

# D1EAD – Análise Estatística para Ciência de Dados



## Logistic Regression

Lecture based on  
Machine Learning course  
by Andrew Ng  
(Coursera)

Prof. Samuel Martins (Samuka)  
[samuel.martins@ifsp.edu.br](mailto:samuel.martins@ifsp.edu.br)



# Binary Classification

**Admission Exam:** Pass / Fail

**Email:** Spam / Non-spam

**Tumor:** Malignant / Benign

**Signature:** Real / False

**App:** Free / Paid

....

# Binary Classification

**Admission Exam:** Pass / Fail

**Email:** Spam / Non-spam

**Tumor:** Malignant / Benign

**Signature:** Real / False

**App:** Free / Paid

....



Categorical variables  
with only 2 values

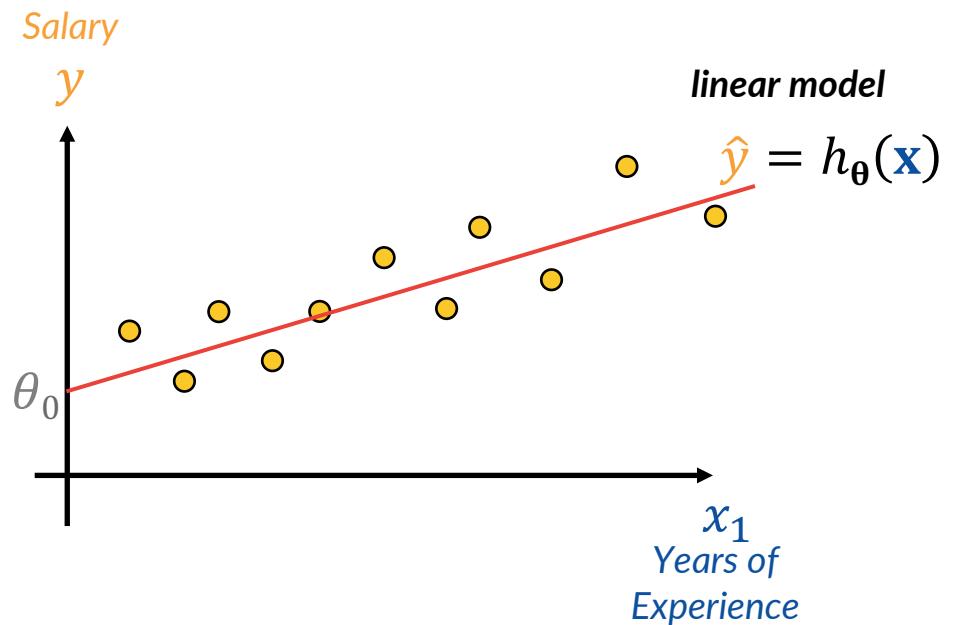
$$y \in \{0, 1\}$$

0: Negative class

1: Positive class

# Linear Regression

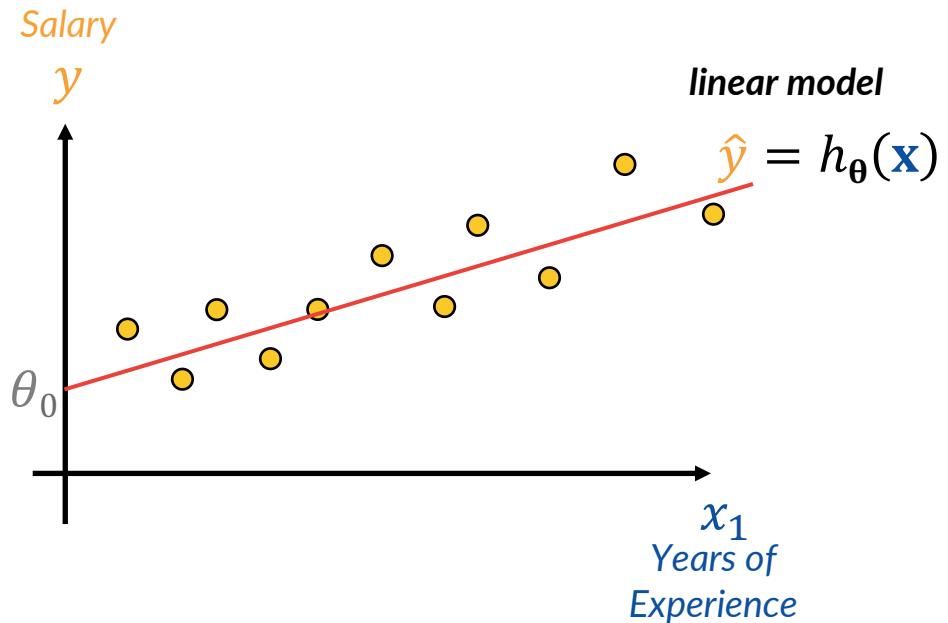
$$\hat{y} = h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 * x_1 + \dots + \theta_n * x_n$$



# Linear Regression

$$\hat{y} = h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 * x_1 + \dots + \theta_n * x_n$$

In Machine Learning, vectors are typically represented as **column vectors**.



$$\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}_{n \times 1} \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}_{n \times 1}$$

# Linear Regression

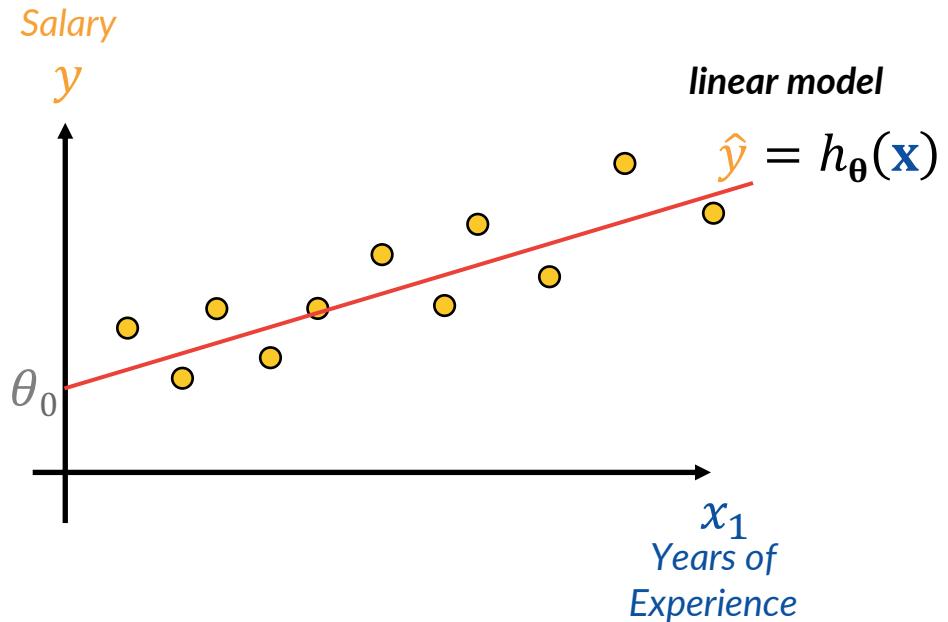
$$\hat{y} = h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 * x_1 + \cdots + \theta_n * x_n$$

$$\hat{y} = h_{\theta}(\mathbf{x}) = \theta_0 + \boldsymbol{\theta}^T * \mathbf{x}$$

In Machine Learning, vectors are typically represented as **column vectors**.

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}_{n \times 1} \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}_{n \times 1}$$

$$\boldsymbol{\theta}^T = [\theta_1 \ \theta_2 \cdots \theta_n]_{1 \times n}$$

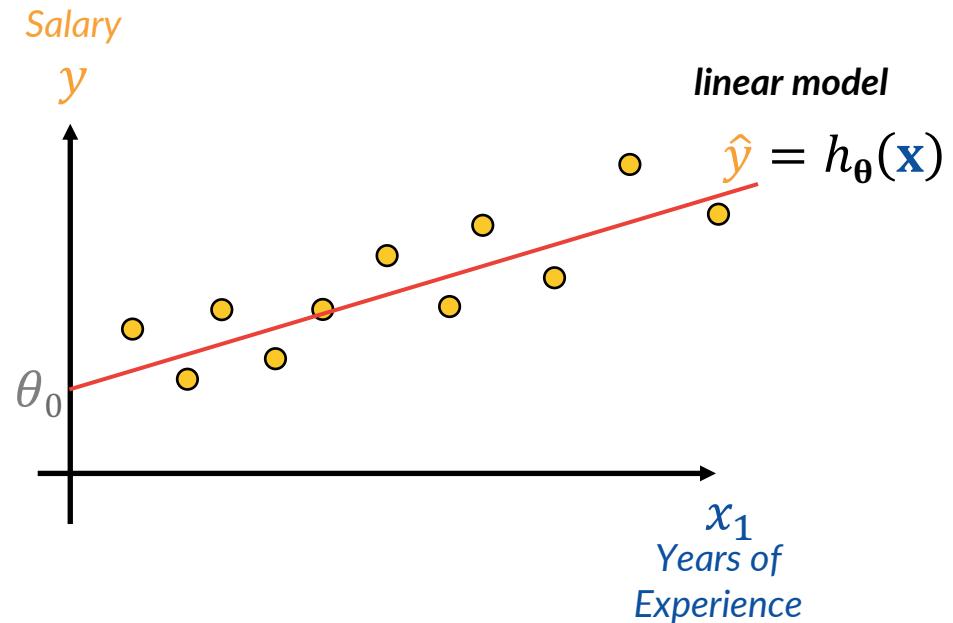


# Linear Regression: Bias Trick

$$\hat{y} = h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 * x_1 + \cdots + \theta_n * x_n$$

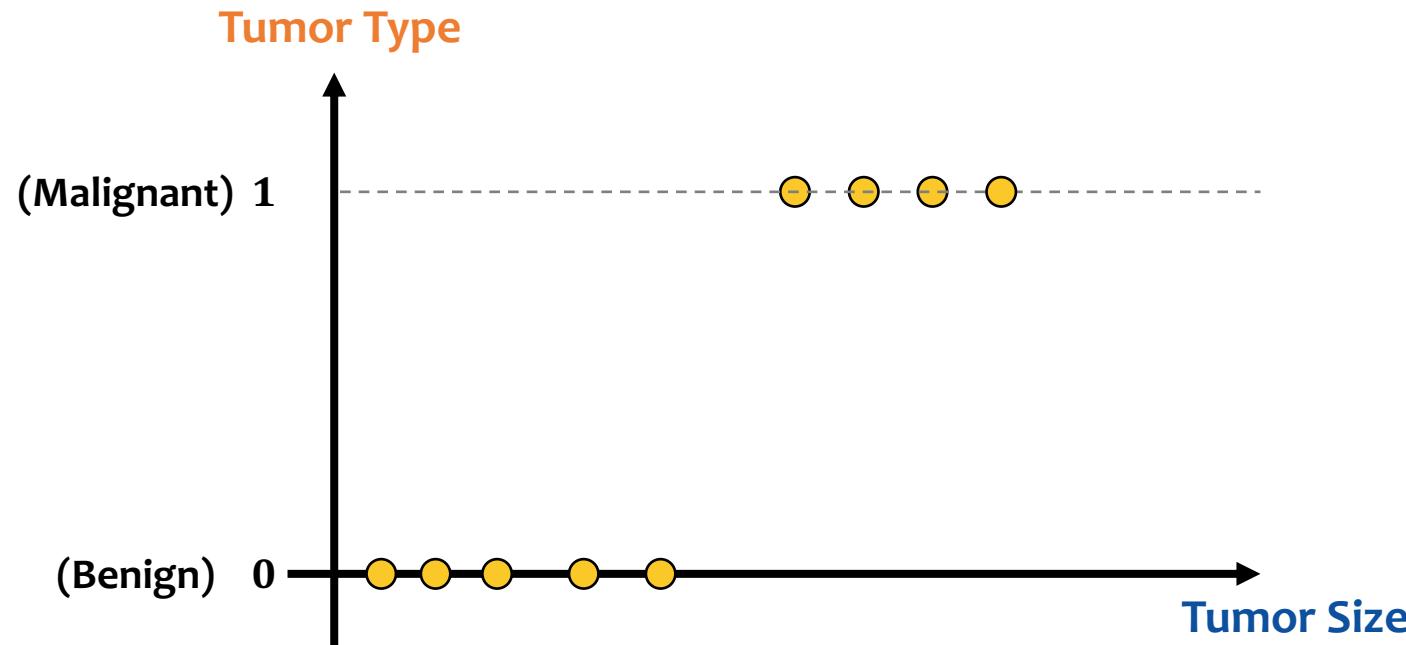
$$\hat{y} = h_{\theta}(\mathbf{x}) = \boldsymbol{\theta}^T * \mathbf{x}$$

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}_{(n+1) \times 1} \quad \mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}_{(n+1) \times 1}$$

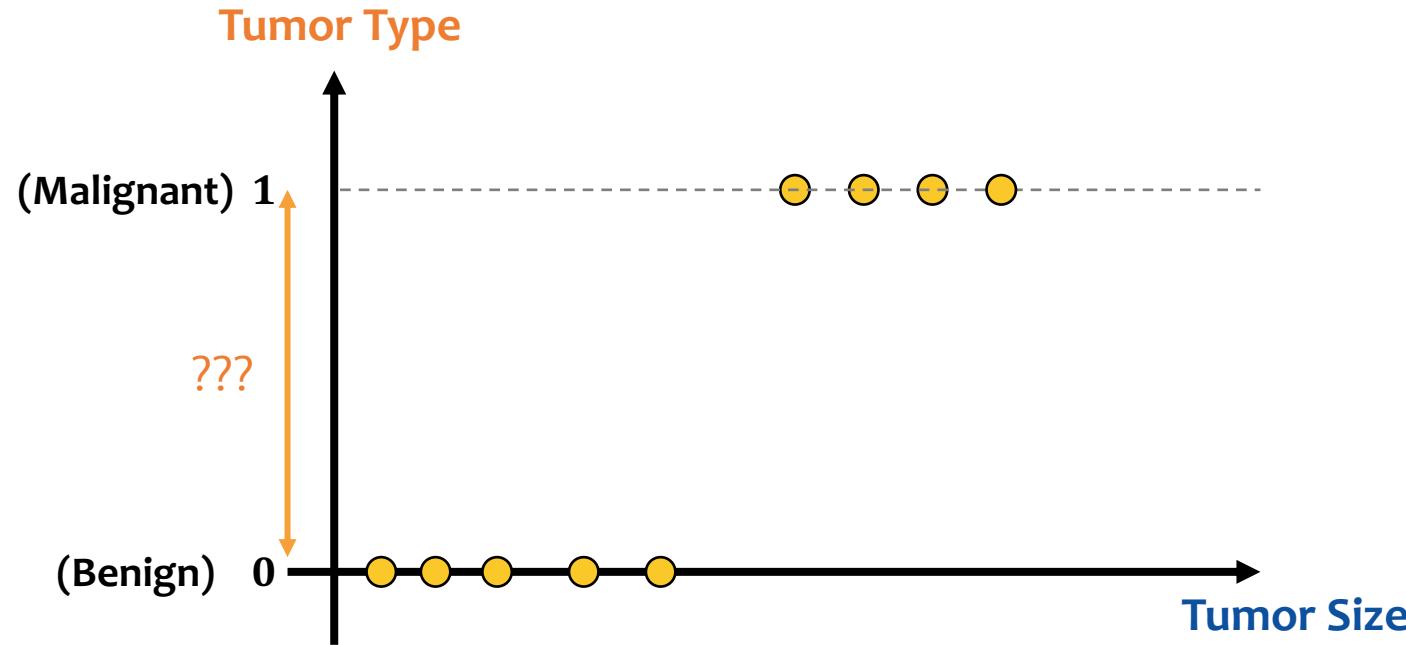


$$\boldsymbol{\theta}^T = [\theta_0 \ \theta_1 \ \theta_2 \cdots \theta_n]_{1 \times (n+1)}$$

