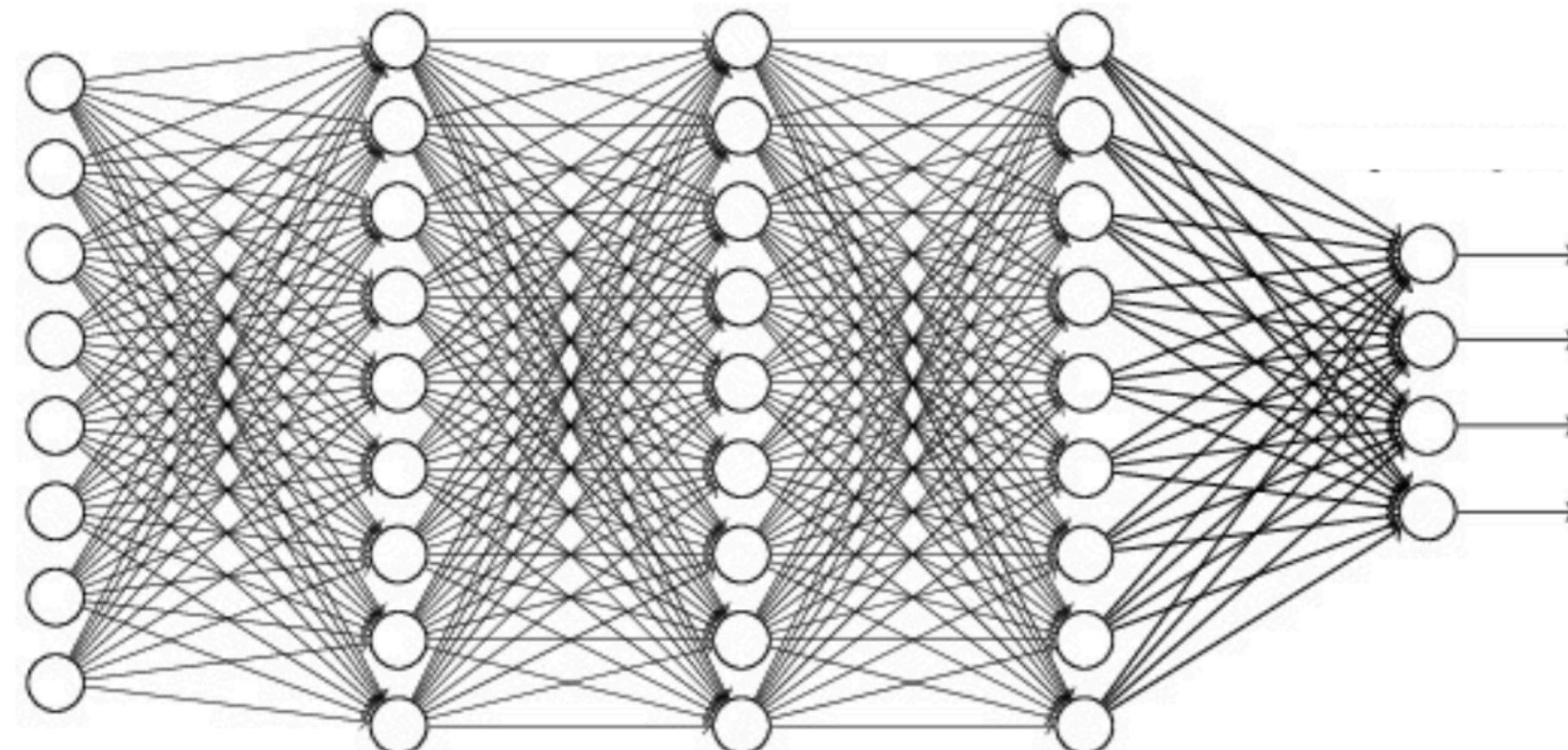
$$L(\mathbf{w}, b) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

$$\frac{\partial L}{\partial \hat{y}} = -\frac{y}{\hat{y}} + \frac{1 - y}{1 - \hat{y}}, \quad \frac{\partial L}{\partial z} = \hat{y} - y, \quad \frac{\partial L}{\partial \mathbf{w}} = (\hat{y} - y) \mathbf{x}^T, \quad \frac{\partial L}{\partial b} = \hat{y} - y$$

# Multiclass classification

$k$  - classes

- Label (of class  $c$ )  $y = [0,0,\dots,0,1,0,0,\dots,0]$  is one hot vector with value 1 at  $c$ -th position
- Network output  $f(x) = \hat{y} \in \mathbb{R}^k$  is probability distribution over  $k$  classes  
 $\hat{y} = \text{softmax}(z)$



$$\hat{y}_i = \frac{\exp(z_i)}{\sum_{j=1}^k \exp(z_j)}$$

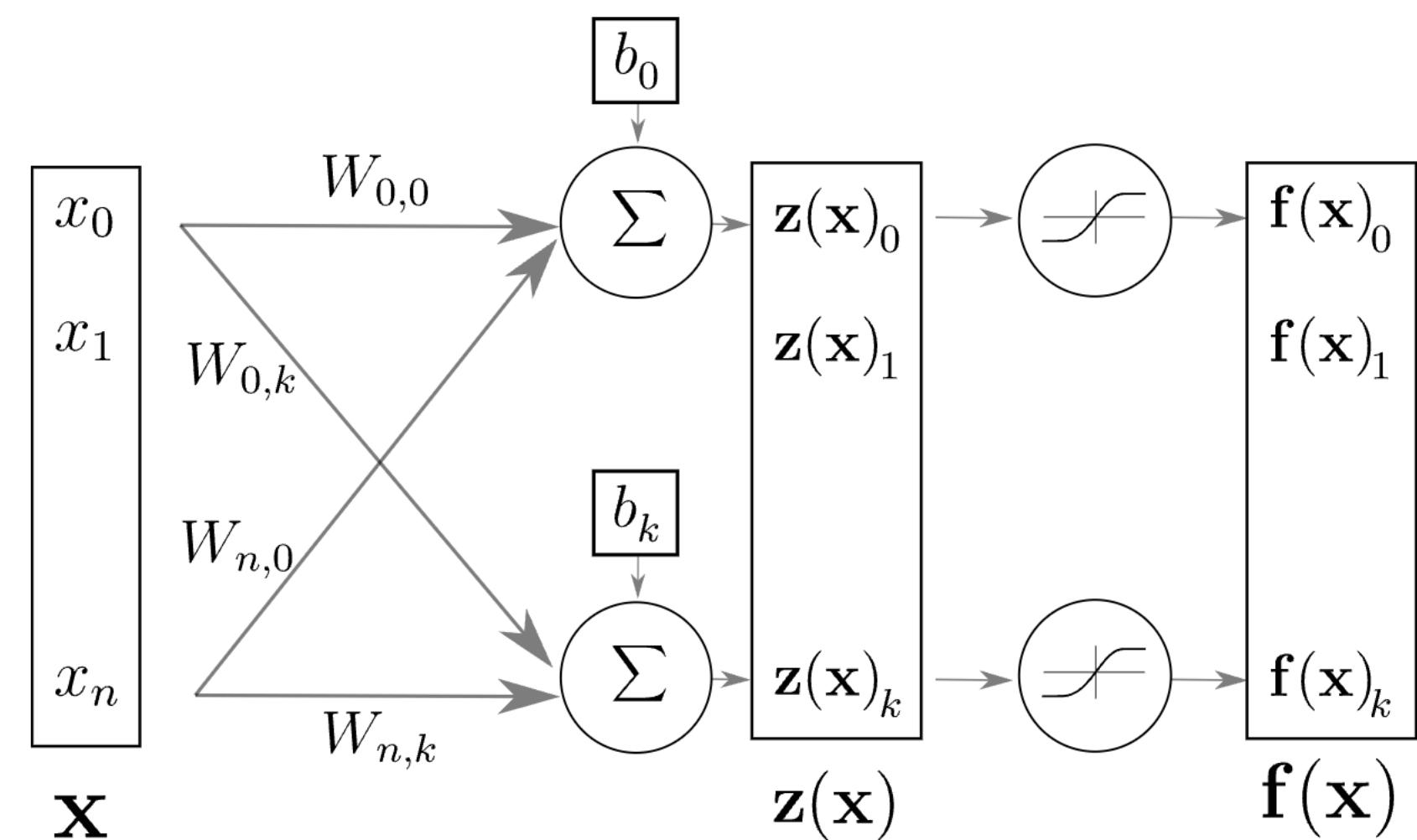
# Deep Learning calculus

$$\hat{y} = \text{softmax}(z), \quad \frac{\partial \hat{y}}{\partial z} = \text{diag}(\hat{y}) - \hat{y}\hat{y}^T$$