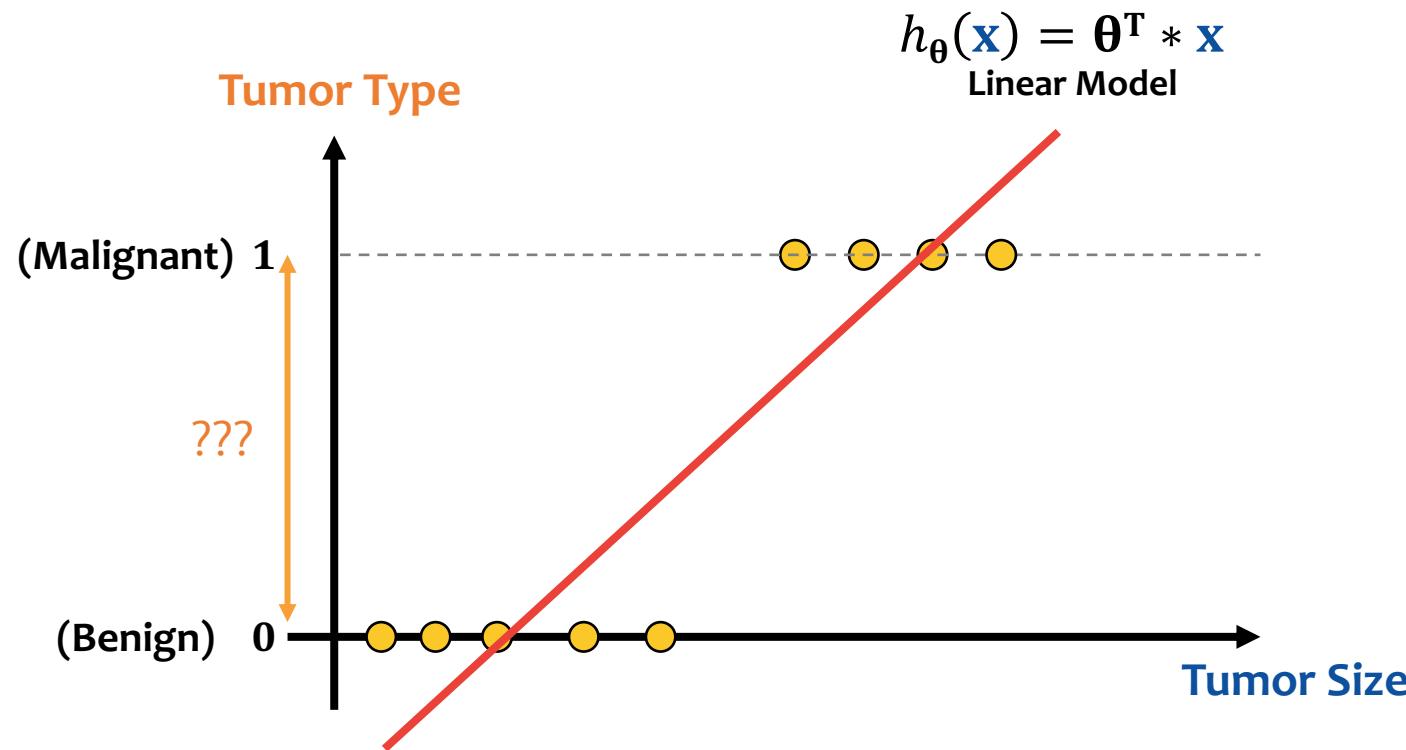
# Binary Classification Problem



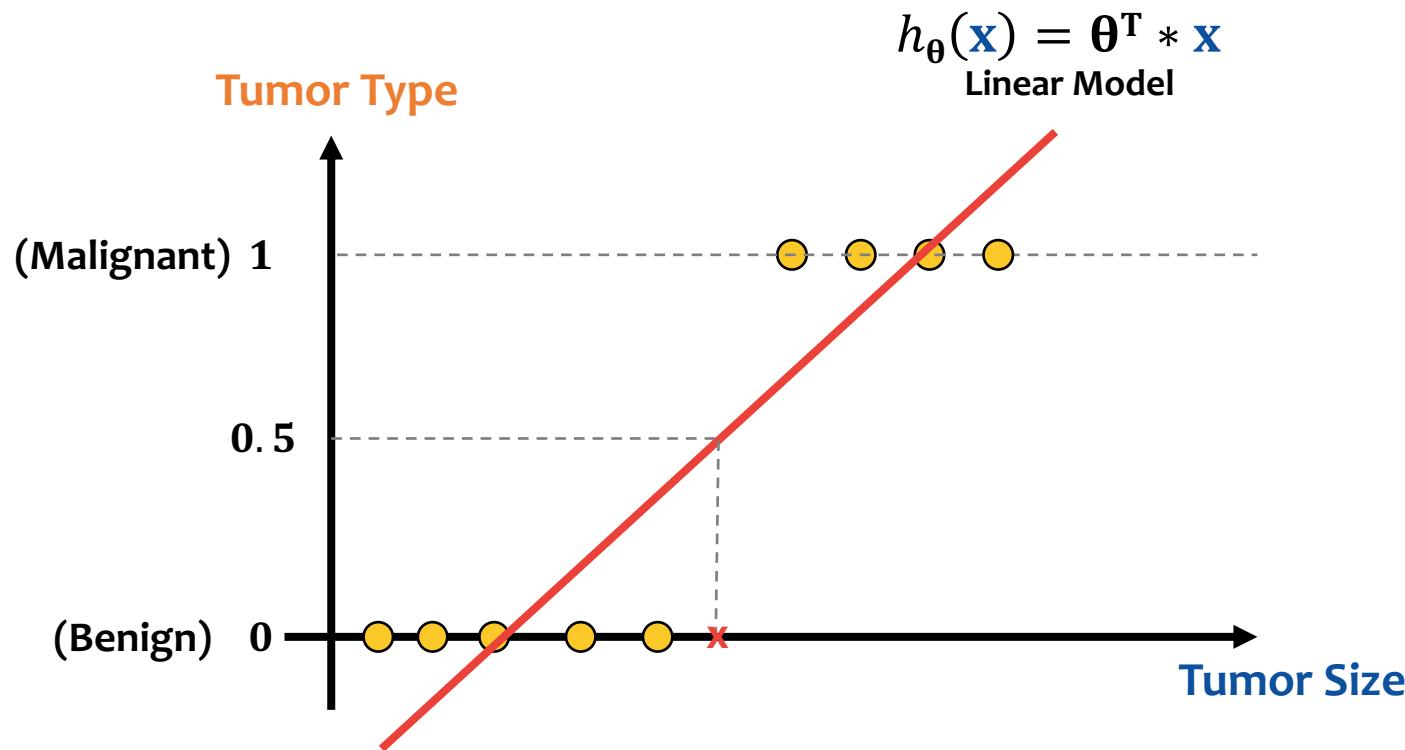
# Binary Classification Problem



# Binary Classification Problem



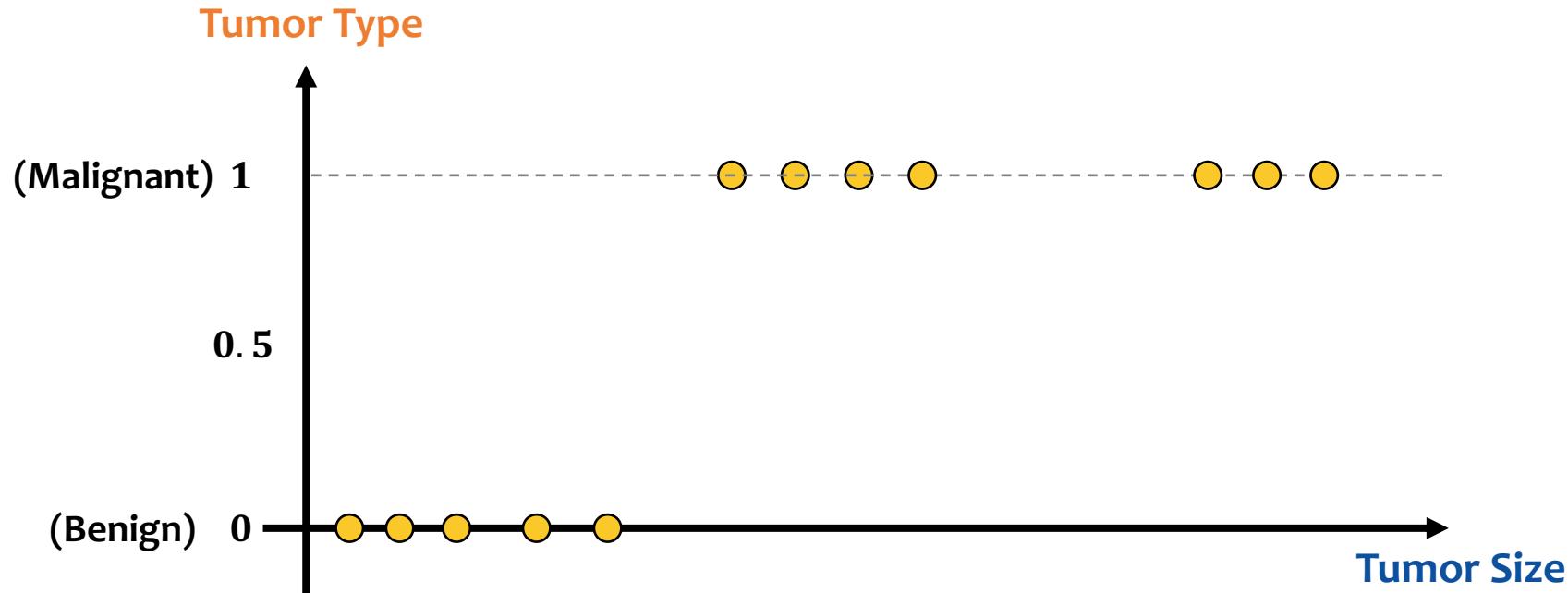
# Binary Classification Problem



If  $h_{\theta}(\mathbf{x}) \geq 0.5$ , classify as  $\hat{y} = 1$

If  $h_{\theta}(\mathbf{x}) < 0$ , classify as  $\hat{y} = 0$

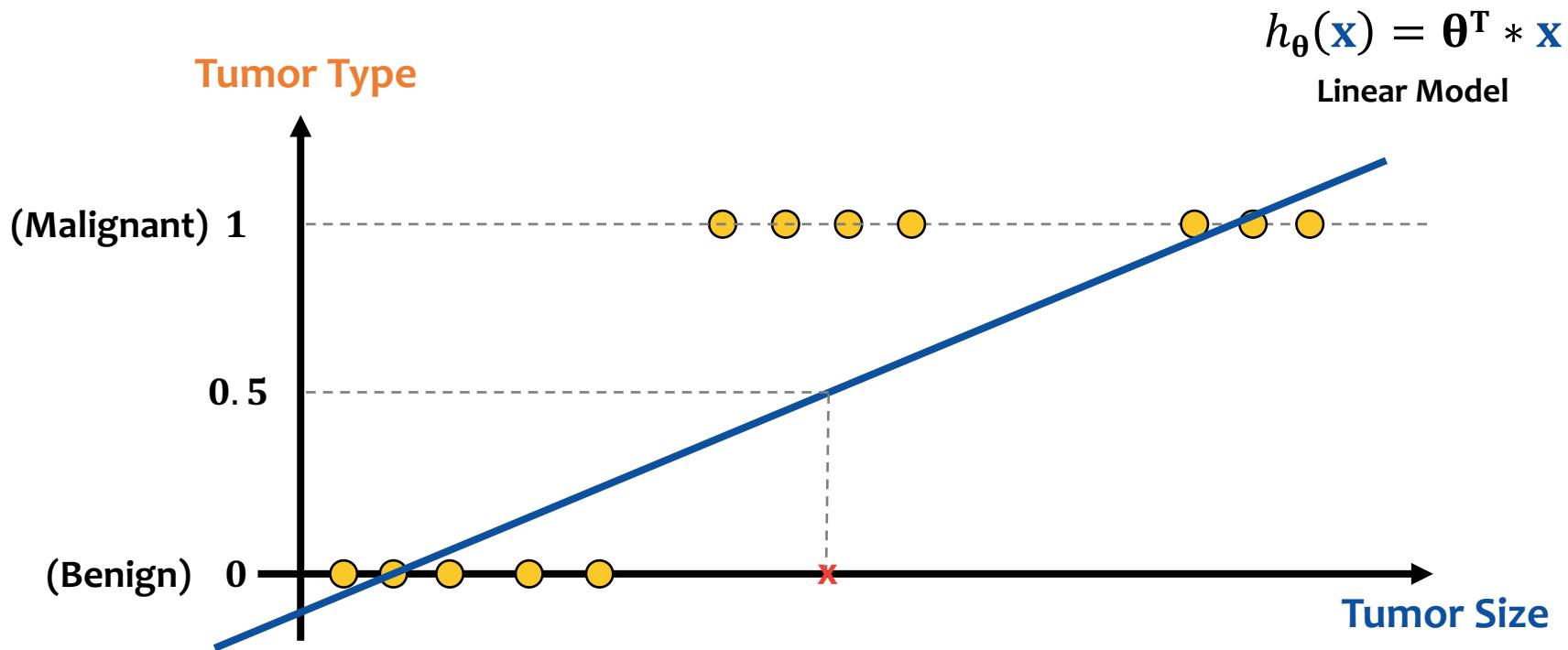
# Binary Classification Problem



If  $h_{\theta}(\mathbf{x}) \geq 0.5$ , classify as  $\hat{y} = 1$

If  $h_{\theta}(\mathbf{x}) < 0$ , classify as  $\hat{y} = 0$

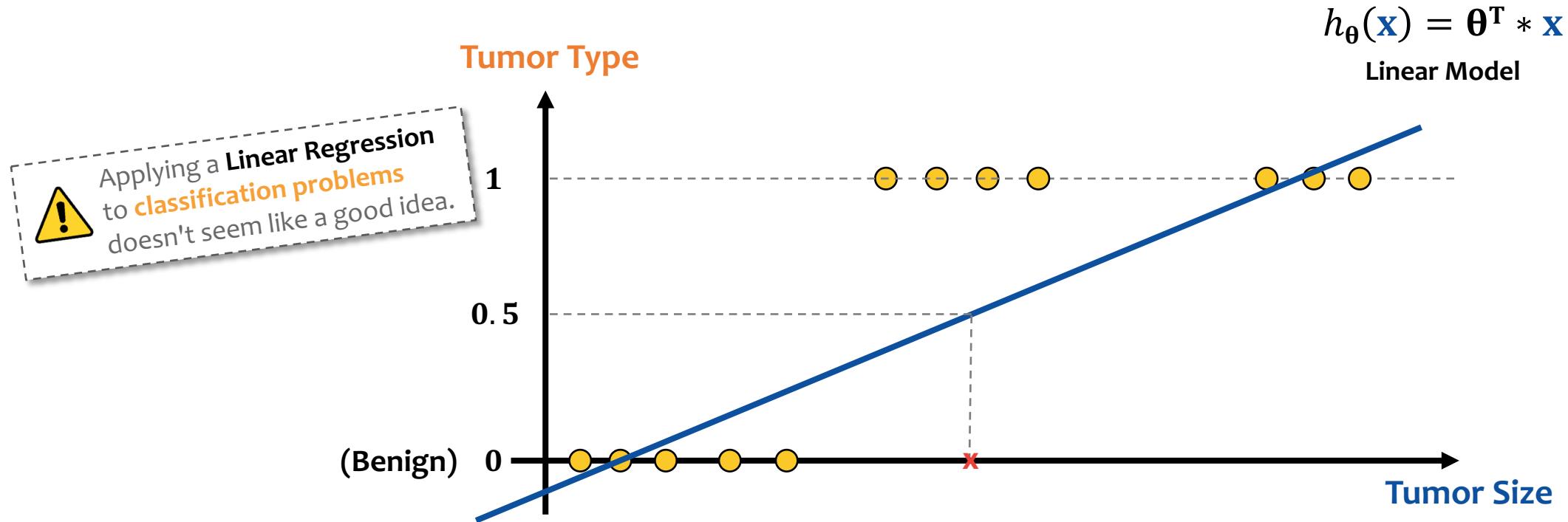
# Binary Classification Problem



If  $h_{\theta}(\mathbf{x}) \geq 0.5$ , classify as  $\hat{y} = 1$

If  $h_{\theta}(\mathbf{x}) < 0$ , classify as  $\hat{y} = 0$

# Binary Classification Problem



If  $h_{\theta}(\mathbf{x}) \geq 0.5$ , classify as  $\hat{y} = 1$

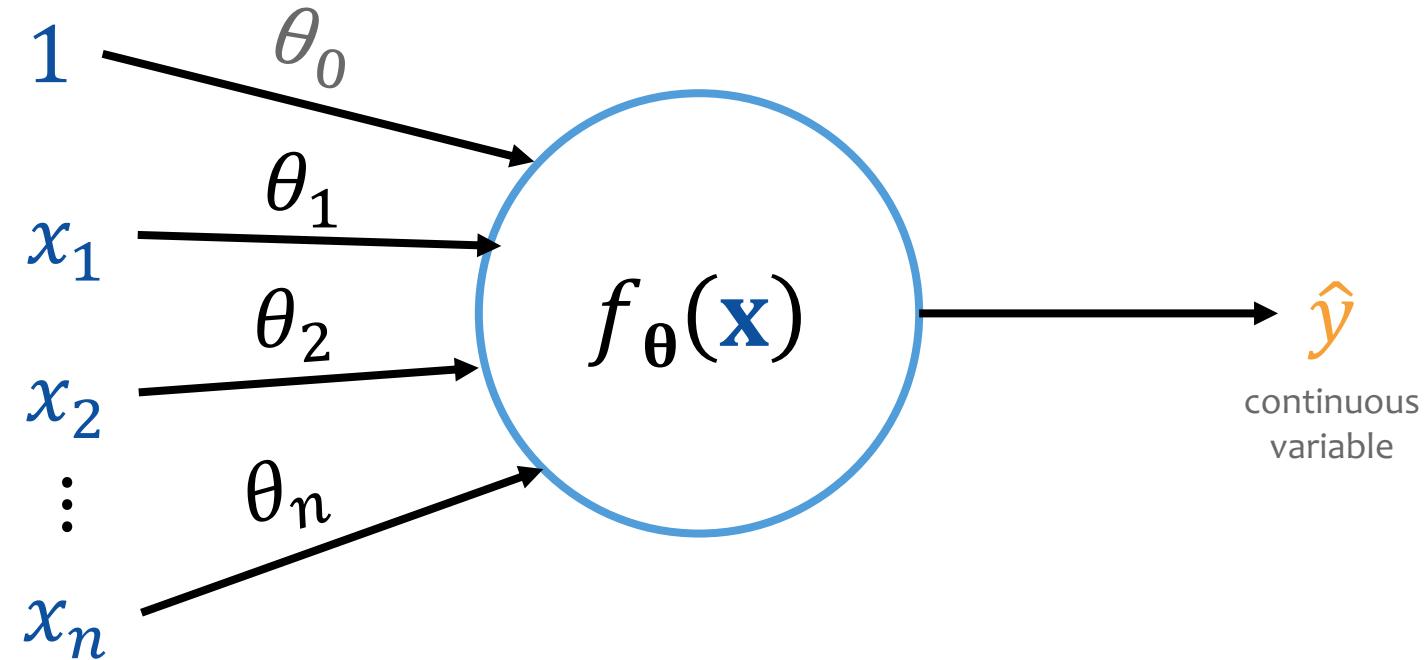
If  $h_{\theta}(\mathbf{x}) < 0$ , classify as  $\hat{y} = 0$

# **From Linear Regression to Binary Linear Classification**

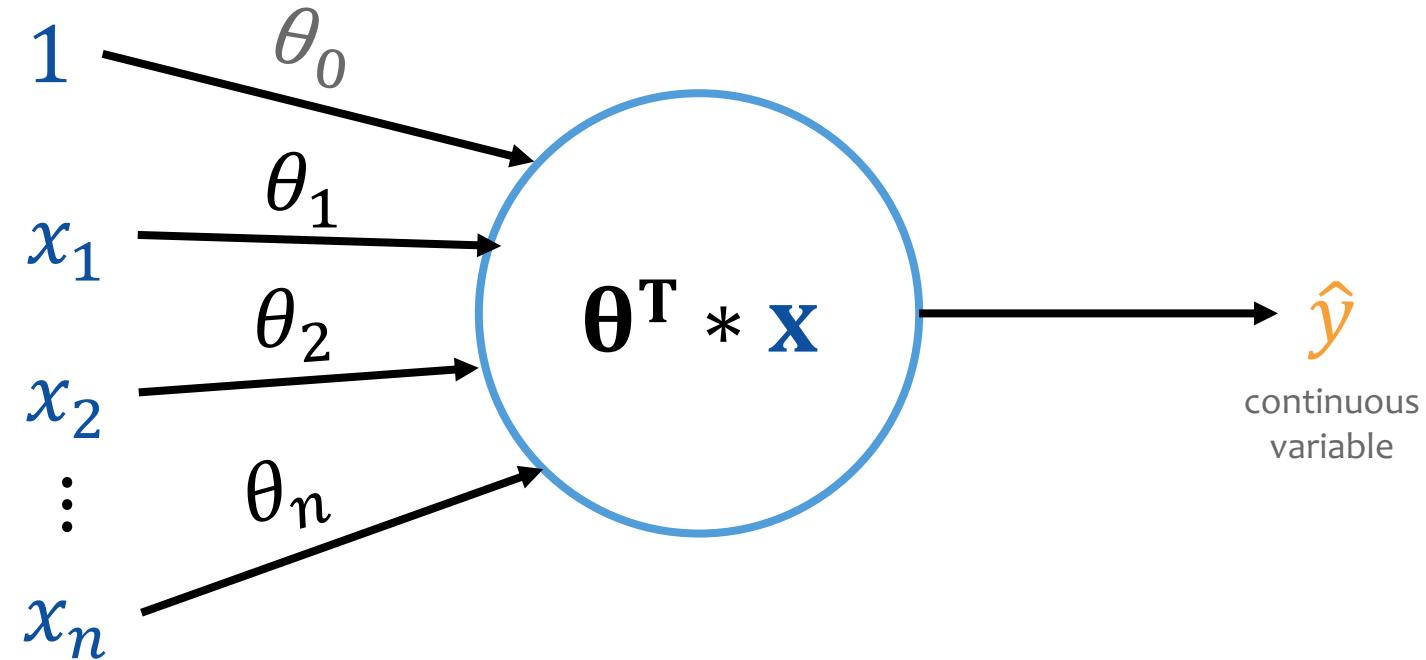
# An Alternative Visualization for Linear Regression

Let's first rename our **Linear Regression model**  $h_{\theta}(\mathbf{x})$  to  $f_{\theta}(\mathbf{x})$ .

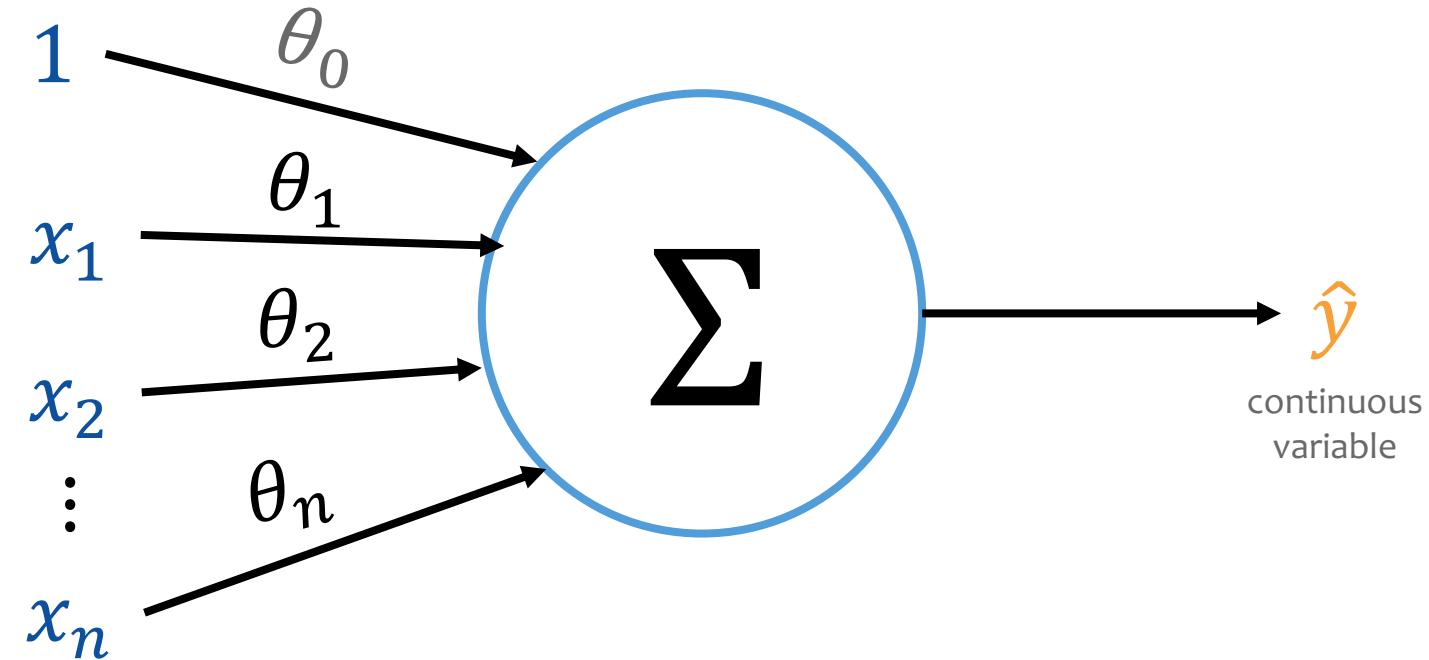
# An Alternative Visualization for Linear Regression



# An Alternative Visualization for Linear Regression

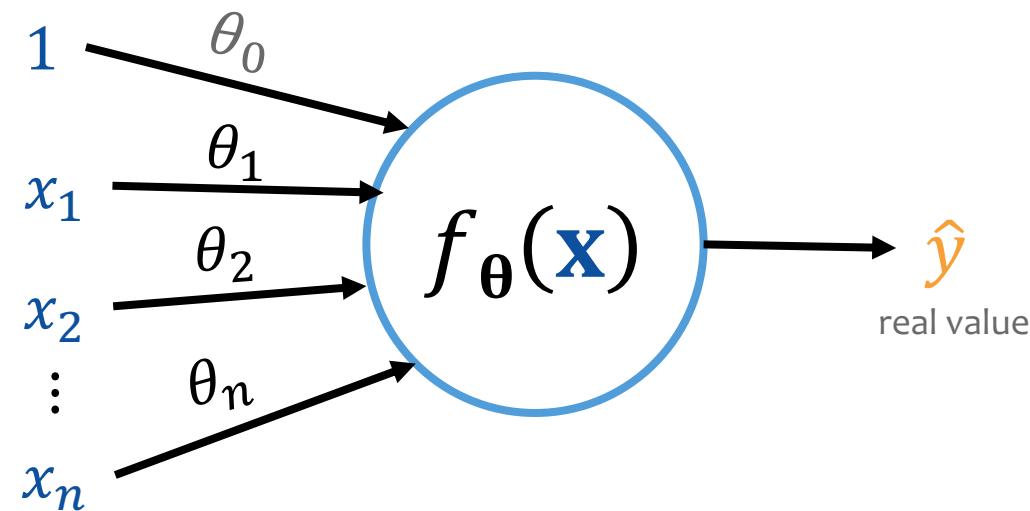


# An Alternative Visualization for Linear Regression



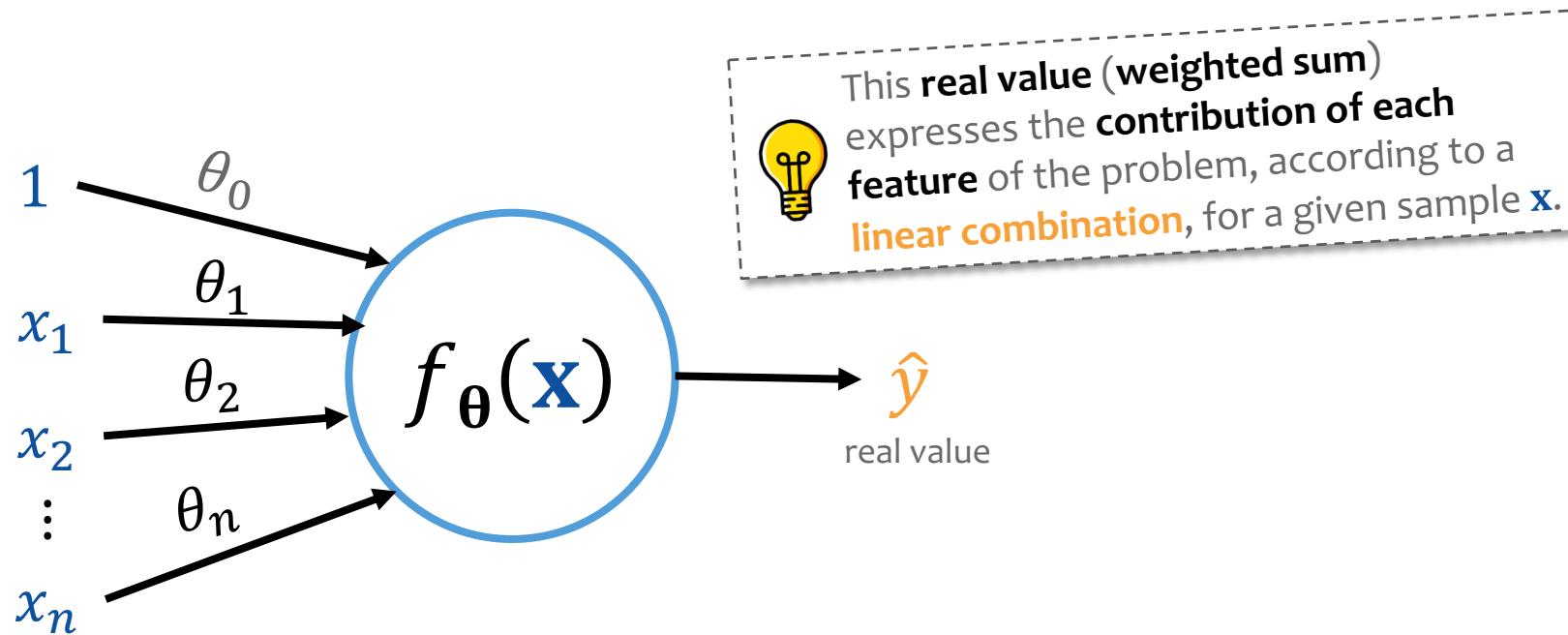
# From Linear Regression to Bin. Linear Classification