... to be continued next week

# Forward propagation

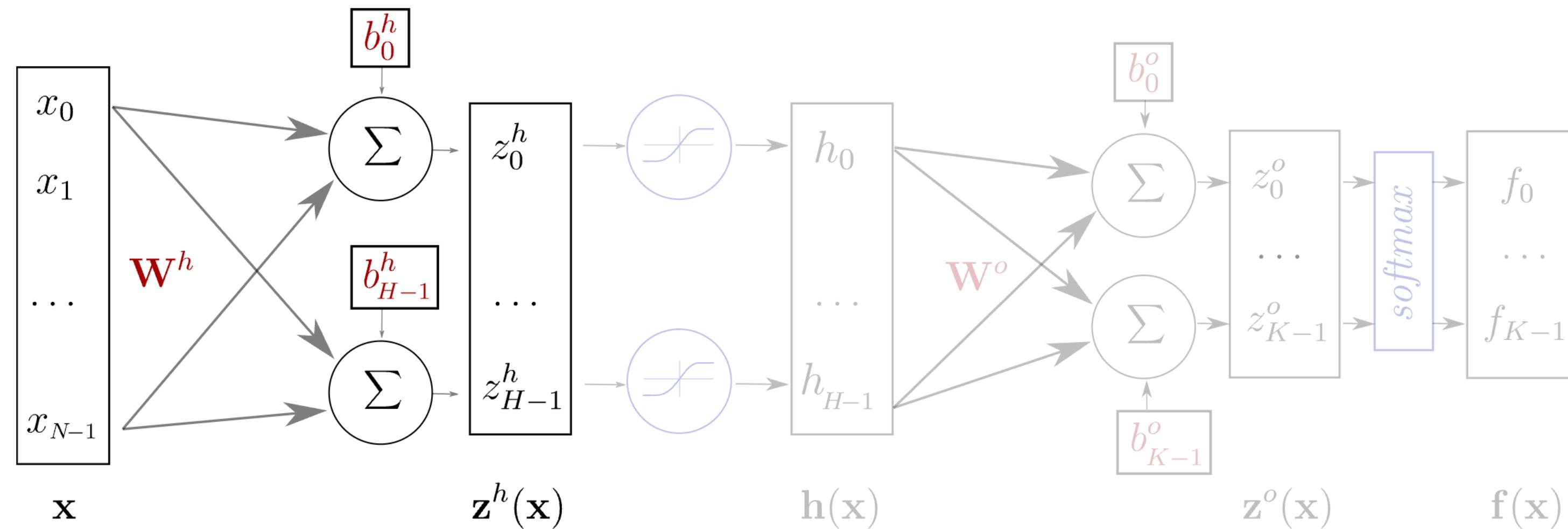
# Layer of Neurons



$$\mathbf{f}(\mathbf{x}) = g(\mathbf{z}(\mathbf{x})) = g(\mathbf{W}\mathbf{x} + \mathbf{b})$$

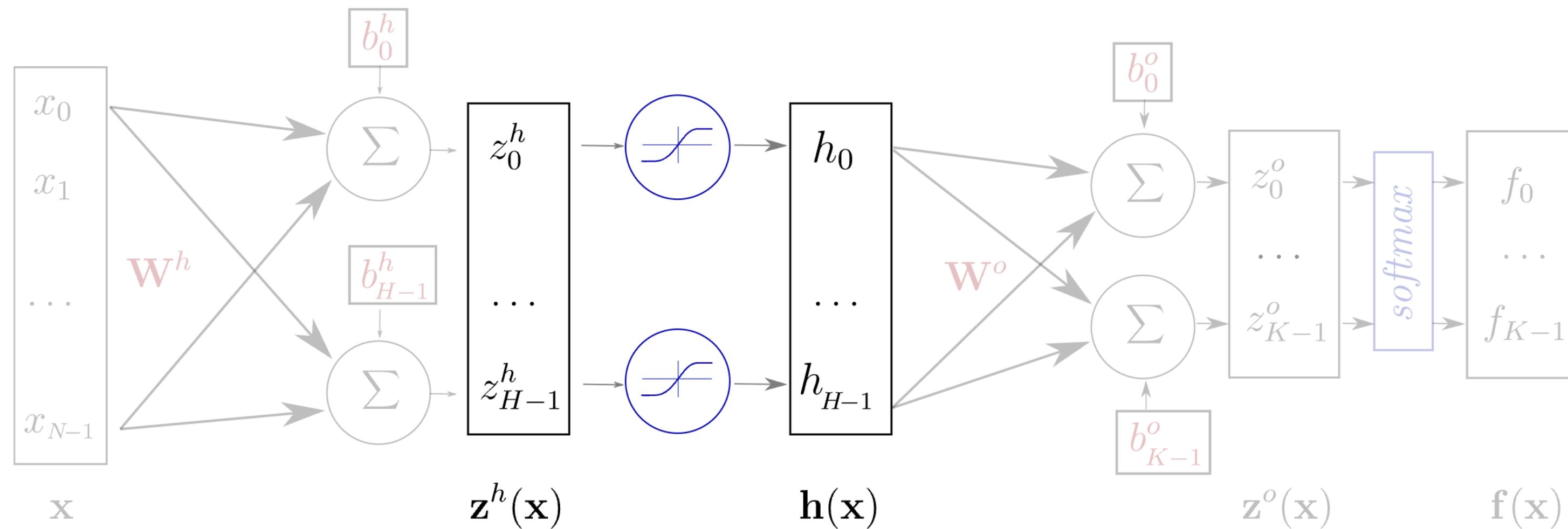
- $\mathbf{W}, \mathbf{b}$  now matrix and vector

# One Hidden Layer Network



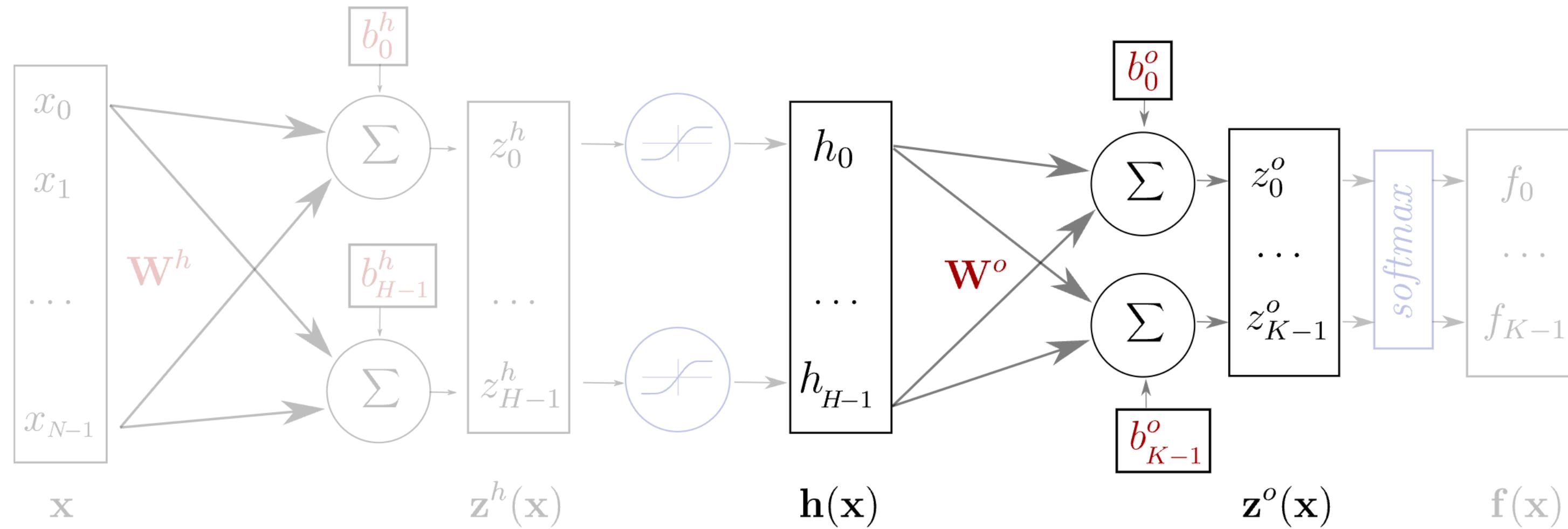
- $\mathbf{z}^h(\mathbf{x}) = \mathbf{W}^h \mathbf{x} + \mathbf{b}^h$
- $\mathbf{h}(\mathbf{x}) = g(\mathbf{z}^h(\mathbf{x})) = g(\mathbf{W}^h \mathbf{x} + \mathbf{b}^h)$
- $\mathbf{z}^o(\mathbf{x}) = \mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o$
- $\mathbf{f}(\mathbf{x}) = \text{softmax}(\mathbf{z}^o) = \text{softmax}(\mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o)$