## Linear Regression



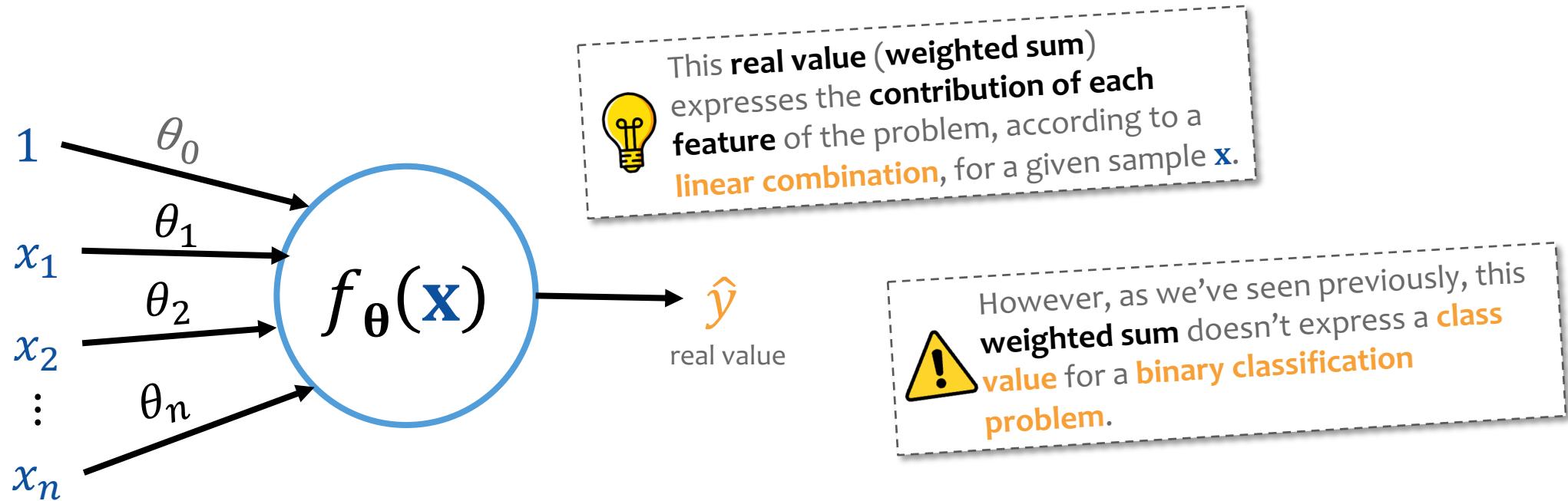
# From Linear Regression to Bin. Linear Classification

## Linear Regression



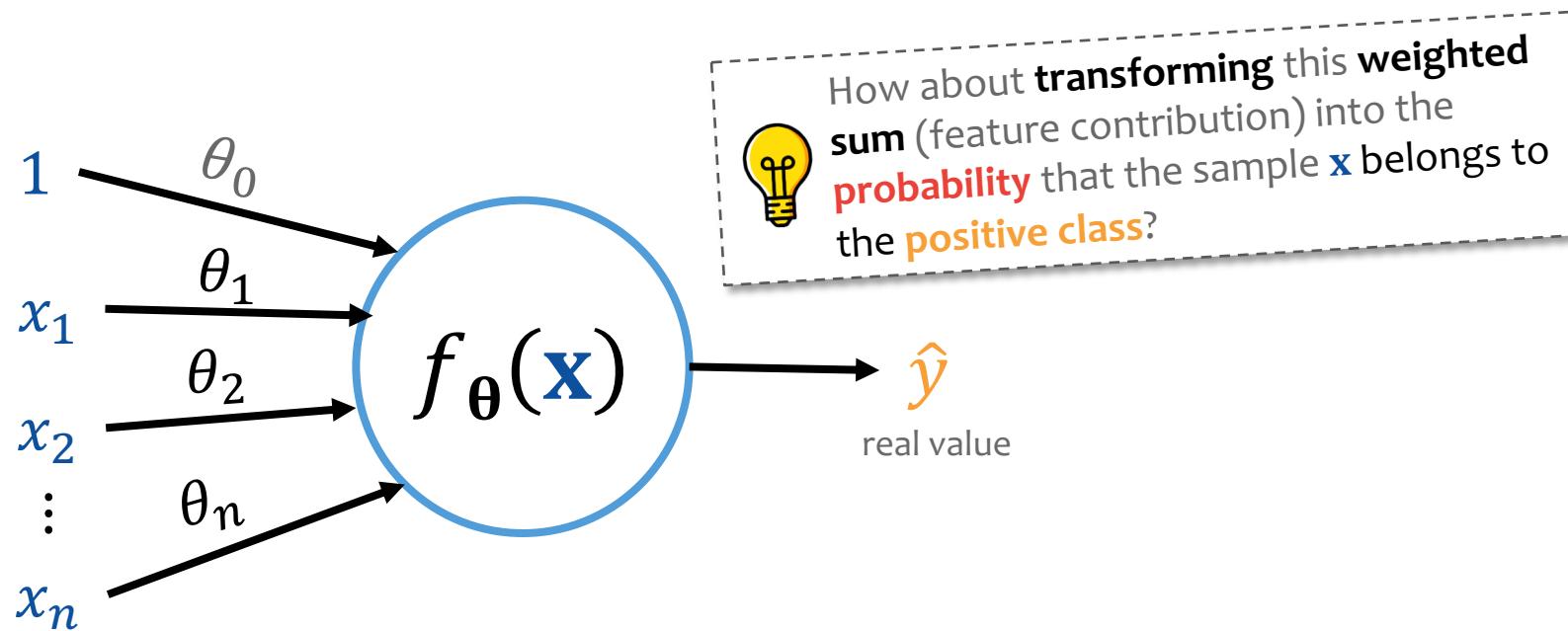
# From Linear Regression to Bin. Linear Classification

## Linear Regression



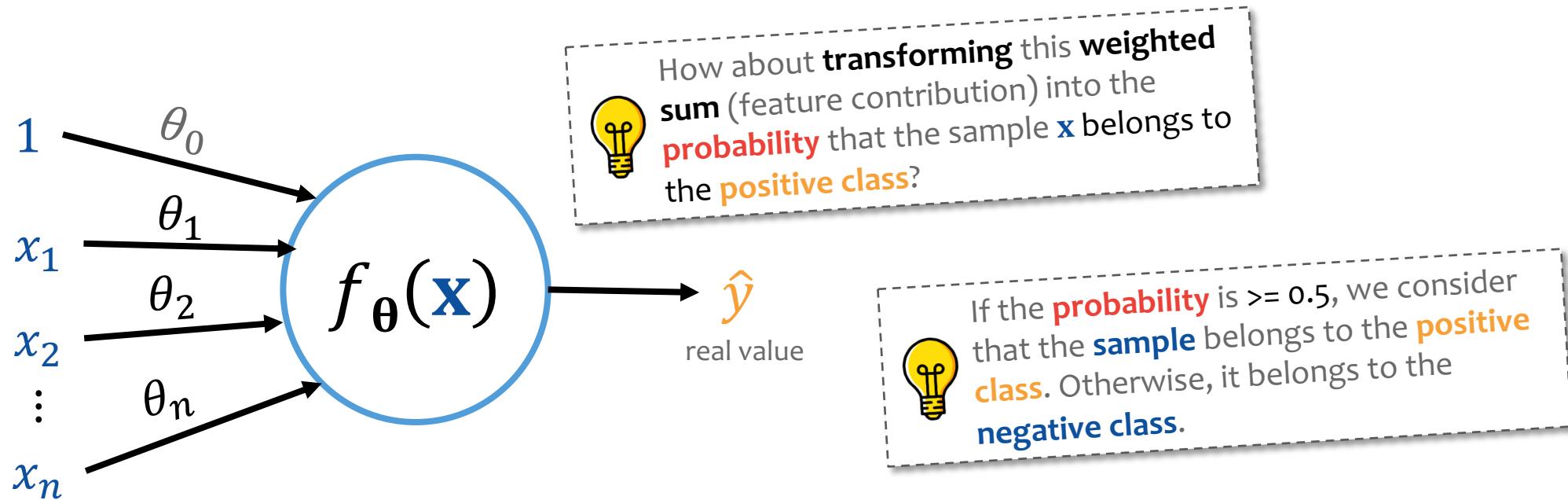
# From Linear Regression to Bin. Linear Classification

## Linear Regression



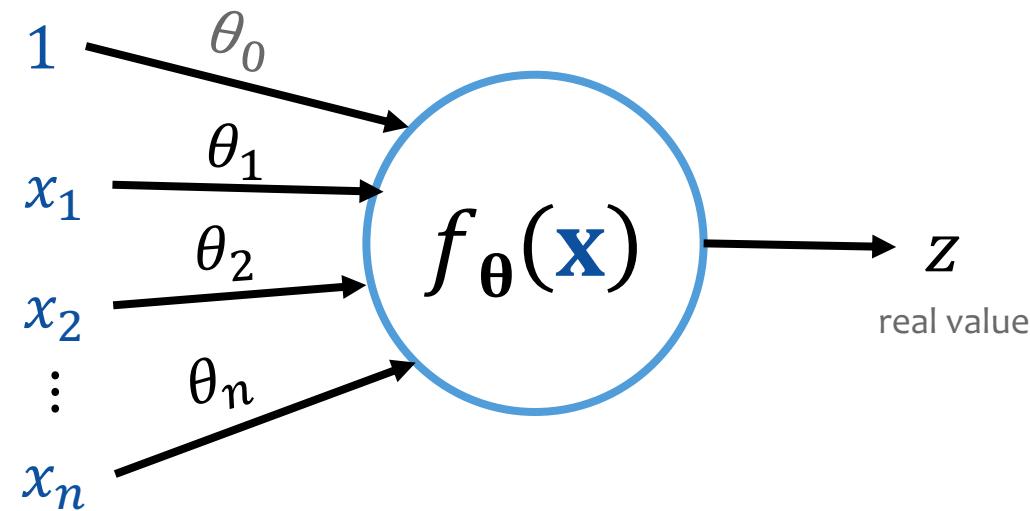
# From Linear Regression to Bin. Linear Classification

## Linear Regression



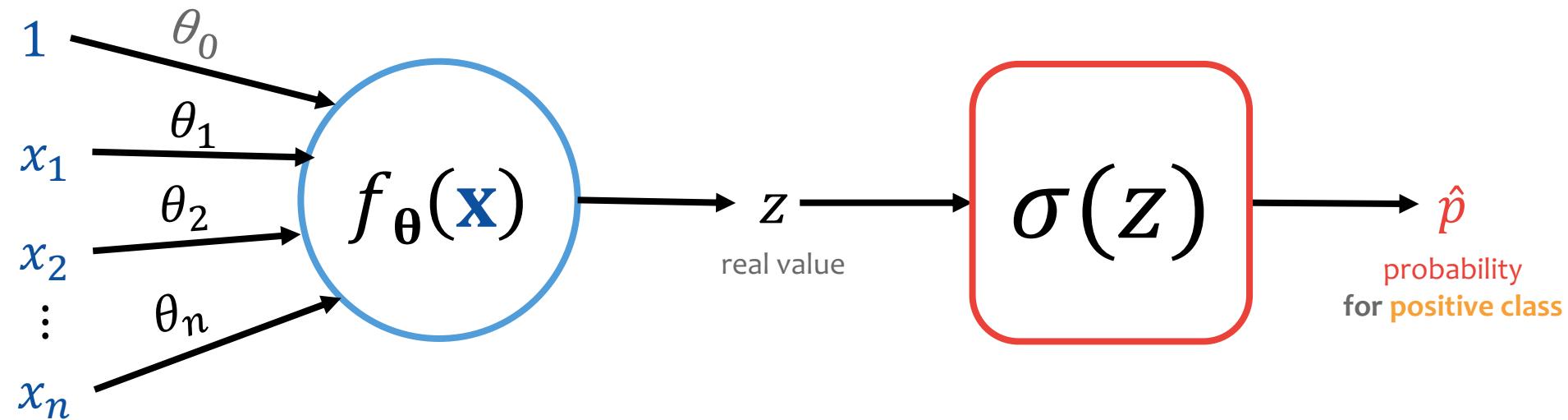
# From Linear Regression to Bin. Linear Classification

## Binary Linear Classifier



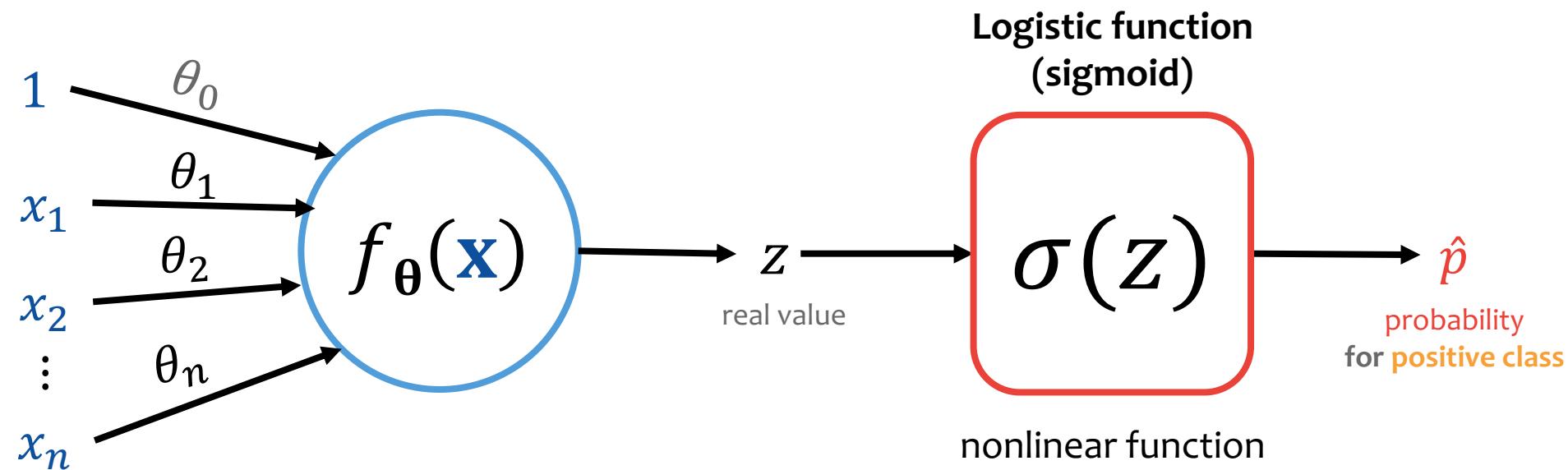
# From Linear Regression to Bin. Linear Classification

## Binary Linear Classifier



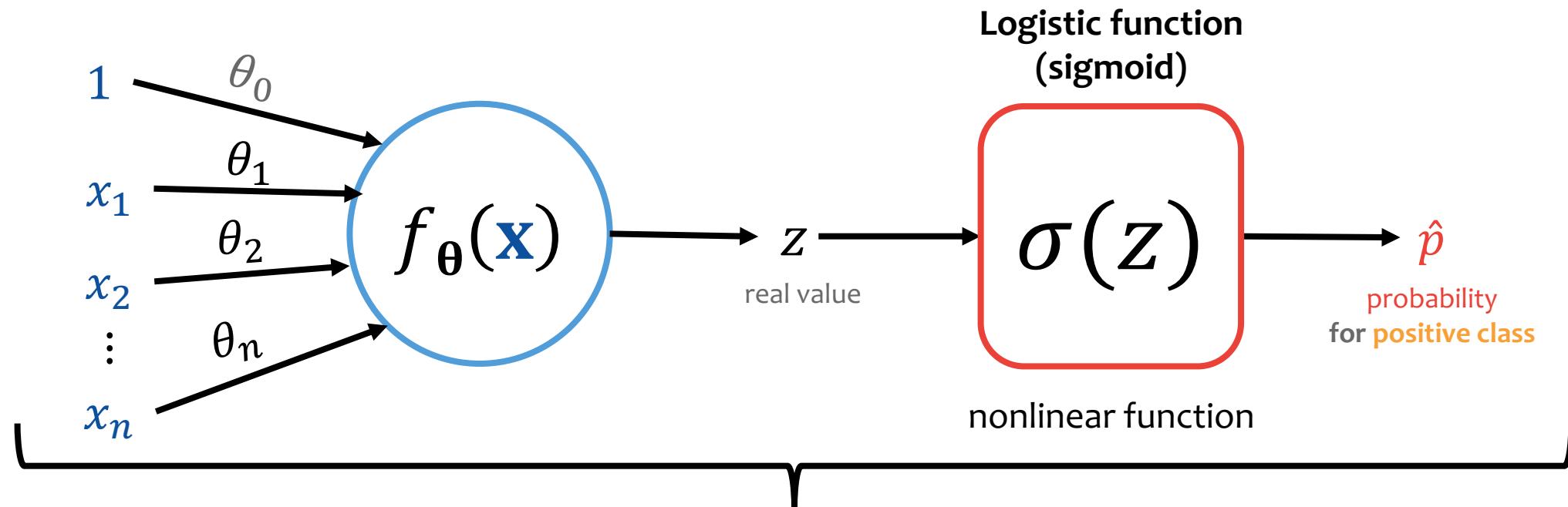
# From Linear Regression to Bin. Linear Classification

## Binary Linear Classifier



# From Linear Regression to Bin. Linear Classification

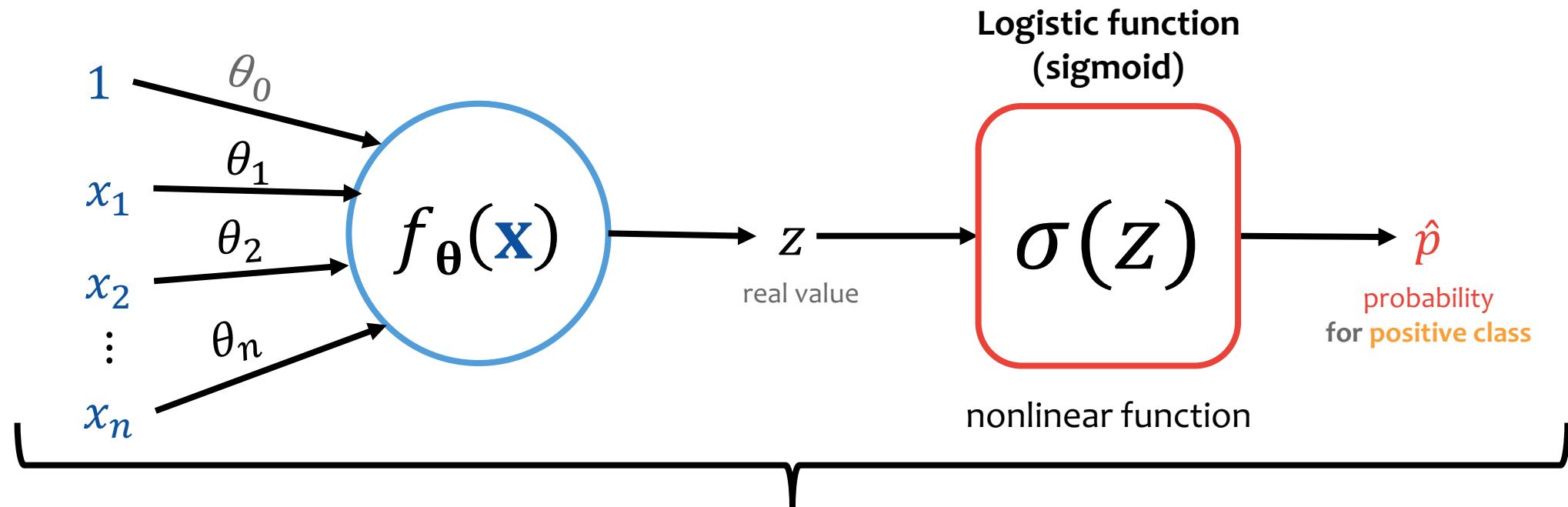
## Binary Linear Classifier



This is our **new model**  $h_{\theta}(\mathbf{x})$  (**binary linear classifier**)

# From Linear Regression to Bin. Linear Classification

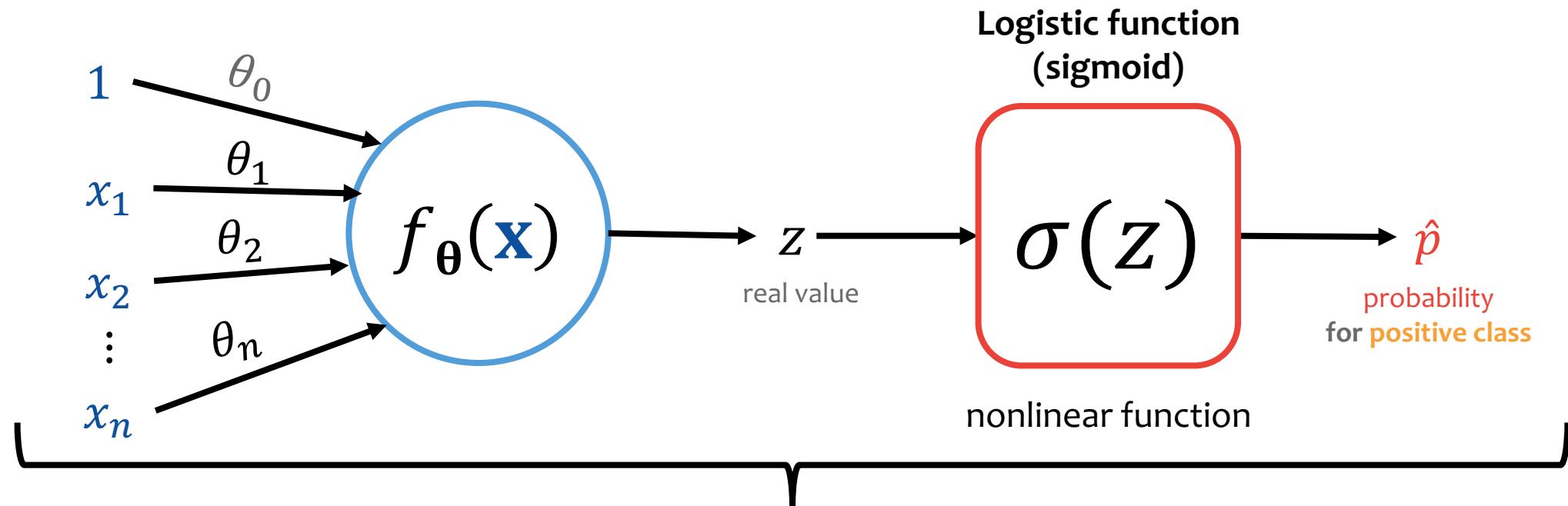
## Binary Linear Classifier



$$\hat{p} = h_{\theta}(\mathbf{x})$$

# From Linear Regression to Bin. Linear Classification

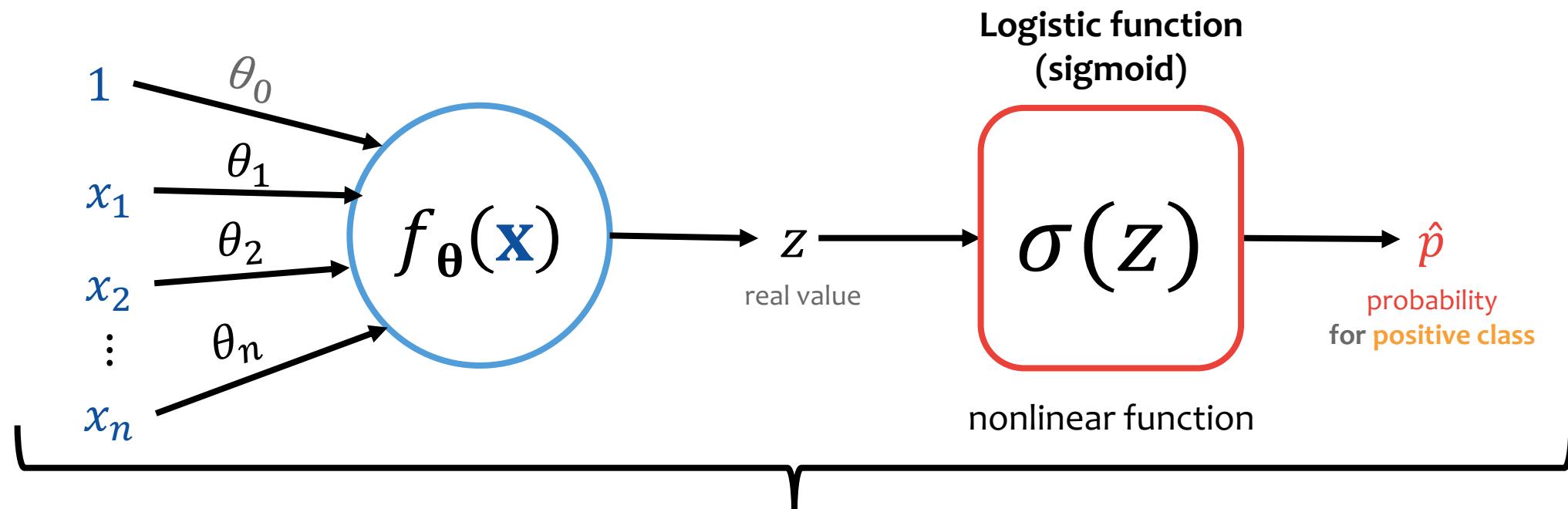
## Binary Linear Classifier



$$\hat{p} = h_{\theta}(\mathbf{x}) = P(\mathbf{y} = 1 | \mathbf{x}; \theta)$$

# From Linear Regression to Bin. Linear Classification

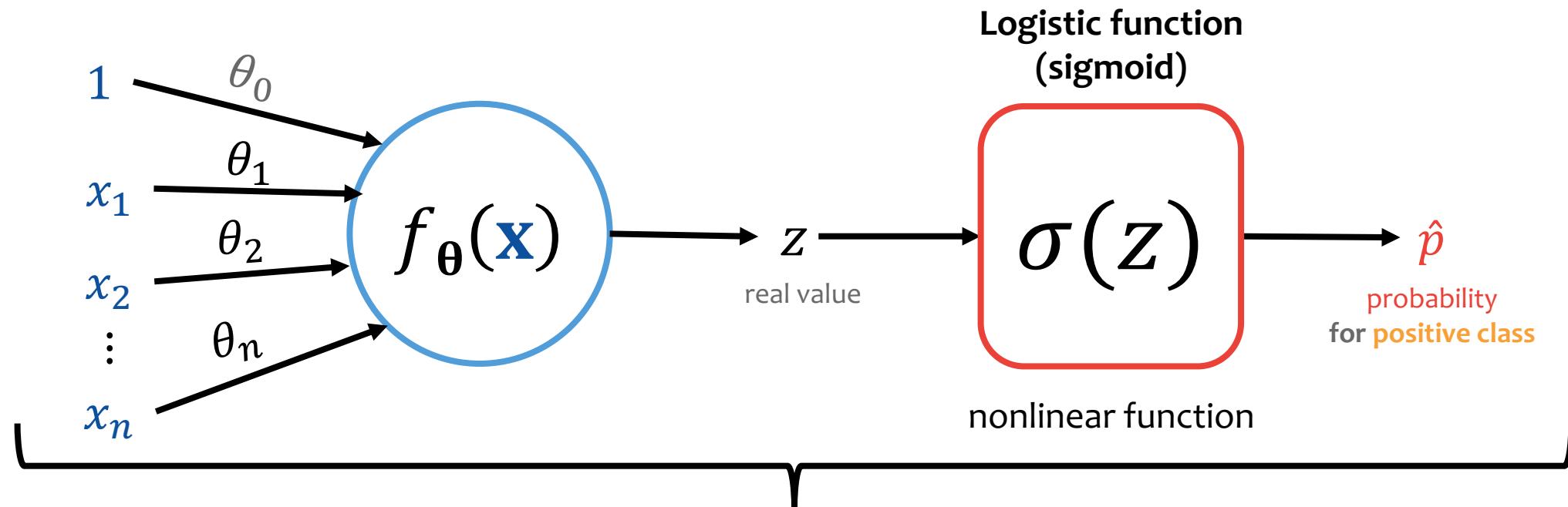
## Binary Linear Classifier



$$\hat{p} = h_{\theta}(\mathbf{x}) = P(\mathbf{y} = 1 | \mathbf{x}; \boldsymbol{\theta}) = \sigma(z)$$

# From Linear Regression to Bin. Linear Classification

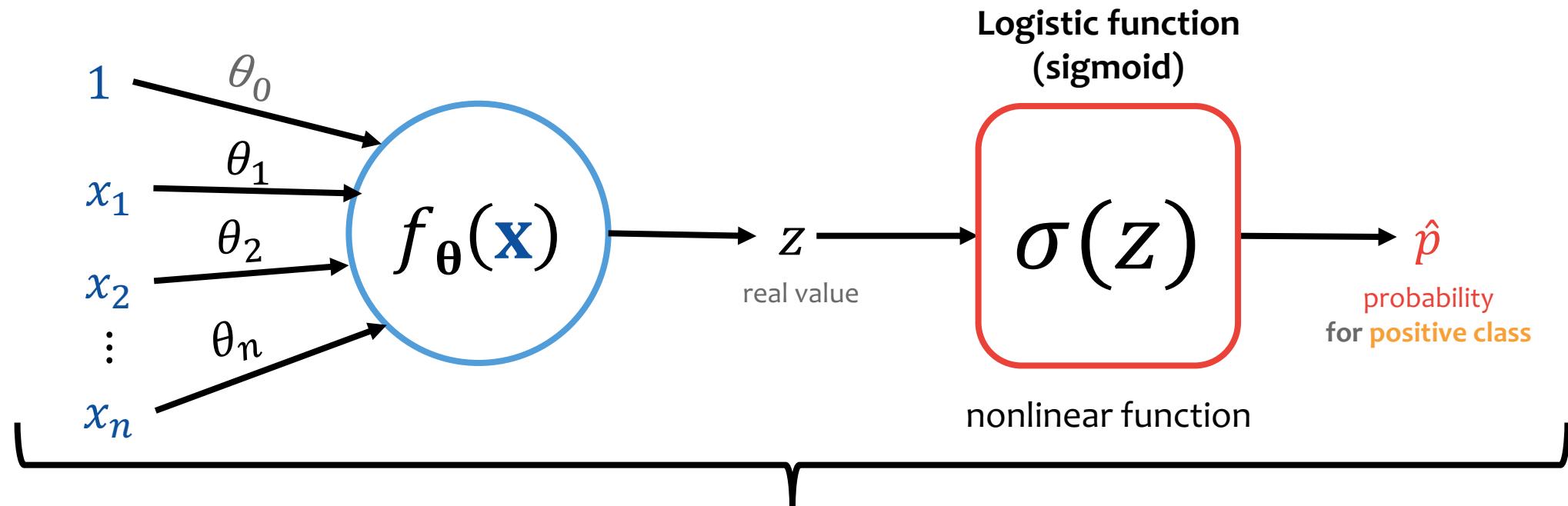
## Binary Linear Classifier



$$\hat{p} = h_{\theta}(\mathbf{x}) = P(\mathbf{y} = 1 | \mathbf{x}; \boldsymbol{\theta}) = \sigma(z) = \sigma(f_{\theta}(\mathbf{x}))$$

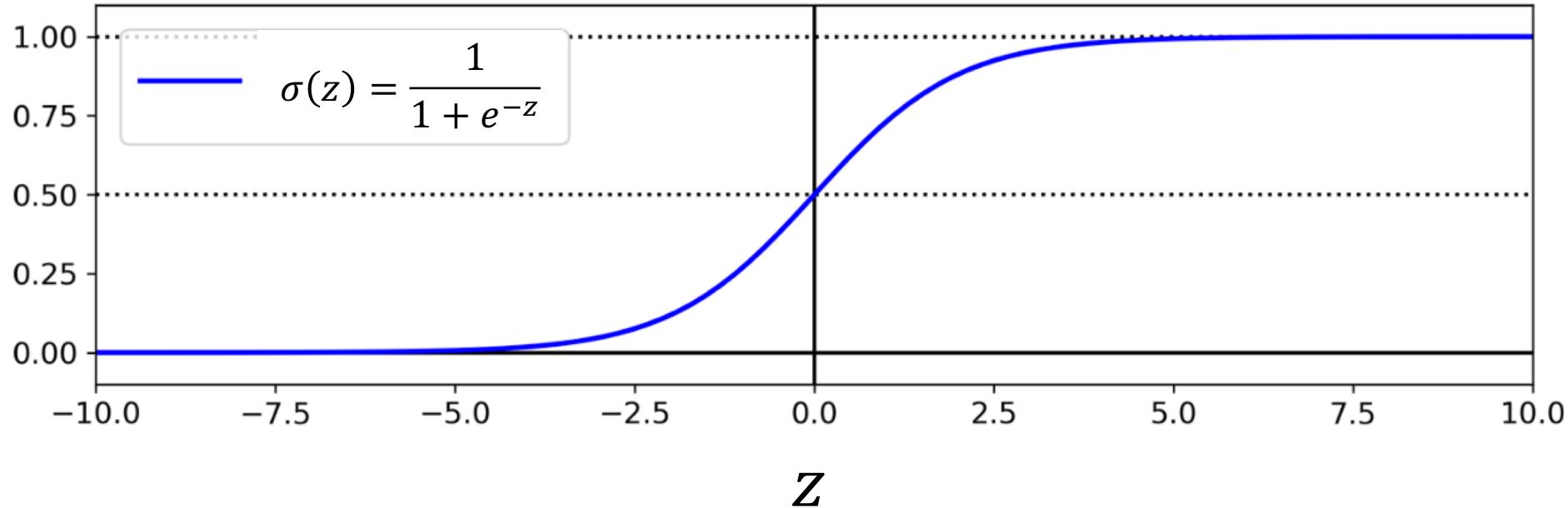
# From Linear Regression to Bin. Linear Classification

## Binary Linear Classifier



$$\hat{p} = h_{\theta}(\mathbf{x}) = P(\mathbf{y} = 1 | \mathbf{x}; \boldsymbol{\theta}) = \sigma(z) = \sigma(f_{\theta}(\mathbf{x})) = \sigma(\boldsymbol{\theta}^T * \mathbf{x})$$

# Logistic Function (Sigmoid)



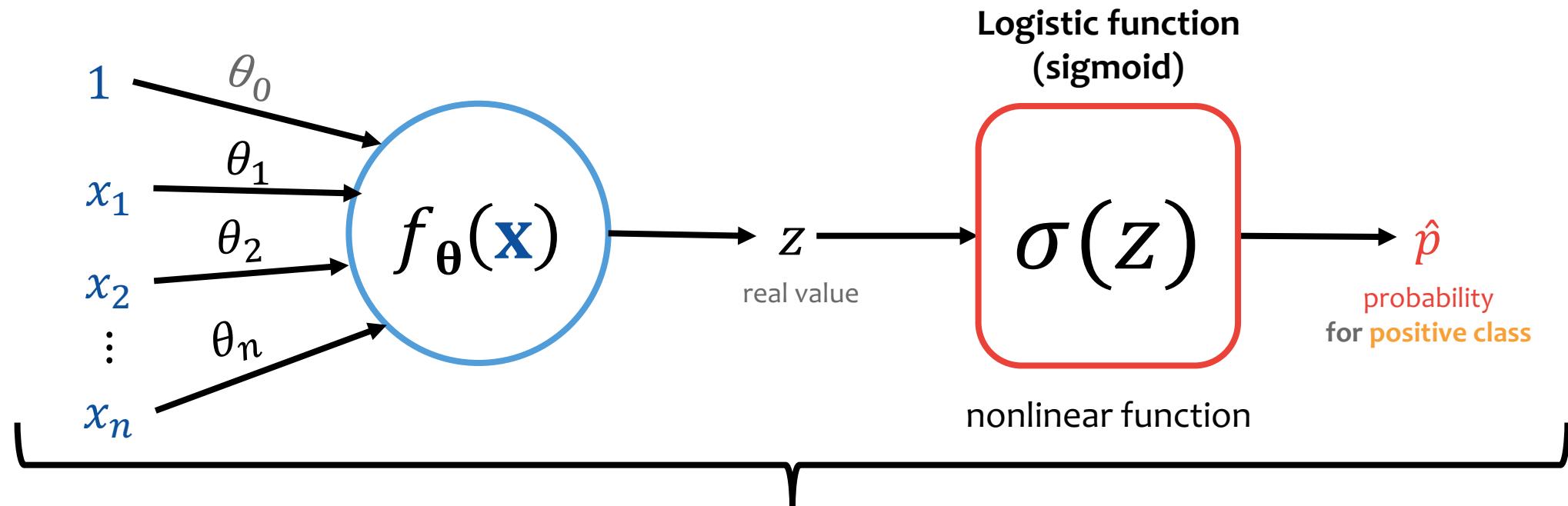
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

nonlinear function

$\sigma(z)$  returns a value between **0** and **1**  
 $\sigma(z) < 0.5$  if  $z < 0$   
 $\sigma(z) \geq 0.5$  if  $z \geq 0$

# From Linear Regression to Bin. Linear Classification

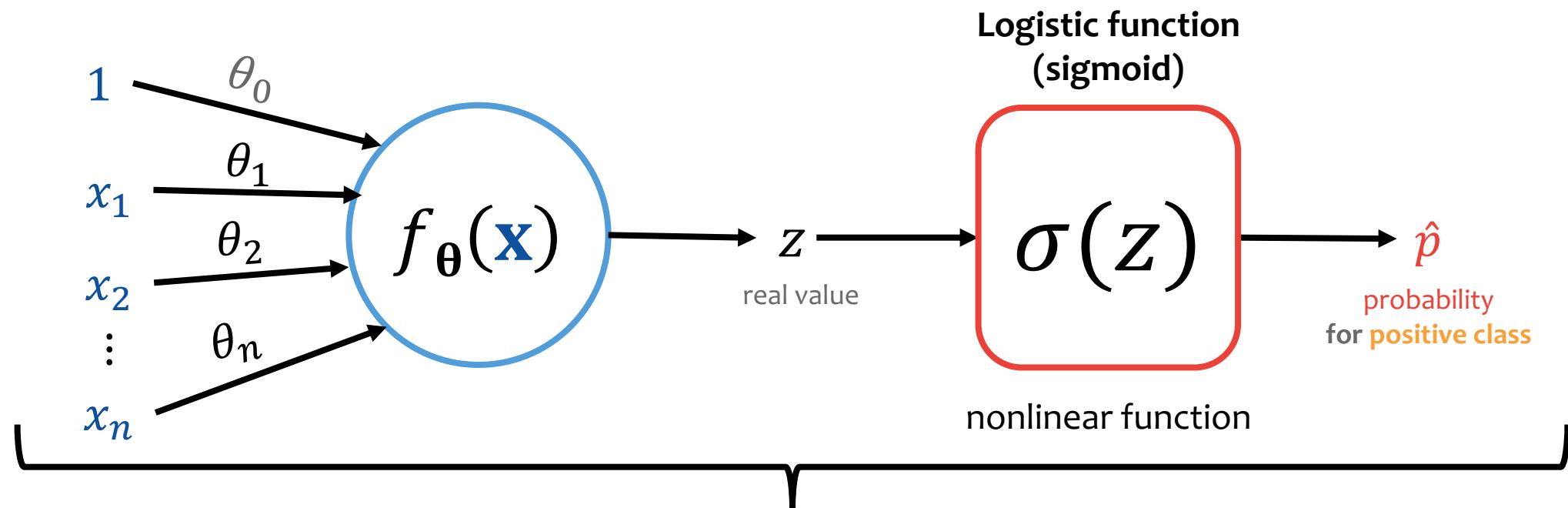
## Binary Linear Classifier



$$\hat{p} = h_{\theta}(\mathbf{x}) = P(\mathbf{y} = 1 | \mathbf{x}; \boldsymbol{\theta}) = \sigma(z) = \sigma(f_{\theta}(\mathbf{x})) = \sigma(\boldsymbol{\theta}^T * \mathbf{x})$$

# From Linear Regression to Bin. Linear Classification

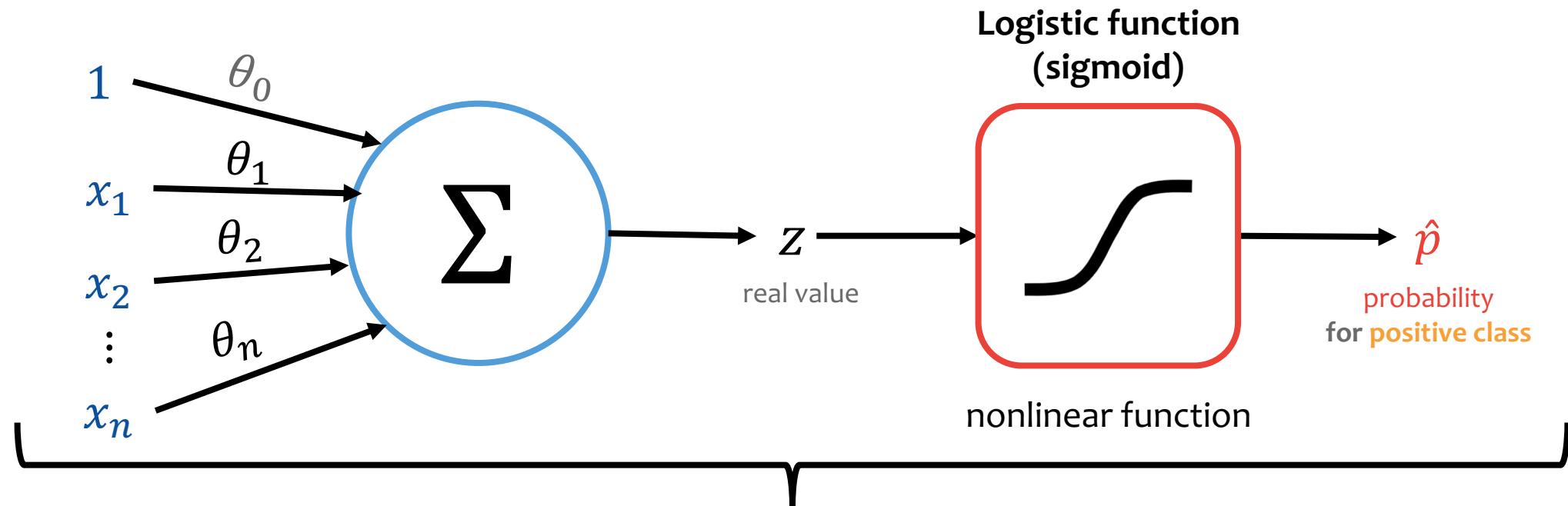
## Binary Linear Classifier



$$\hat{p} = h_{\theta}(\mathbf{x}) = P(\mathbf{y} = 1 | \mathbf{x}; \theta) = \sigma(z) = \sigma(f_{\theta}(\mathbf{x})) = \sigma(\boldsymbol{\theta}^T * \mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T * \mathbf{x}}}$$

# From Linear Regression to Bin. Linear Classification

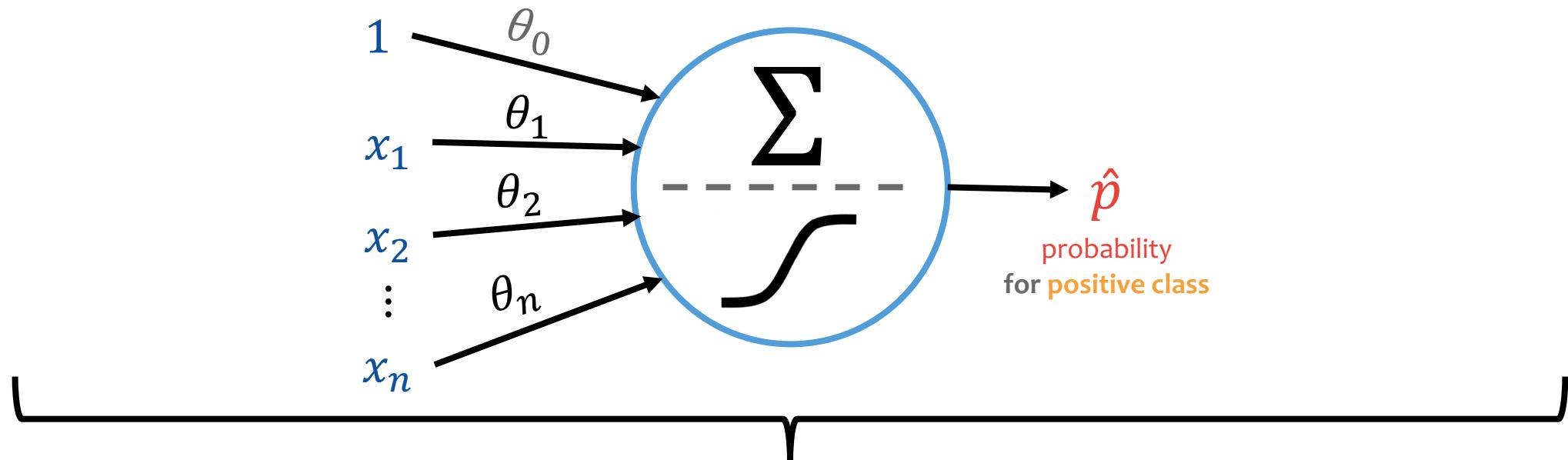
## Binary Linear Classifier



$$\hat{p} = h_{\theta}(\mathbf{x}) = P(\mathbf{y} = 1 | \mathbf{x}; \theta) = \sigma(z) = \sigma(f_{\theta}(\mathbf{x})) = \sigma(\boldsymbol{\theta}^T * \mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T * \mathbf{x}}}$$

# From Linear Regression to Bin. Linear Classification

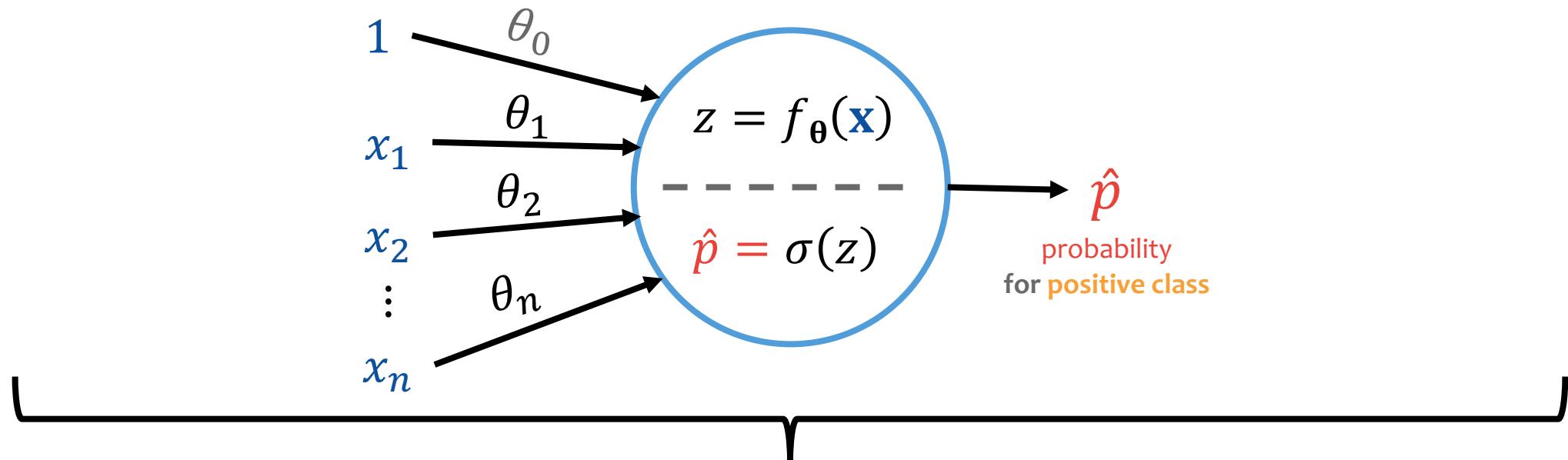
## Binary Linear Classifier



$$\hat{p} = h_{\theta}(\mathbf{x}) = P(\mathbf{y} = 1 | \mathbf{x}; \boldsymbol{\theta}) = \sigma(z) = \sigma(f_{\boldsymbol{\theta}}(\mathbf{x})) = \sigma(\boldsymbol{\theta}^T * \mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T * \mathbf{x}}}$$

# From Linear Regression to Bin. Linear Classification

## Binary Linear Classifier



$$\hat{p} = h_{\theta}(\mathbf{x}) = P(\mathbf{y} = 1 | \mathbf{x}; \theta) = \sigma(z) = \sigma(f_{\theta}(\mathbf{x})) = \sigma(\boldsymbol{\theta}^T * \mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T * \mathbf{x}}}$$