# One Hidden Layer Network



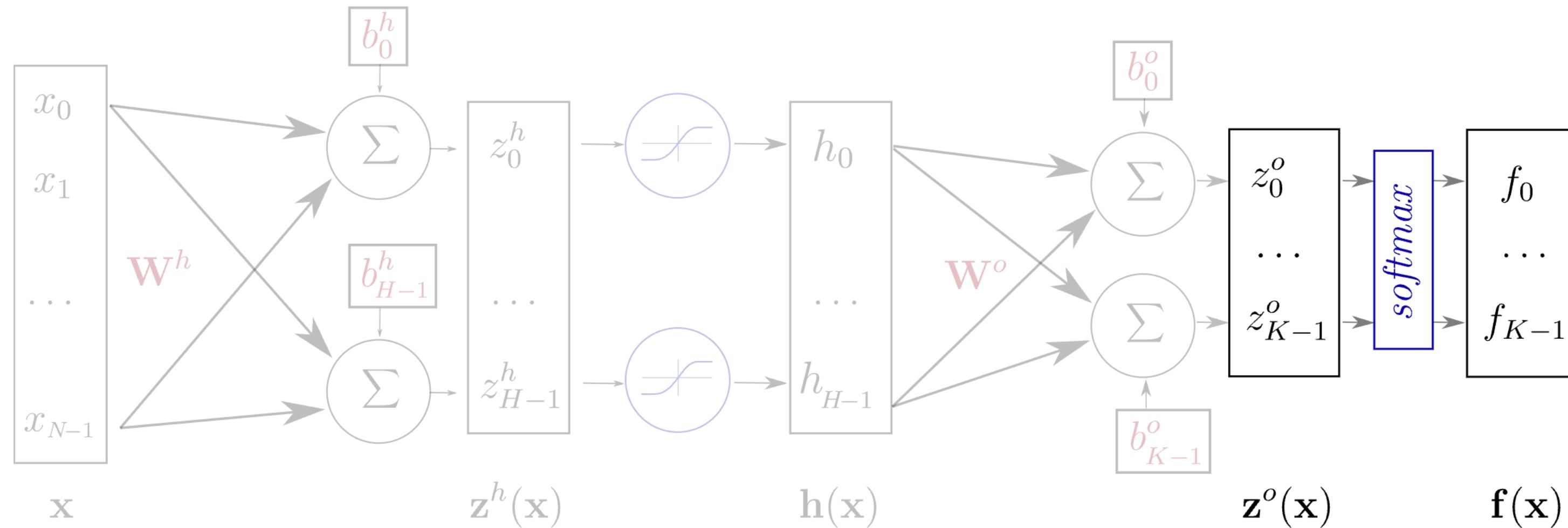
- $\mathbf{z}^h(\mathbf{x}) = \mathbf{W}^h \mathbf{x} + \mathbf{b}^h$
- $\mathbf{h}(\mathbf{x}) = g(\mathbf{z}^h(\mathbf{x})) = g(\mathbf{W}^h \mathbf{x} + \mathbf{b}^h)$
- $\mathbf{z}^o(\mathbf{x}) = \mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o$
- $\mathbf{f}(\mathbf{x}) = \text{softmax}(\mathbf{z}^o) = \text{softmax}(\mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o)$

# One Hidden Layer Network



- $\mathbf{z}^h(\mathbf{x}) = \mathbf{W}^h \mathbf{x} + \mathbf{b}^h$
- $\mathbf{h}(\mathbf{x}) = g(\mathbf{z}^h(\mathbf{x})) = g(\mathbf{W}^h \mathbf{x} + \mathbf{b}^h)$
- $\mathbf{z}^o(\mathbf{x}) = \mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o$
- $\mathbf{f}(\mathbf{x}) = \text{softmax}(\mathbf{z}^o) = \text{softmax}(\mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o)$

# One Hidden Layer Network



- $\mathbf{z}^h(\mathbf{x}) = \mathbf{W}^h \mathbf{x} + \mathbf{b}^h$
- $\mathbf{h}(\mathbf{x}) = g(\mathbf{z}^h(\mathbf{x})) = g(\mathbf{W}^h \mathbf{x} + \mathbf{b}^h)$
- $\mathbf{z}^o(\mathbf{x}) = \mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o$
- $\mathbf{f}(\mathbf{x}) = \text{softmax}(\mathbf{z}^o) = \text{softmax}(\mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o)$