# From Linear Regression to Bin. Linear Classification

## Binary Linear Classifier

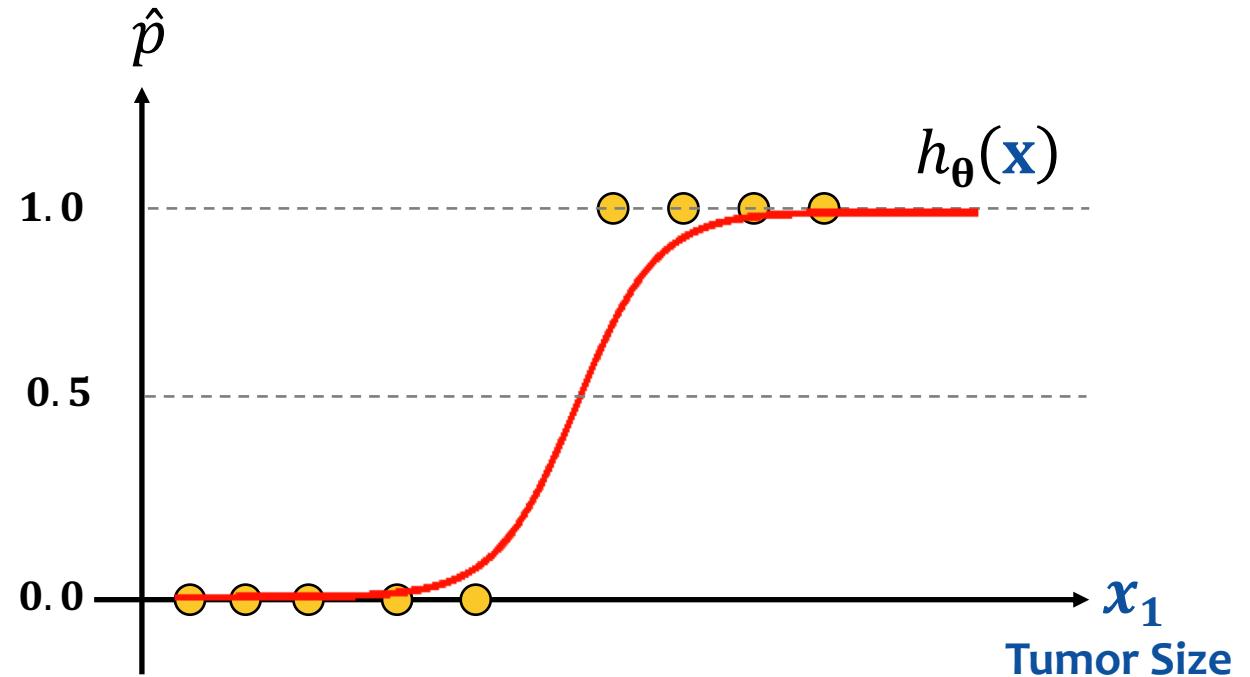
This model (binary linear classifier) is called **Logistic Regression**.

$$\hat{p} = h_{\theta}(\mathbf{x}) = P(\mathbf{y} = 1 | \mathbf{x}; \theta) = \sigma(z) = \sigma(f_{\theta}(\mathbf{x})) = \sigma(\boldsymbol{\theta}^T * \mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T * \mathbf{x}}}$$

## Logistic Regression Model

$$\hat{p} = h_{\theta}(\mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}}$$

Malignant Tumor  
Probability



## Logistic Regression Model

$$\hat{p} = h_{\theta}(\mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}}$$

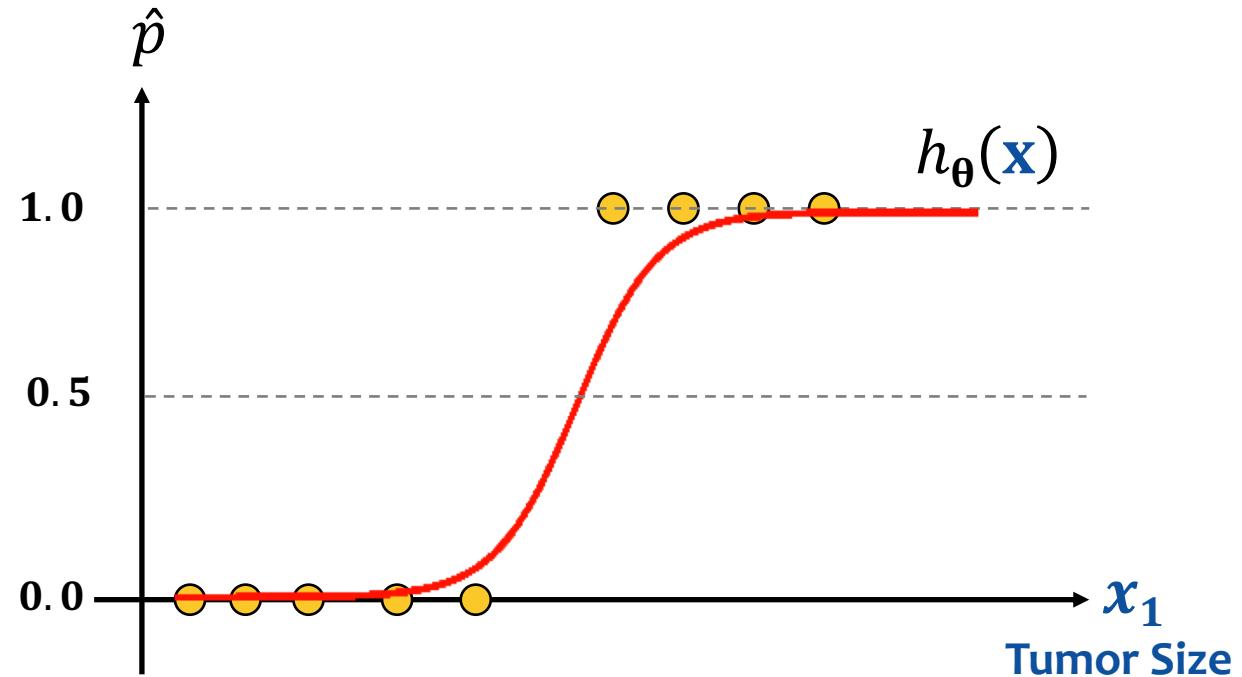
E.g:

$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumor size} \end{bmatrix}$$

$$\hat{p} = h_{\theta}(\mathbf{x}) = 0.7$$

The **probability** of the tumor being **malignant (positive class)** is **70%**.

Malignant Tumor Probability



## Logistic Regression Model

$$\hat{p} = h_{\theta}(\mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}}$$

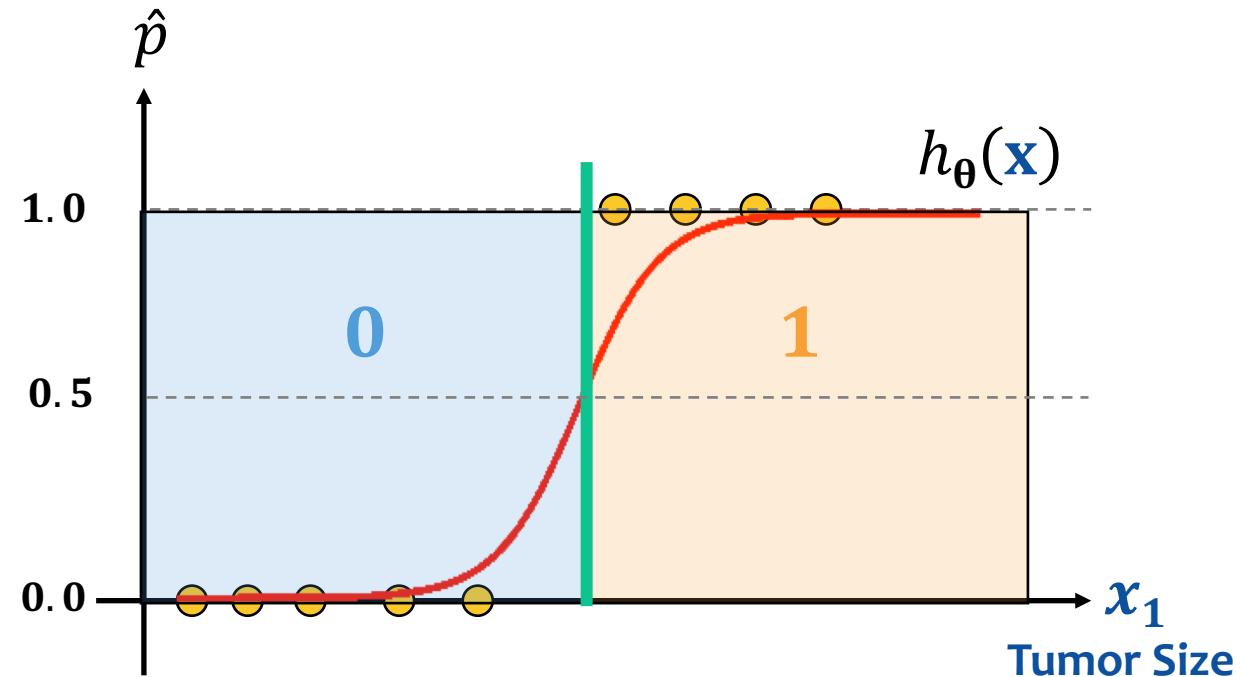
E.g:

$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumor size} \end{bmatrix}$$

$$\hat{p} = h_{\theta}(\mathbf{x}) = 0.7$$

The **probability** of the tumor being **malignant (positive class)** is **70%**.

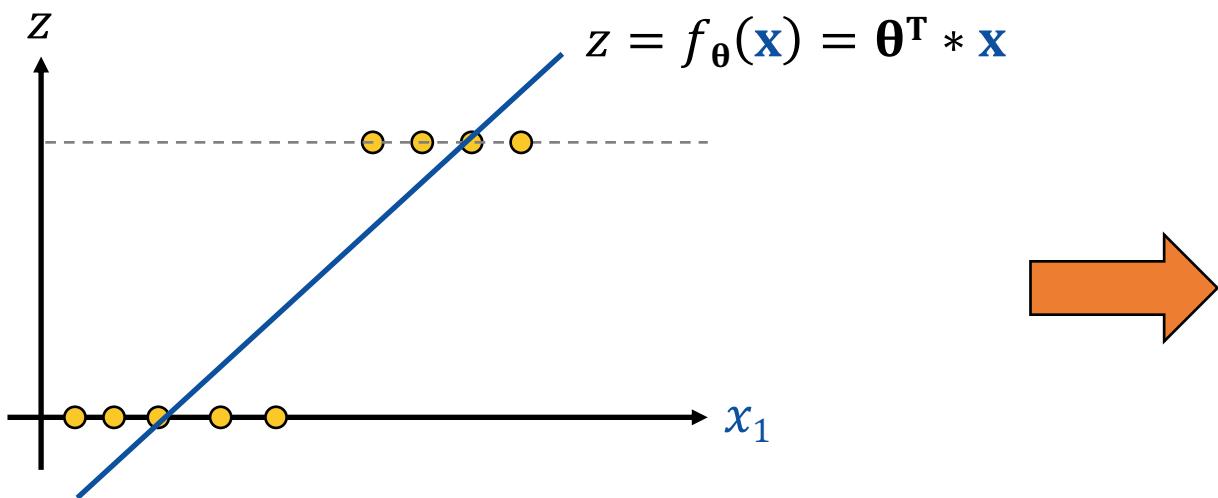
## Malignant Tumor Probability



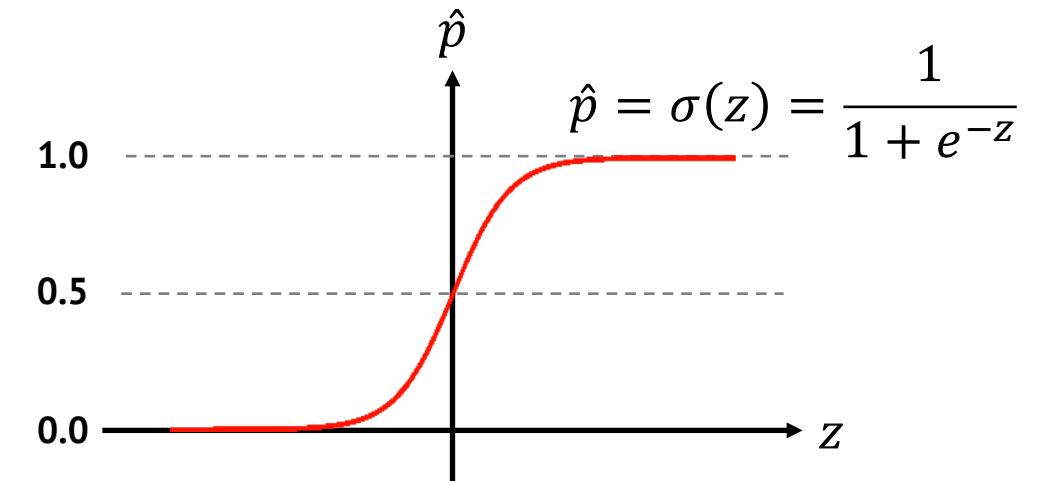
## Classification

$$\hat{y} = \begin{cases} 1 & \text{if } \hat{p} \geq 0.5 \\ 0 & \text{if } \hat{p} < 0.5 \end{cases}$$

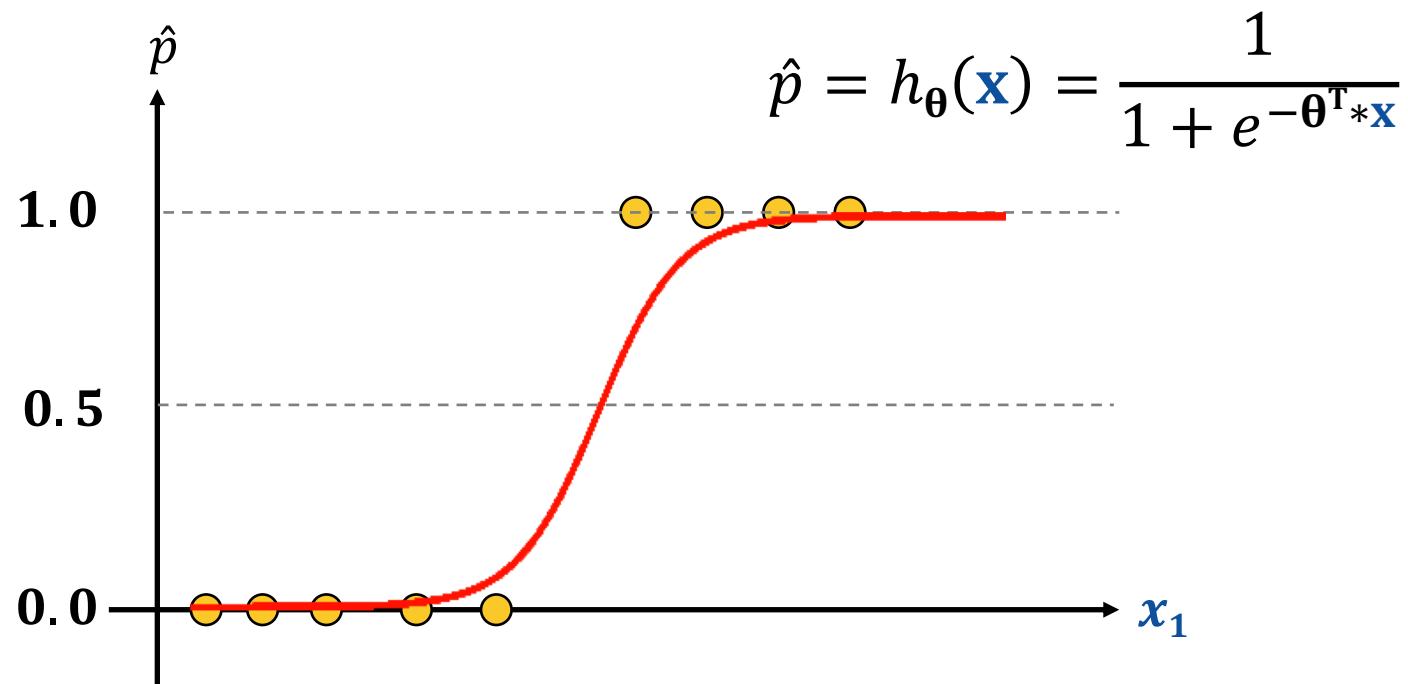
## Linear Model



## nonlinear Logistic Function (Sigmoid)



## Logistic Regression Model



# Training

**Goal:**  $\min_{\theta} J(\theta)$

# Training

objective cost function

Goal:  $\min_{\theta} J(\theta)$

# Training

Goal:  $\min_{\theta} J(\theta)$

↑  
objective cost function

Find the **vector  $\theta$**  (weights and bias) that **minimizes** a given **cost function  $J$**

# Training

Goal:  $\min_{\theta} J(\theta)$

↑  
objective cost function



Find the **vector  $\theta$**  (weights and bias) that **minimizes** a given **cost function  $J$**

**high probabilities** for **positive instances ( $y = 1$ )**;

**low probabilities** for **negative instances ( $y = 0$ )**;

# Training

Goal:  $\min_{\theta} J(\theta)$

↑  
objective cost function

Find the **vector  $\theta$**  (weights and bias) that **minimizes** a given **cost function  $J$**

→ **Gradient Descent!**

**high probabilities** for **positive instances ( $y = 1$ )**;

**low probabilities** for **negative instances ( $y = 0$ )**;

# Cost Function

The common chosen loss function for **probabilities** is the **log loss cost function**, defined as:

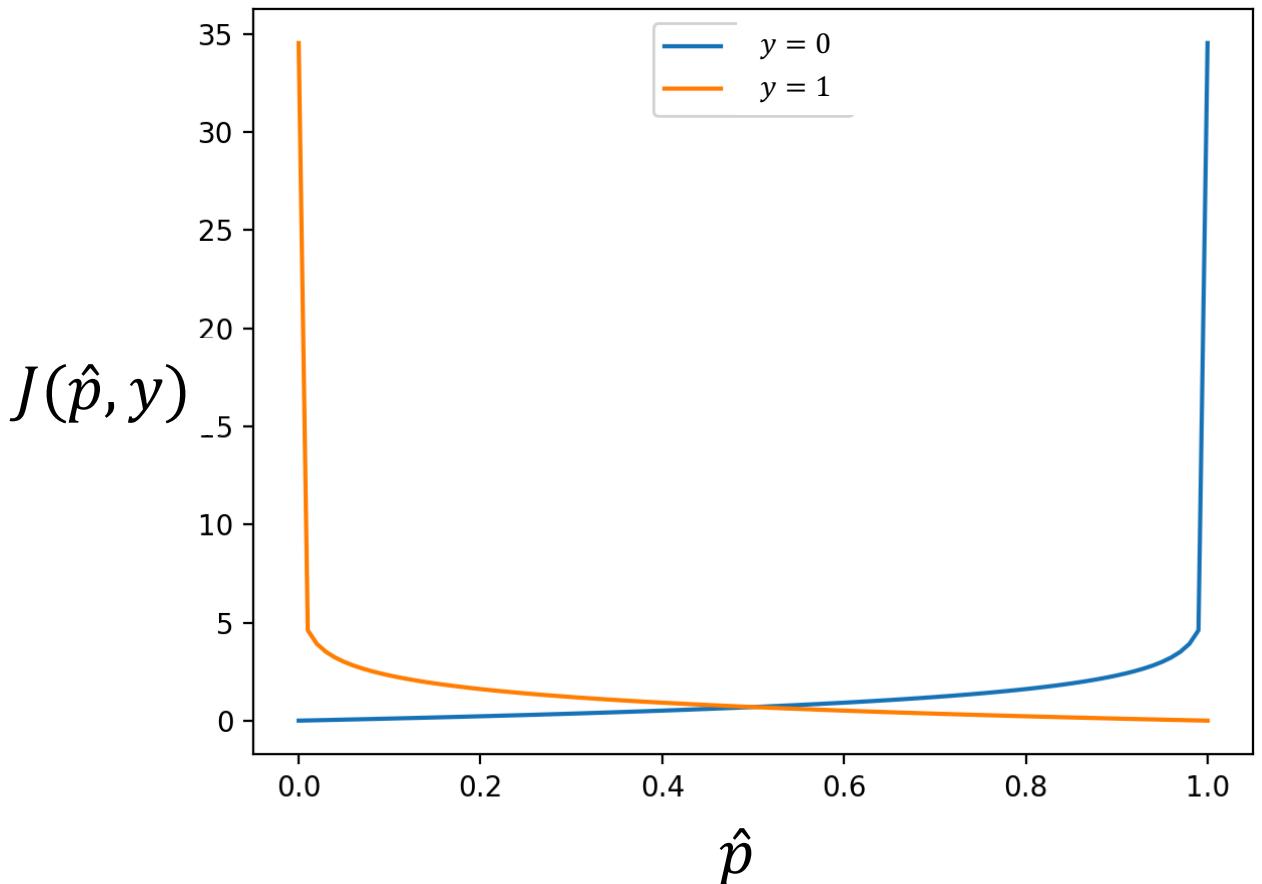
# Cost Function

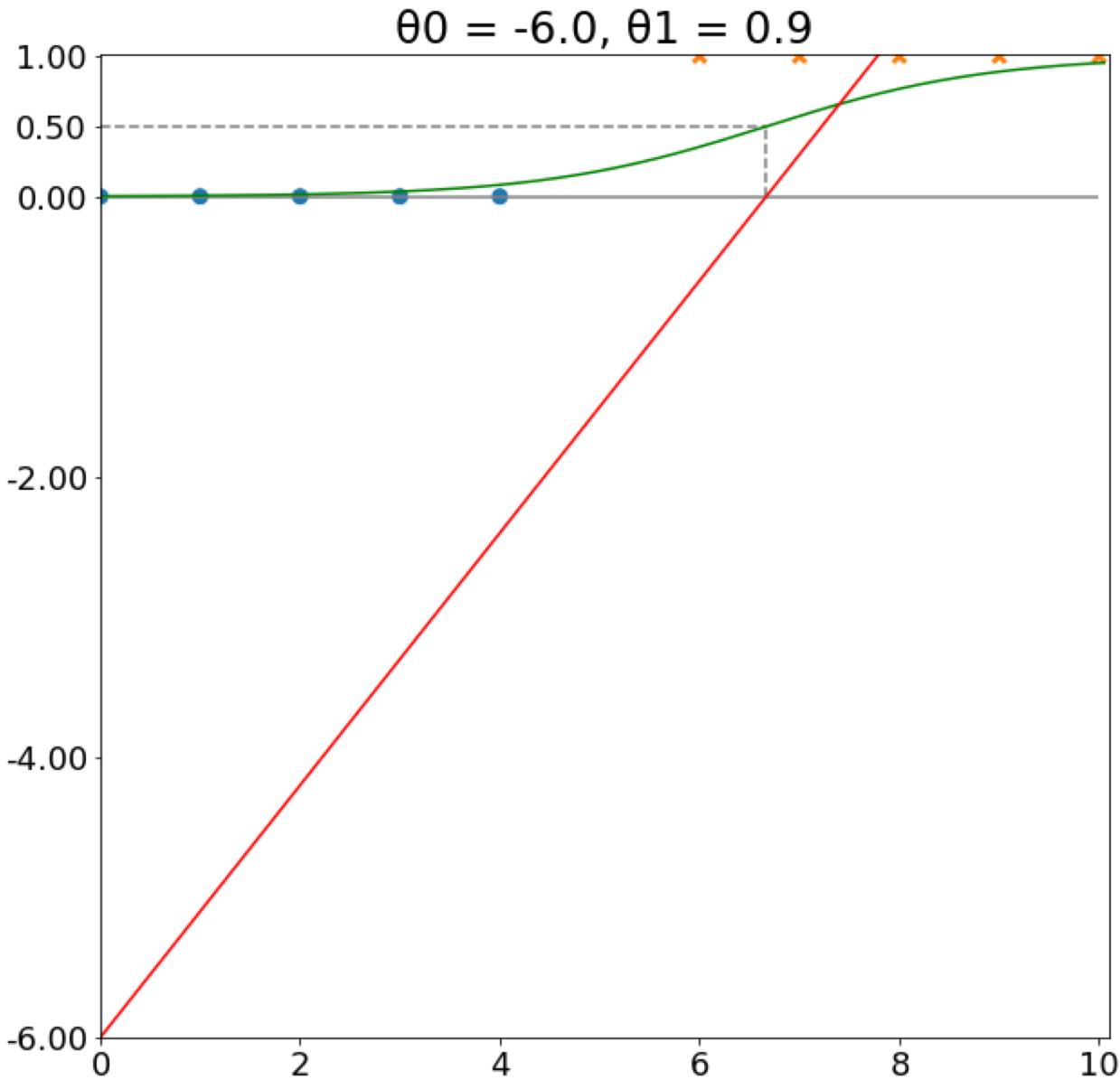
The common chosen loss function for **probabilities** is the **log loss cost function**, defined as:

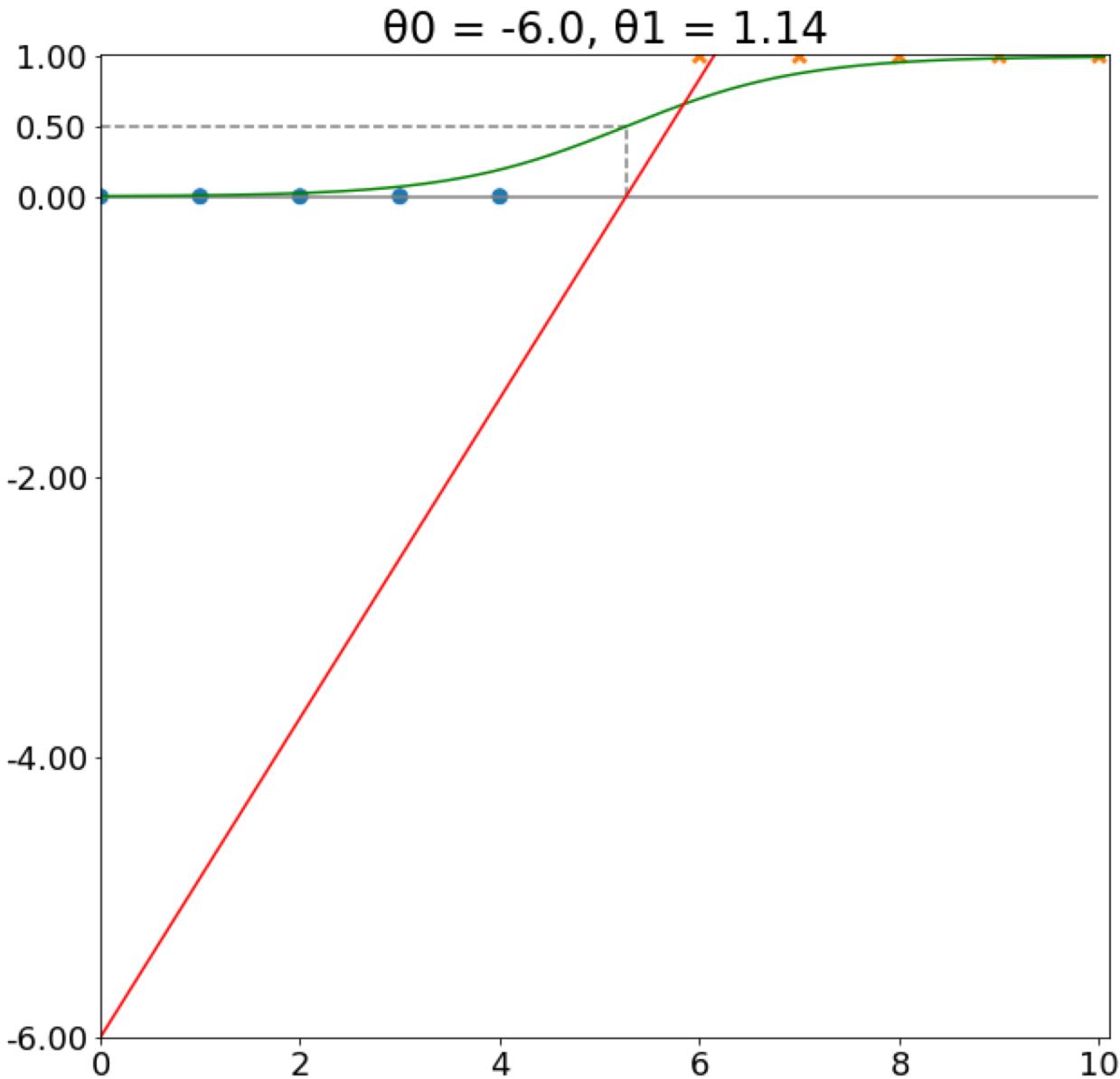
$$J(\hat{p}, y) = \begin{cases} -\log(\hat{p}) & \text{if } y = 1 \\ -\log(1 - \hat{p}) & \text{if } y = 0 \end{cases}$$

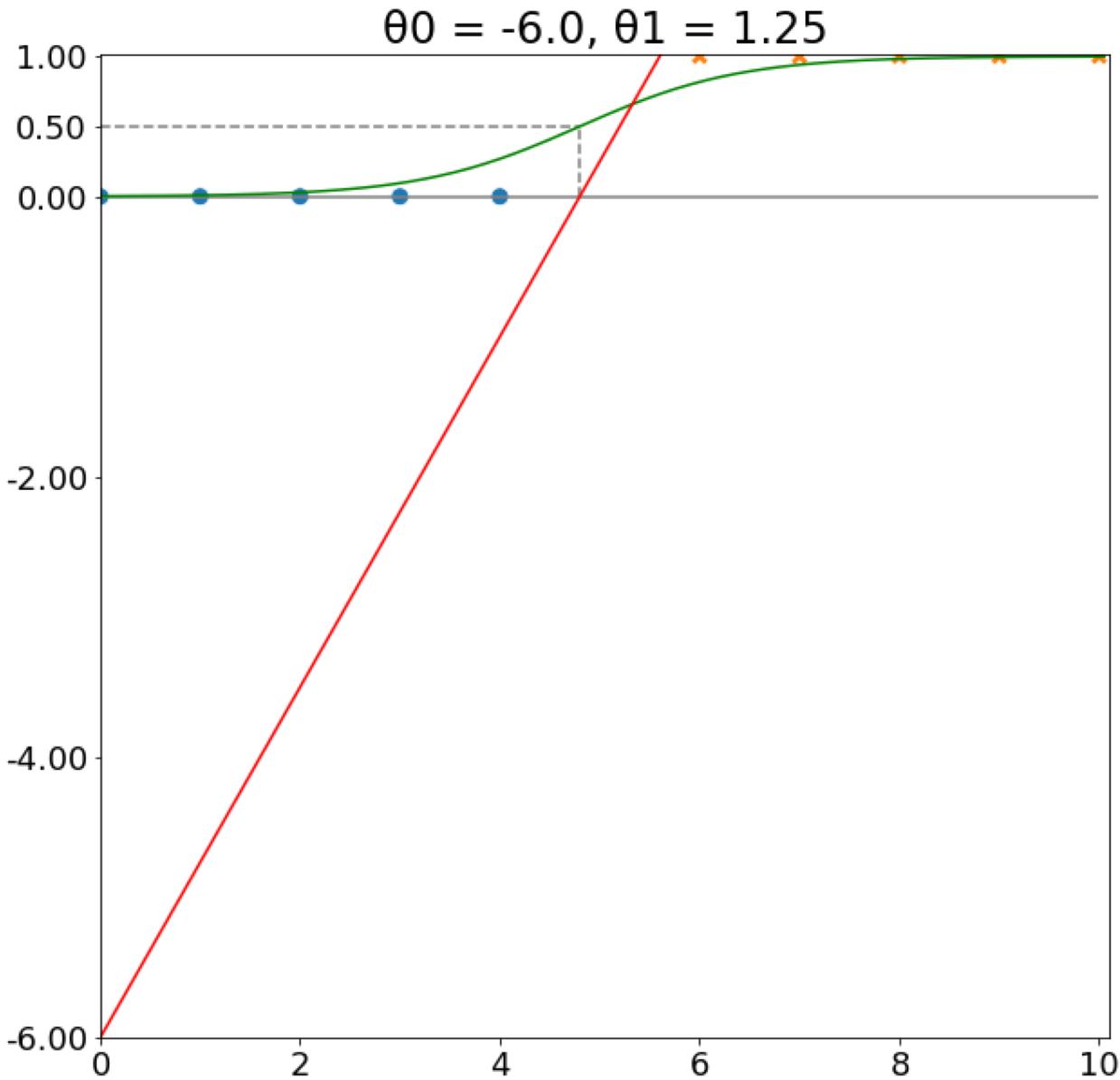
ground-truth (true label)

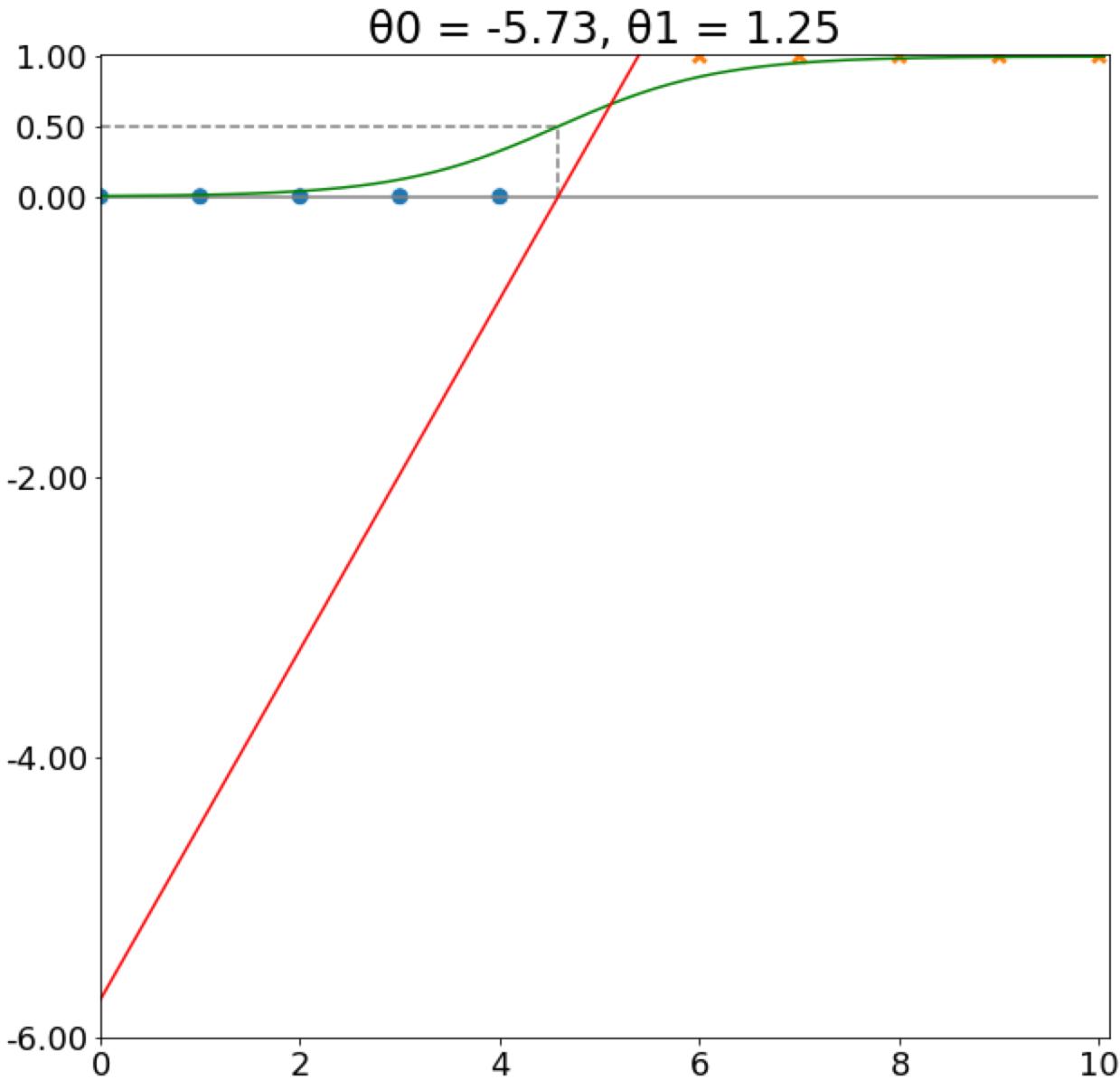
*y* = 1

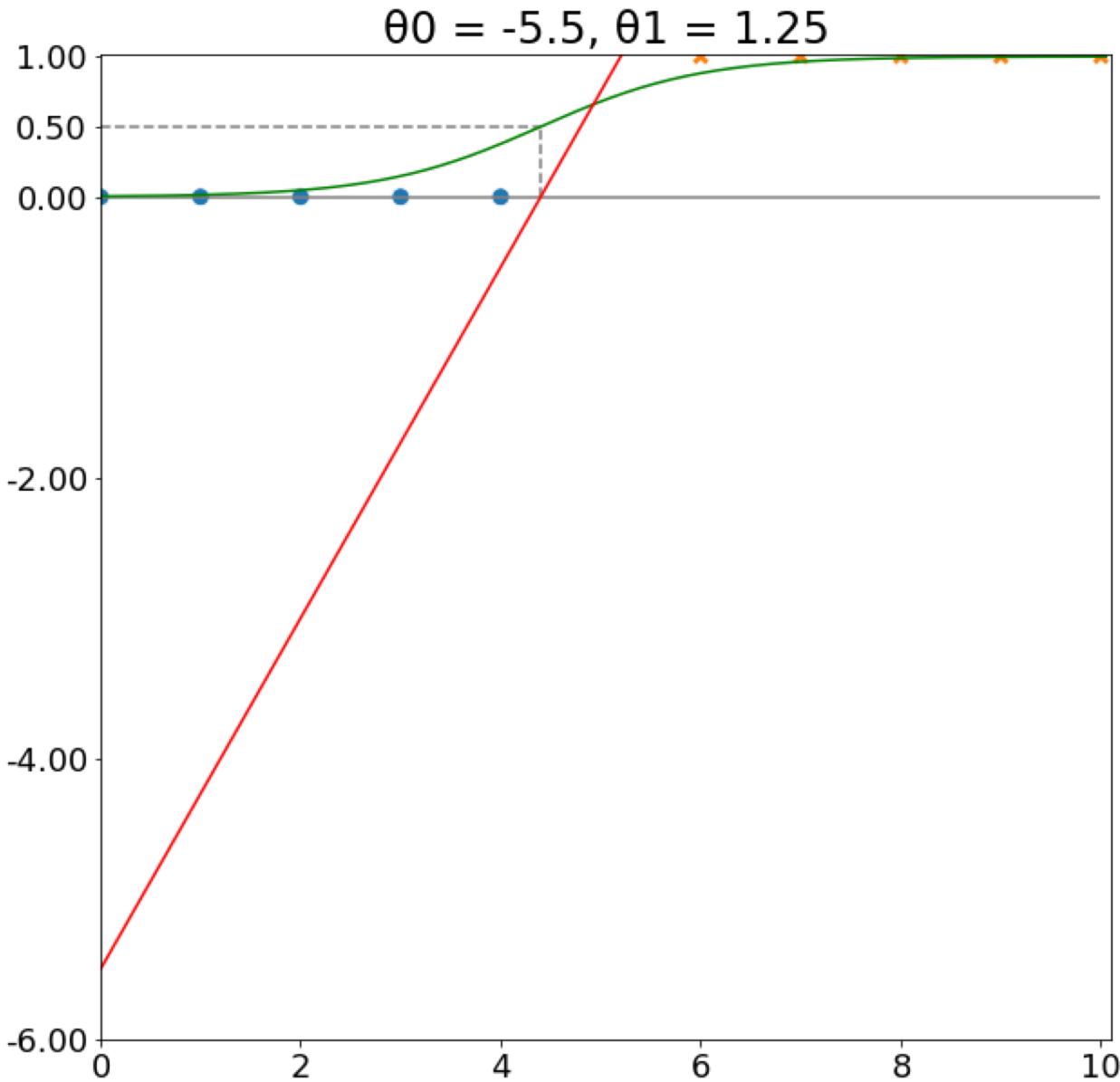




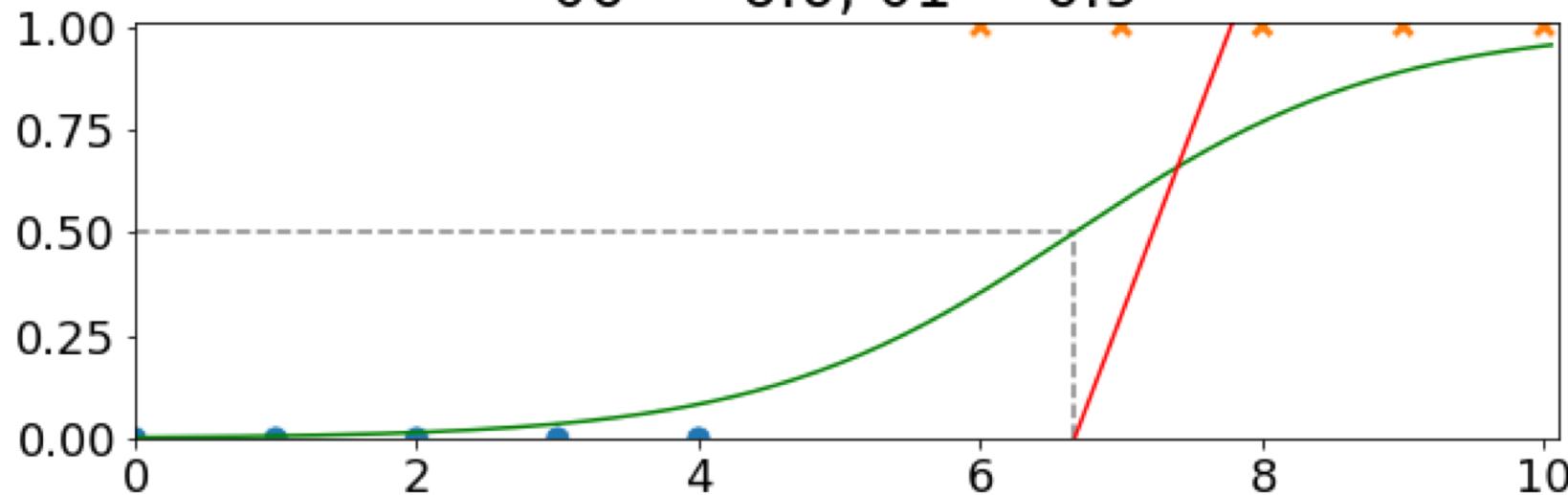




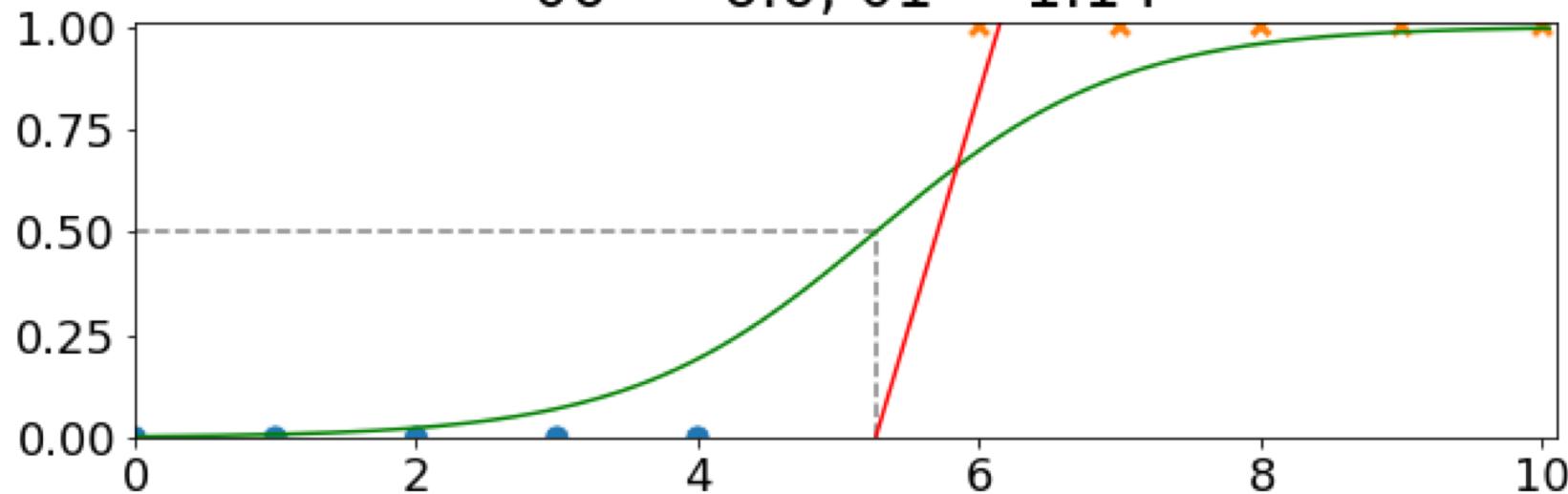




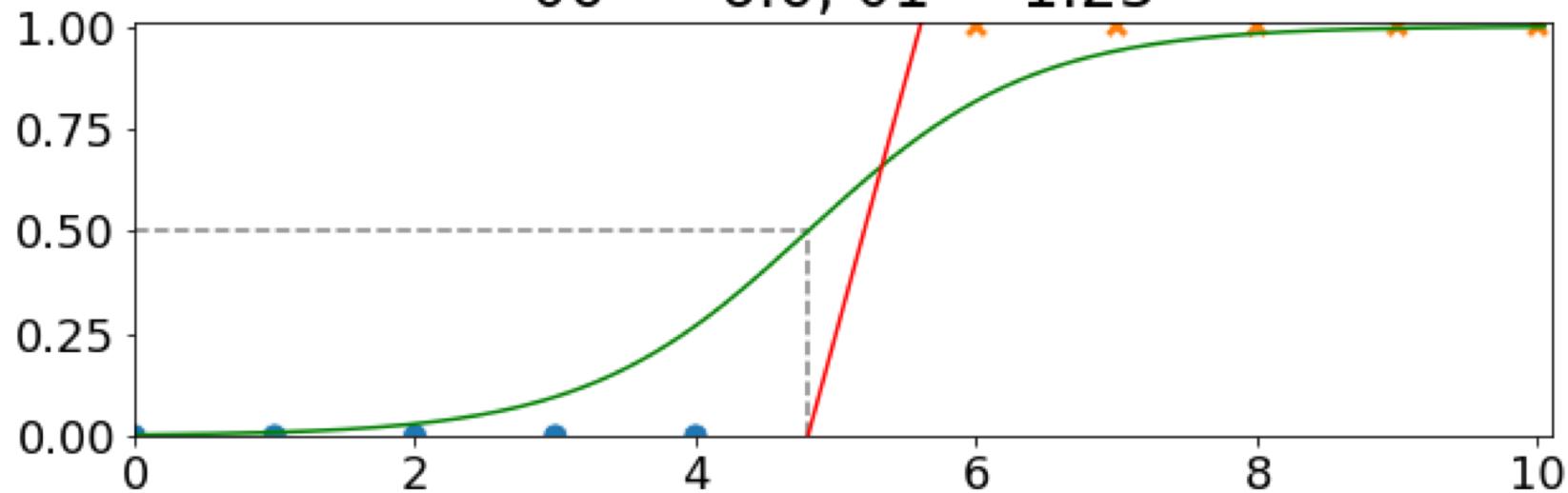
$$\theta_0 = -6.0, \theta_1 = 0.9$$



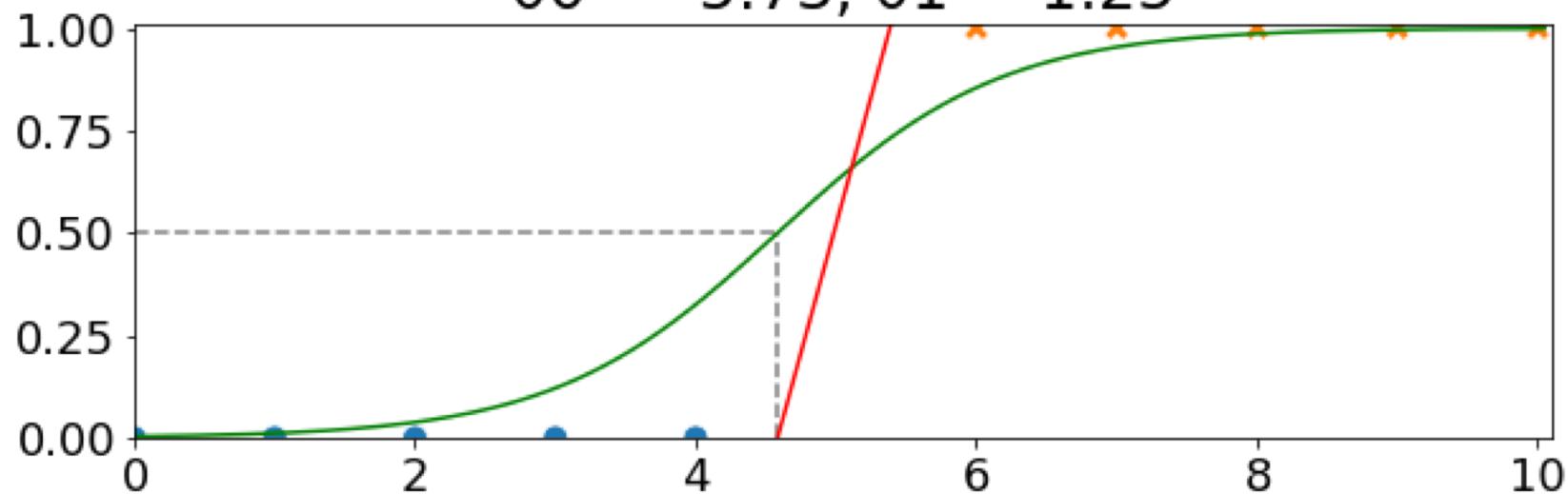
$$\theta_0 = -6.0, \theta_1 = 1.14$$

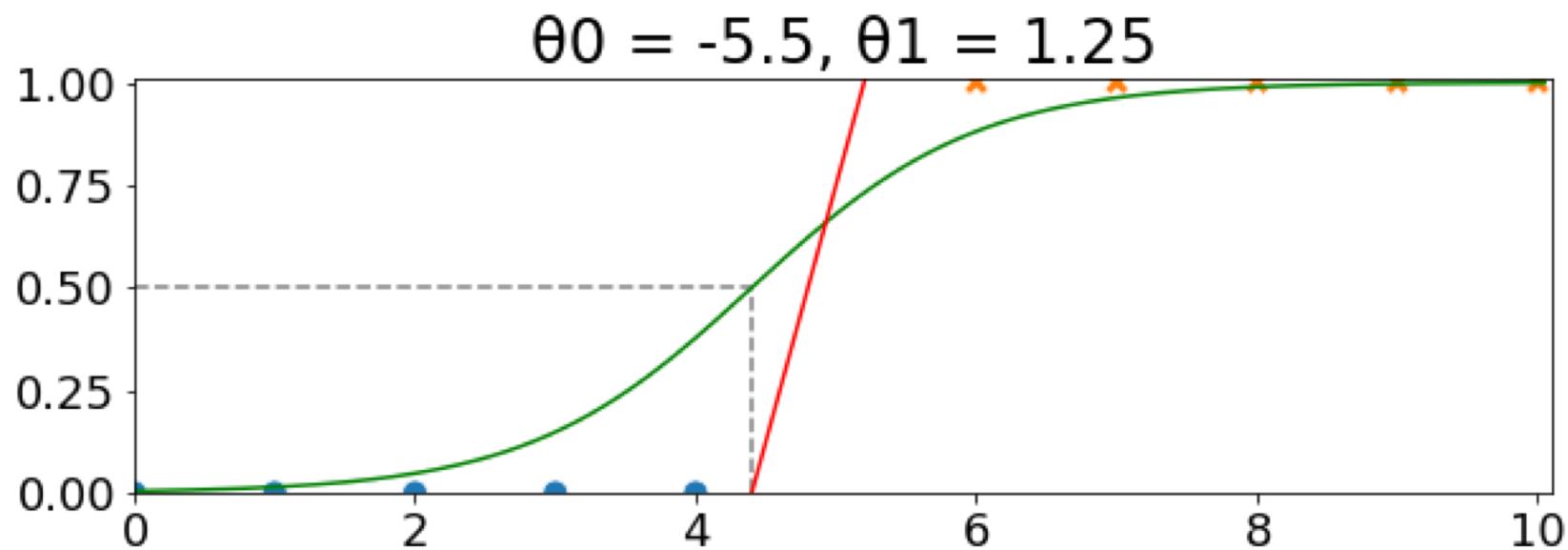


$$\theta_0 = -6.0, \theta_1 = 1.25$$

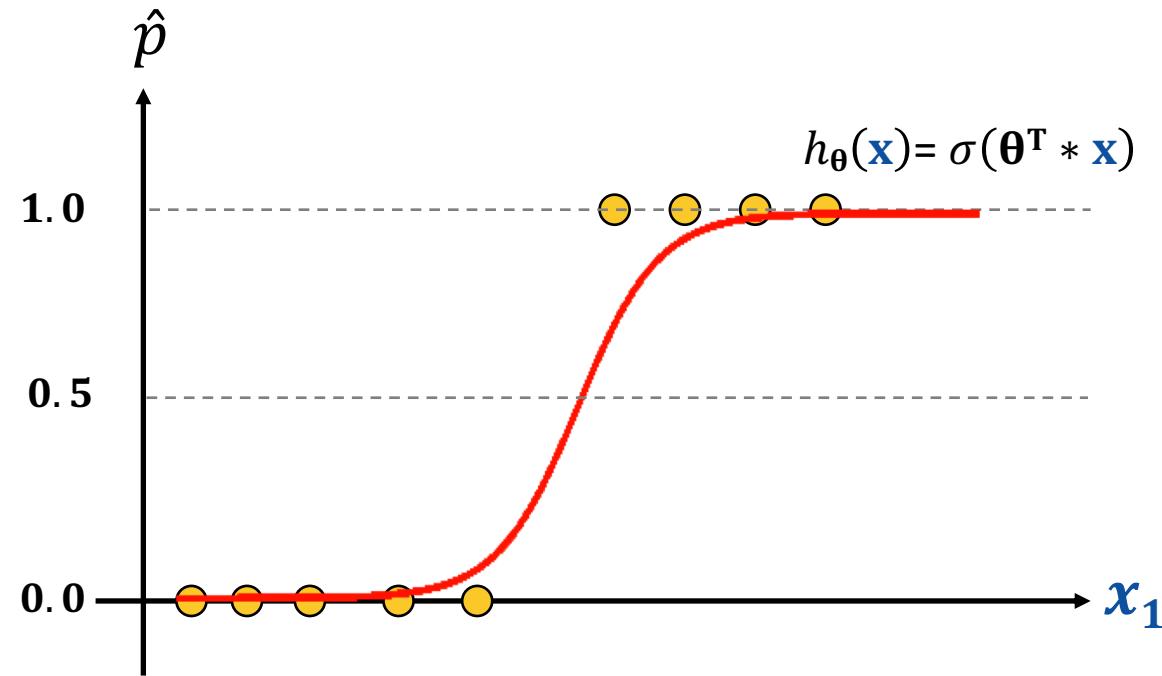


$$\theta_0 = -5.73, \theta_1 = 1.25$$

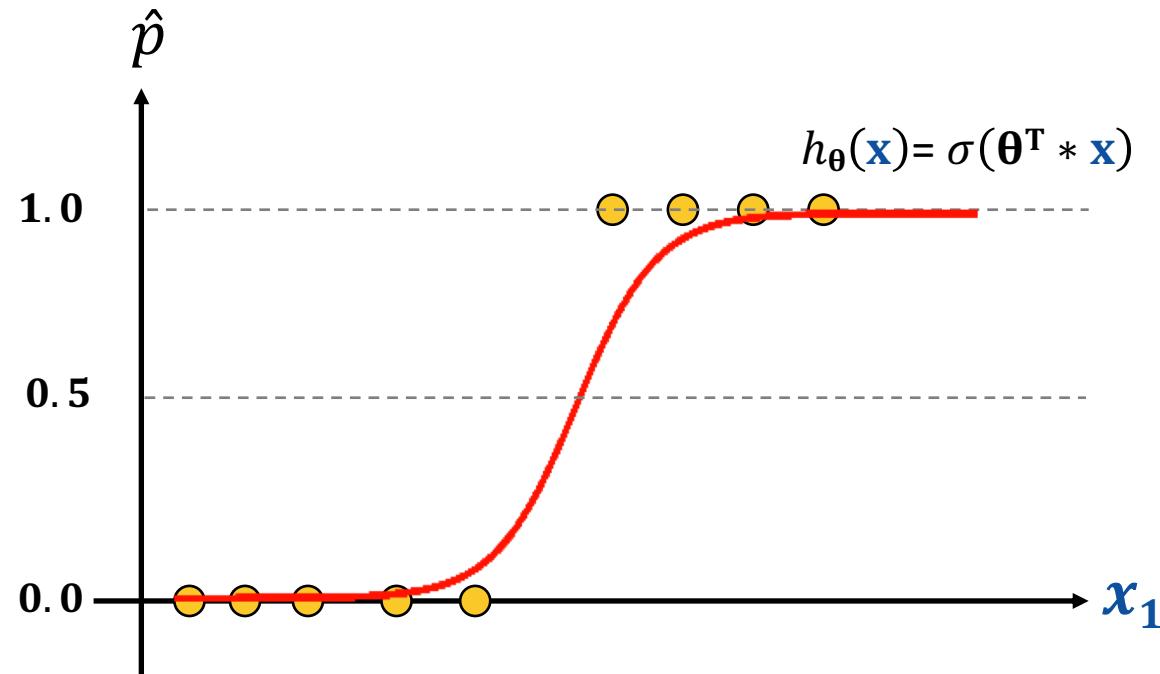




# Decision Boundary



# Decision Boundary



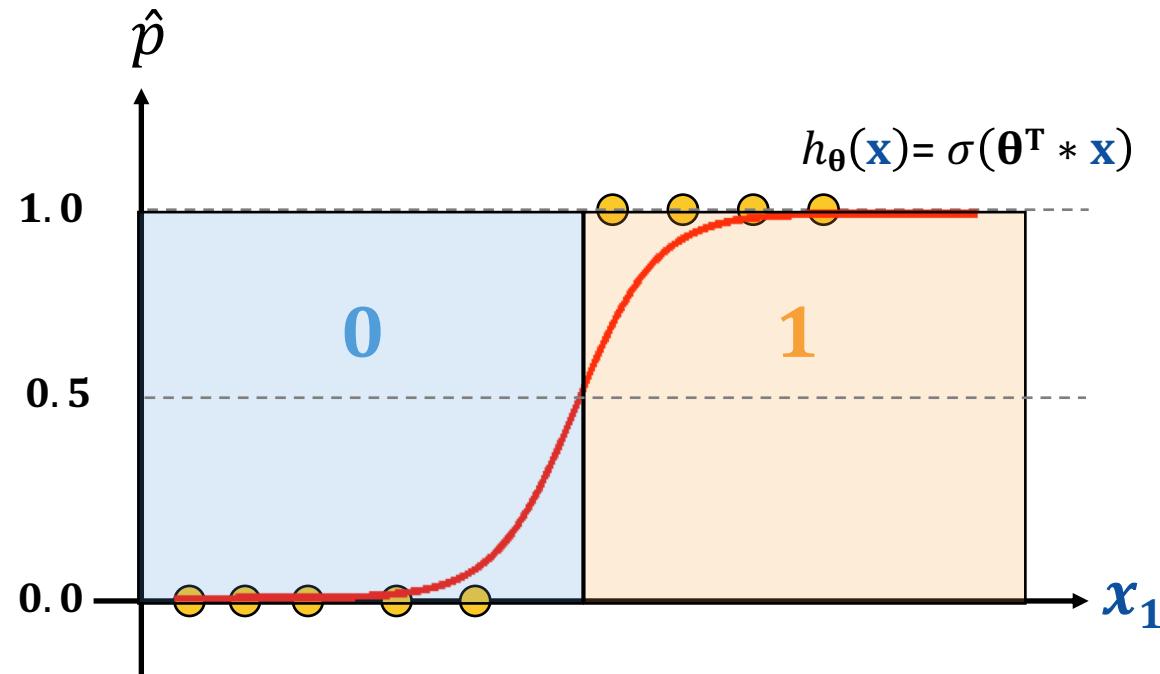
$\sigma(z)$  returns a value between **0** and **1**

$\sigma(z) < 0.5$  if  $z < 0$  (**negative class**)

$\sigma(z) \geq 0.5$  if  $z \geq 0$  (**positive class**)

$\sigma(z) = 0.5$  if  $z = 0$

# Decision Boundary



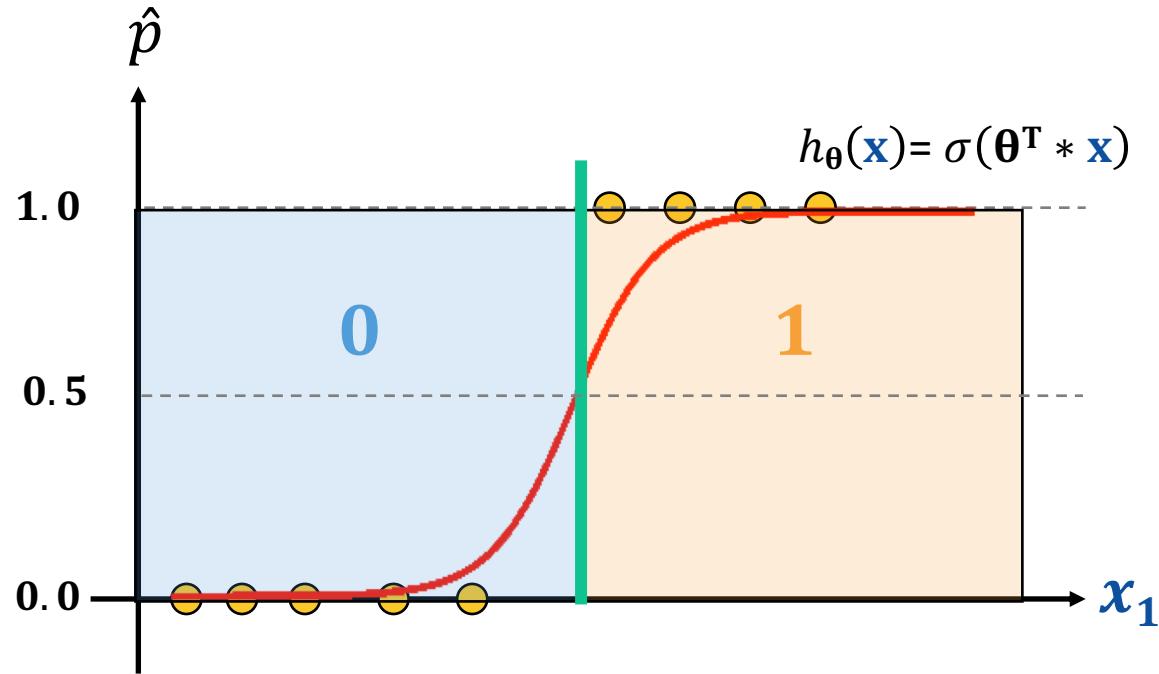
$\sigma(z)$  returns a value between 0 and 1

$\sigma(z) < 0.5$  if  $z < 0$  (negative class)

$\sigma(z) \geq 0.5$  if  $z \geq 0$  (positive class)

$\sigma(z) = 0.5$  if  $z = 0$

# Decision Boundary



$\sigma(z)$  returns a value between **0** and **1**

$\sigma(z) < 0.5$  if  $z < 0$  (**negative class**)

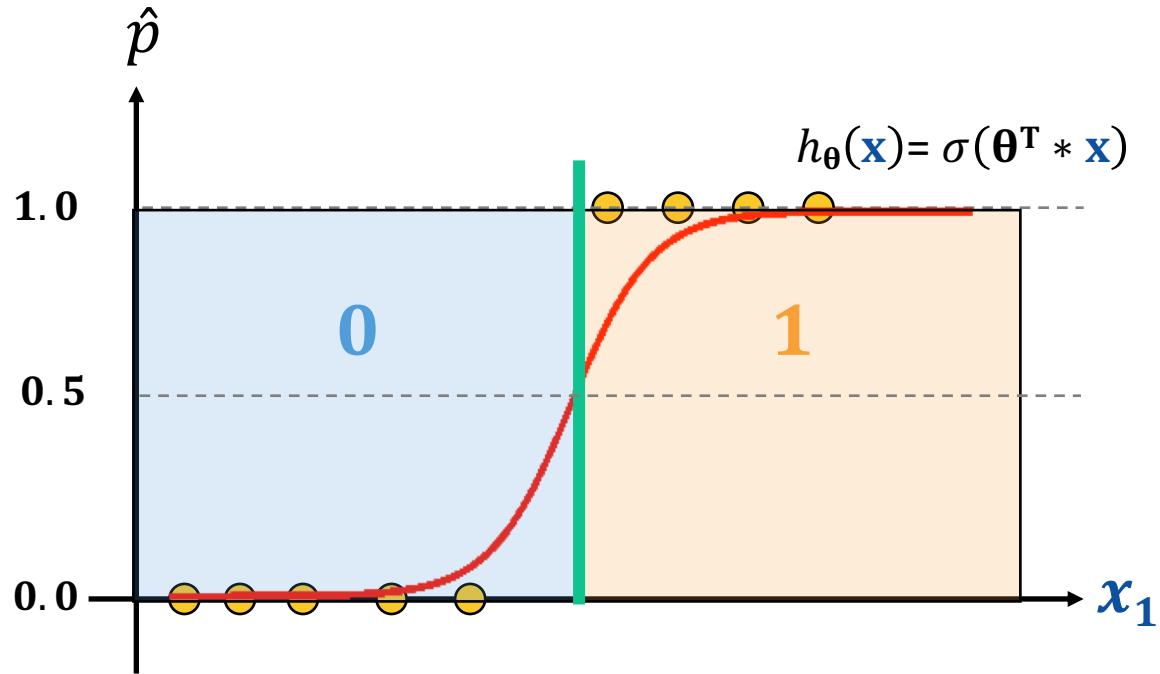
$\sigma(z) \geq 0.5$  if  $z \geq 0$  (**positive class**)

$\sigma(z) = 0.5$  if  $z = 0$



Thus, the hyperplane  $(\boldsymbol{\theta}^T * \mathbf{x}) = 0$  is the **decision boundary** of the classifier..

# Decision Boundary



$\sigma(z)$  returns a value between **0** and **1**

$\sigma(z) < 0.5$  if  $z < 0$  (**negative class**)

$\sigma(z) \geq 0.5$  if  $z \geq 0$  (**positive class**)

$\sigma(z) = 0.5$  if  $z = 0$

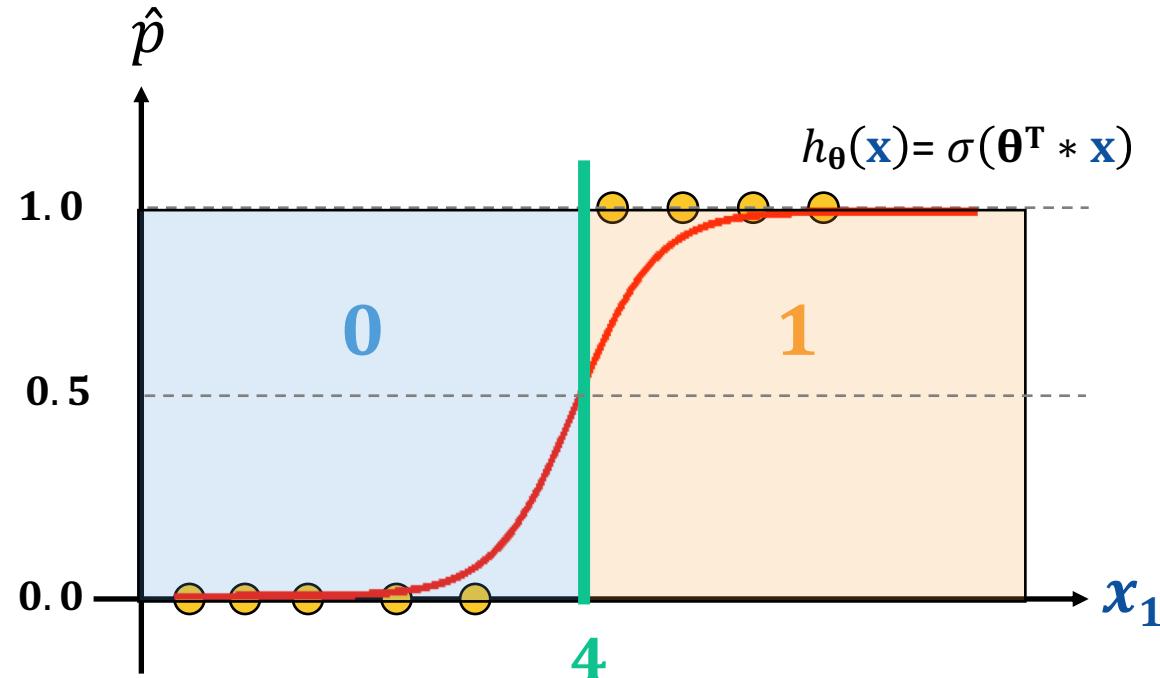


Thus, the hyperplane  $(\theta^T * x) = 0$  is the **decision boundary** of the classifier..

Suppose the **best parameter vector  $\theta$**  is:

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} = \begin{bmatrix} -4 \\ 1 \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ x_1 \end{bmatrix}$$

# Decision Boundary



$\sigma(z)$  returns a value between **0** and **1**

$\sigma(z) < 0.5$  if  $z < 0$  (**negative class**)

$\sigma(z) \geq 0.5$  if  $z \geq 0$  (**positive class**)

$\sigma(z) = 0.5$  if  $z = 0$



Thus, the hyperplane  $(\theta^T * x) = 0$  is the **decision boundary** of the classifier..

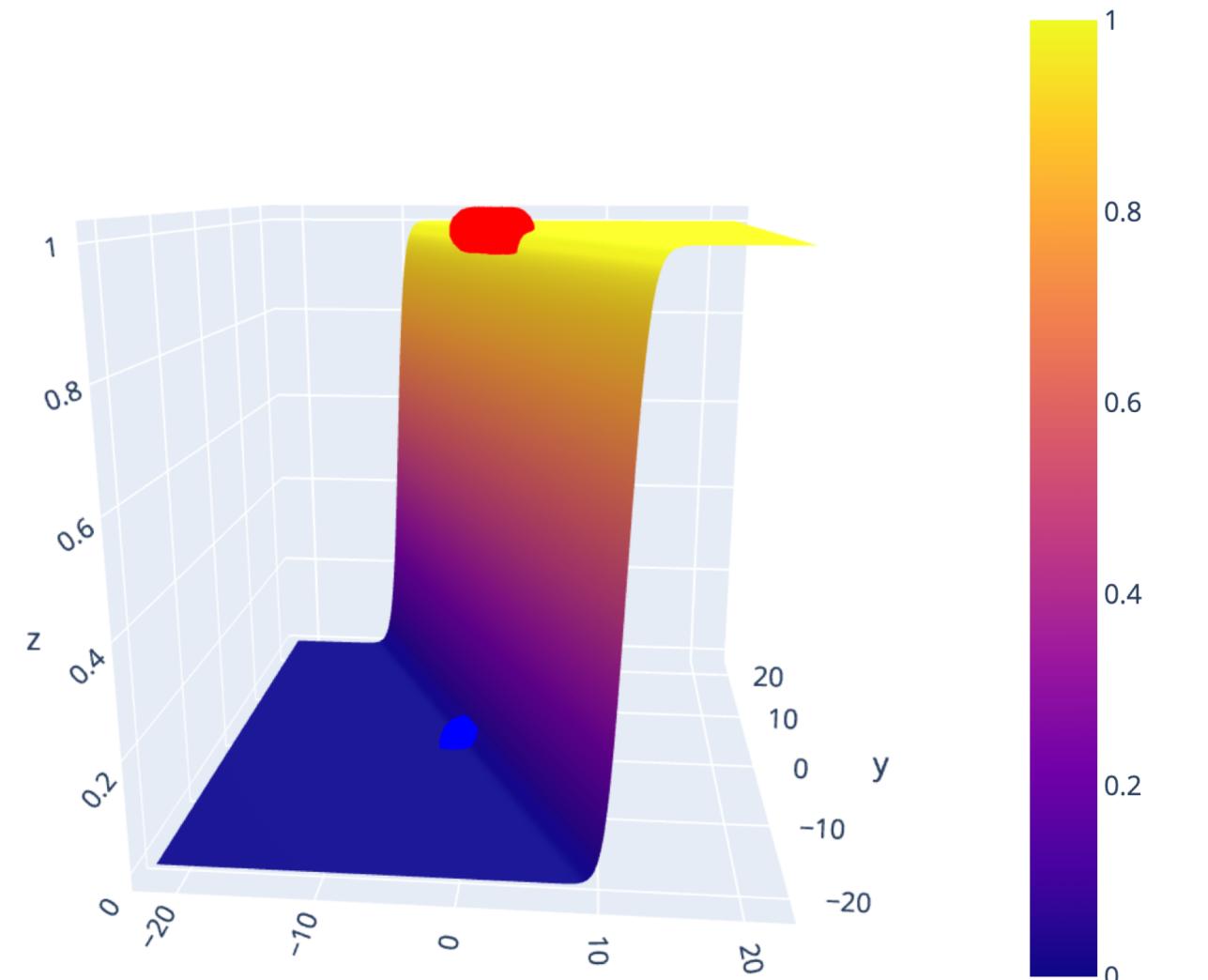
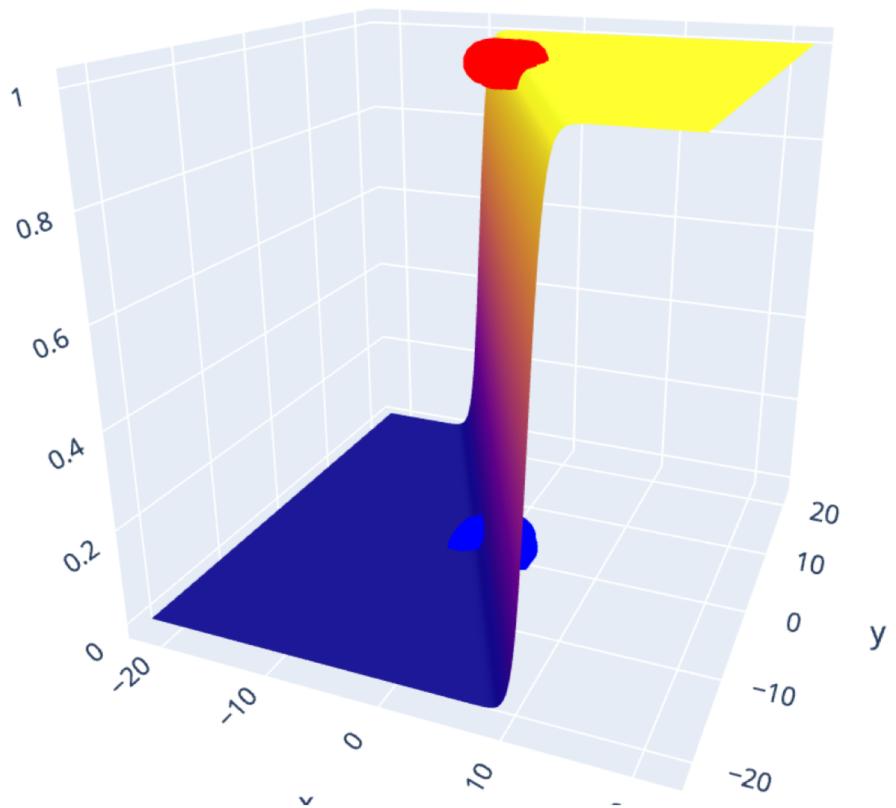
Suppose the **best parameter vector  $\theta$**  is:

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} = \begin{bmatrix} -4 \\ 1 \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ x_1 \end{bmatrix}$$

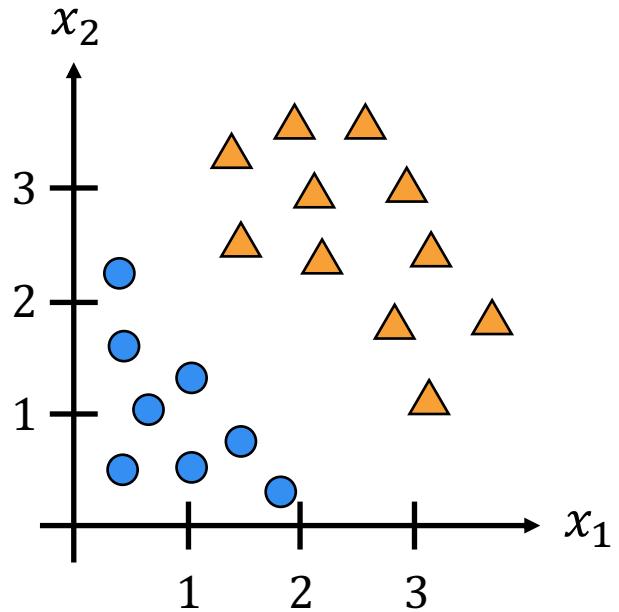
$$\begin{aligned} \theta^T * x &= 0 \\ -4 + 1 * x_1 &= 0 \\ x_1 &= 4 \end{aligned}$$

**decision boundary**

# Decision Boundary for two features



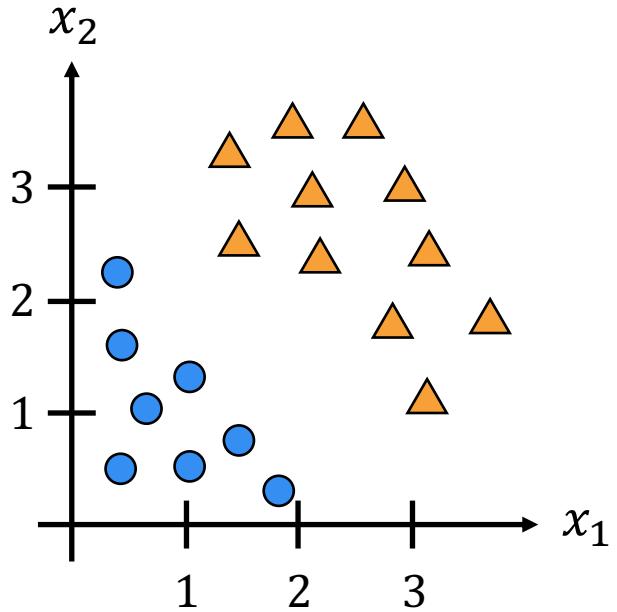
# Decision Boundary for two features



The hyperplane ( $\theta^T * \mathbf{x}$ ) = 0 is the **decision boundary** of the classifier.

$$h_{\theta}(\mathbf{x}) = \sigma(\theta^T * \mathbf{x}) = \sigma(\theta_0 + \theta_1 * x_1 + \theta_2 * x_2)$$

# Decision Boundary for two features



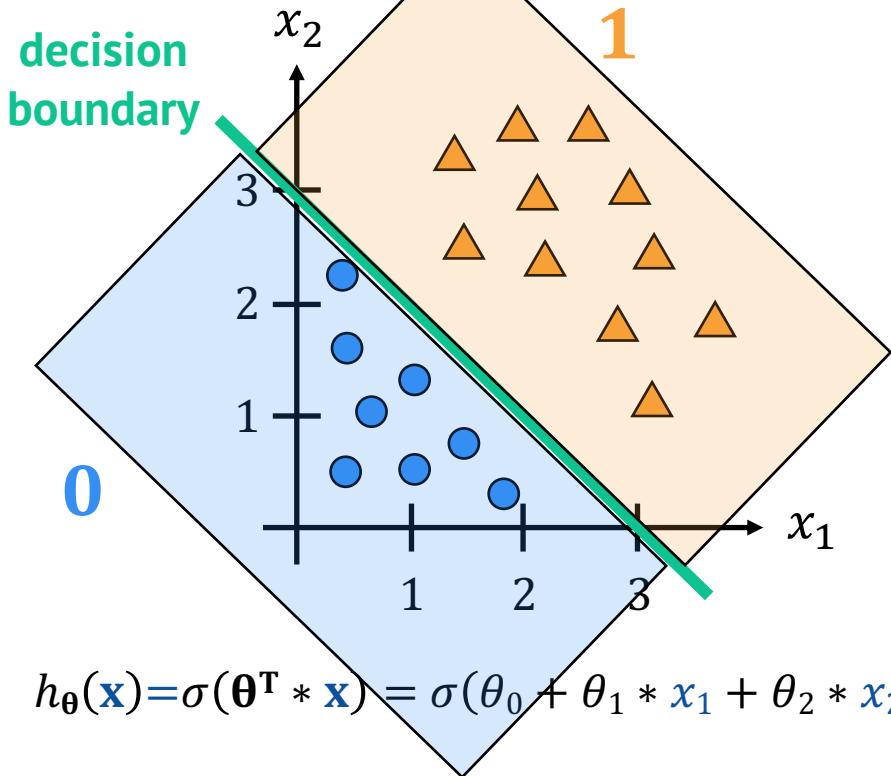
The hyperplane  $(\boldsymbol{\theta}^T * \mathbf{x}) = 0$  is the **decision boundary** of the classifier.

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \sigma(\boldsymbol{\theta}^T * \mathbf{x}) = \sigma(\theta_0 + \theta_1 * x_1 + \theta_2 * x_2)$$

Suppose the **best parameter vector  $\boldsymbol{\theta}$**  is:

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

# Decision Boundary for two features



$$h_{\theta}(\mathbf{x}) = \sigma(\boldsymbol{\theta}^T * \mathbf{x}) = \sigma(\theta_0 + \theta_1 * x_1 + \theta_2 * x_2)$$

Suppose the **best parameter vector  $\theta$**  is:

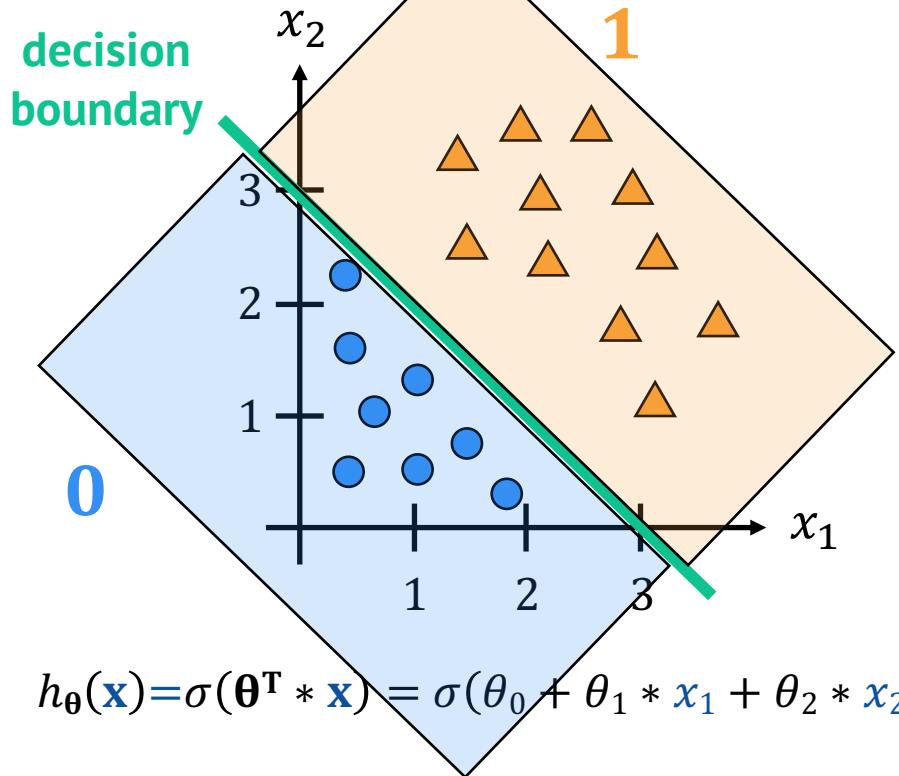
$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

The hyperplane  $(\boldsymbol{\theta}^T * \mathbf{x}) = 0$  is the **decision boundary** of the classifier.

$$\begin{aligned} \boldsymbol{\theta}^T * \mathbf{x} &= 0 \\ -3 + 1 * x_1 + 1 * x_2 &= 0 \\ \boxed{x_1 + x_2} &= 3 \end{aligned}$$

**decision boundary**

# Decision Boundary for two features



Suppose the **best parameter vector**  $\theta$  is:

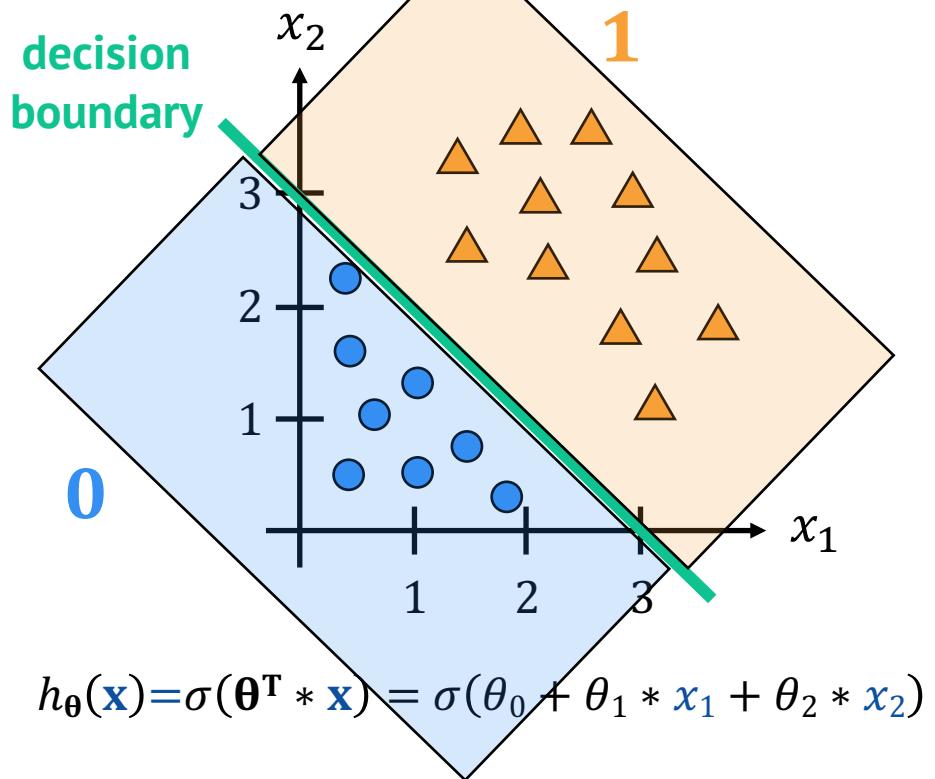
$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

**Logistic Regression** is a **linear classifier** since its estimated **decision boundary** is **linear**.

$$\begin{aligned} \boldsymbol{\theta}^T * \mathbf{x} &= 0 \\ -3 + 1 * x_1 + 1 * x_2 &= 0 \\ x_1 + x_2 &= 3 \end{aligned}$$

*decision boundary*

# Decision Boundary for two features



Suppose the **best parameter vector  $\theta$**  is:

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

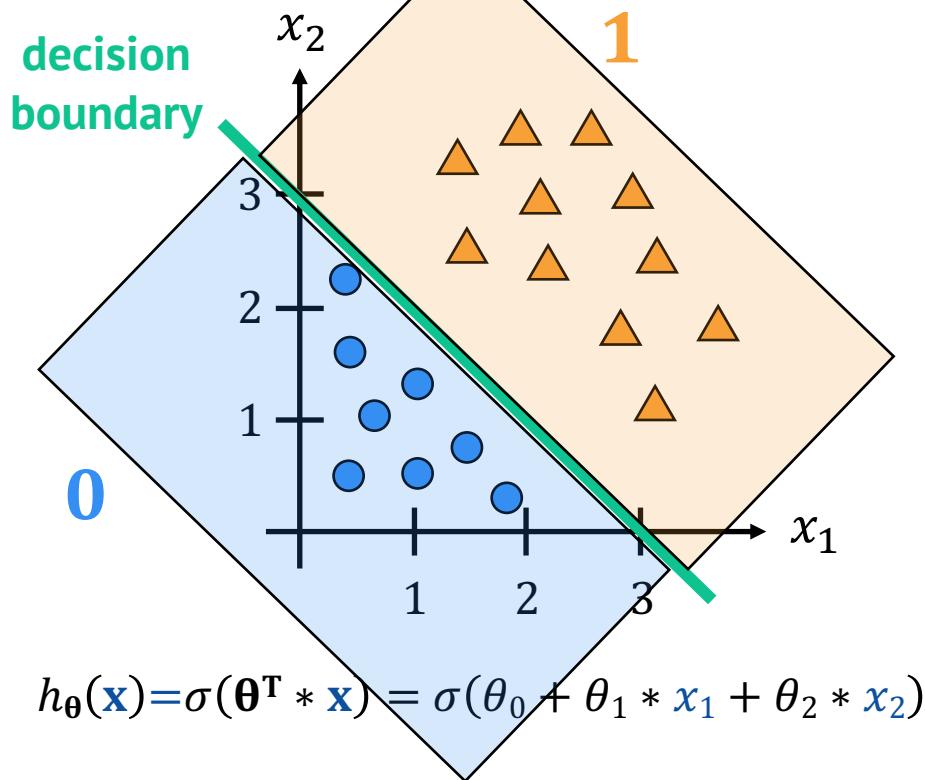
**Logistic Regression** is a **linear classifier** since its estimated **decision boundary** is **linear**.

**Logistic Regression** is used for **BINARY classification problems**.

$$\begin{aligned} \boldsymbol{\theta}^T \cdot \mathbf{x} &= 0 \\ -3 + 1 * x_1 + 1 * x_2 &= 0 \\ x_1 + x_2 &= 3 \end{aligned}$$

*decision boundary*

# Decision Boundary for two features



Suppose the best parameter vector  $\theta$  is:

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

**Logistic Regression** is a **linear classifier** since its estimated **decision boundary** is **linear**.

**Logistic Regression** is used for **BINARY classification problems**.

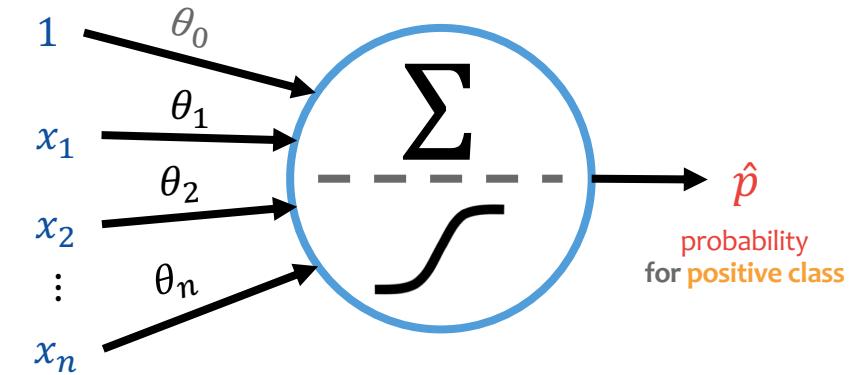
There are some **strategies** to extend **Logistic Regression** for **multiclass classification**. We'll see them soon!

# Logistic Regression and Artificial Neurons

# Logistic Regression and Artificial Neurons



A **Logistic Regression** classifier is exactly an **Artificial Neuron (Perceptron)** with the **sigmoid** as activation function.



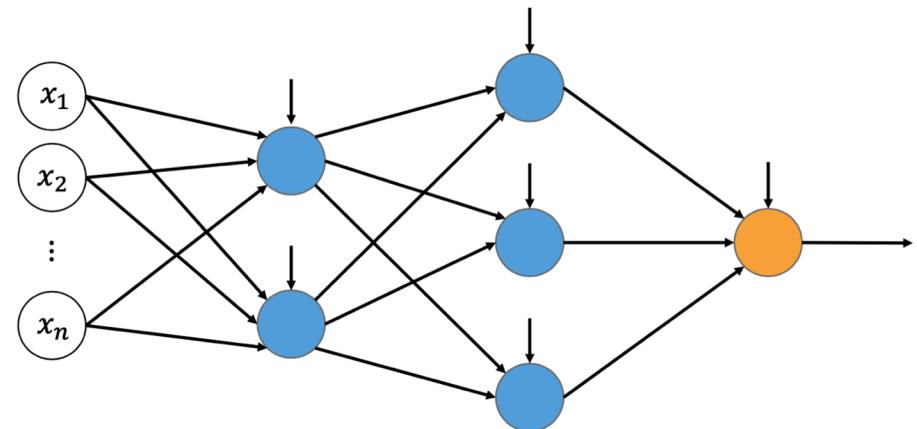
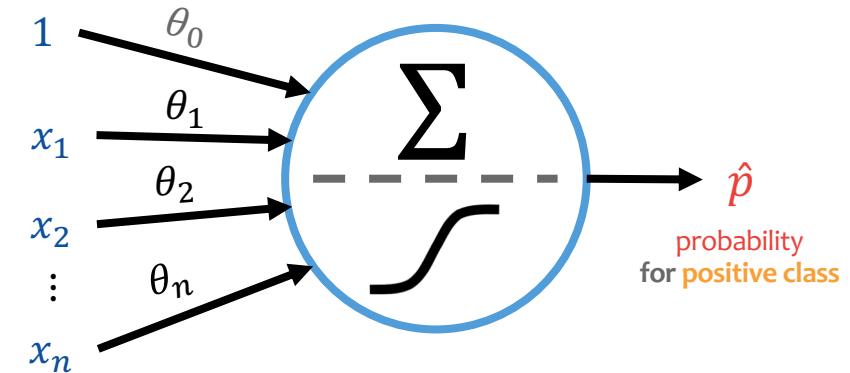
# Logistic Regression and Artificial Neurons



A **Logistic Regression** classifier is exactly an **Artificial Neuron (Perceptron)** with the **sigmoid** as activation function.



**Artificial Neurons** are the basic units of **Artificial Neural Networks**.



# Hands-on

# D1EAD – Análise Estatística para Ciência de Dados



## Logistic Regression

Lecture based on  
Machine Learning course  
by Andrew Ng  
(Coursera)

Prof. Samuel Martins (Samuka)  
[samuel.martins@ifsp.edu.br](mailto:samuel.martins@ifsp.edu.br)

