# Gradient-Based Optimization

Elena Congeduti, 14-11-2024

# Lecture's Agenda

1. Loss Functions

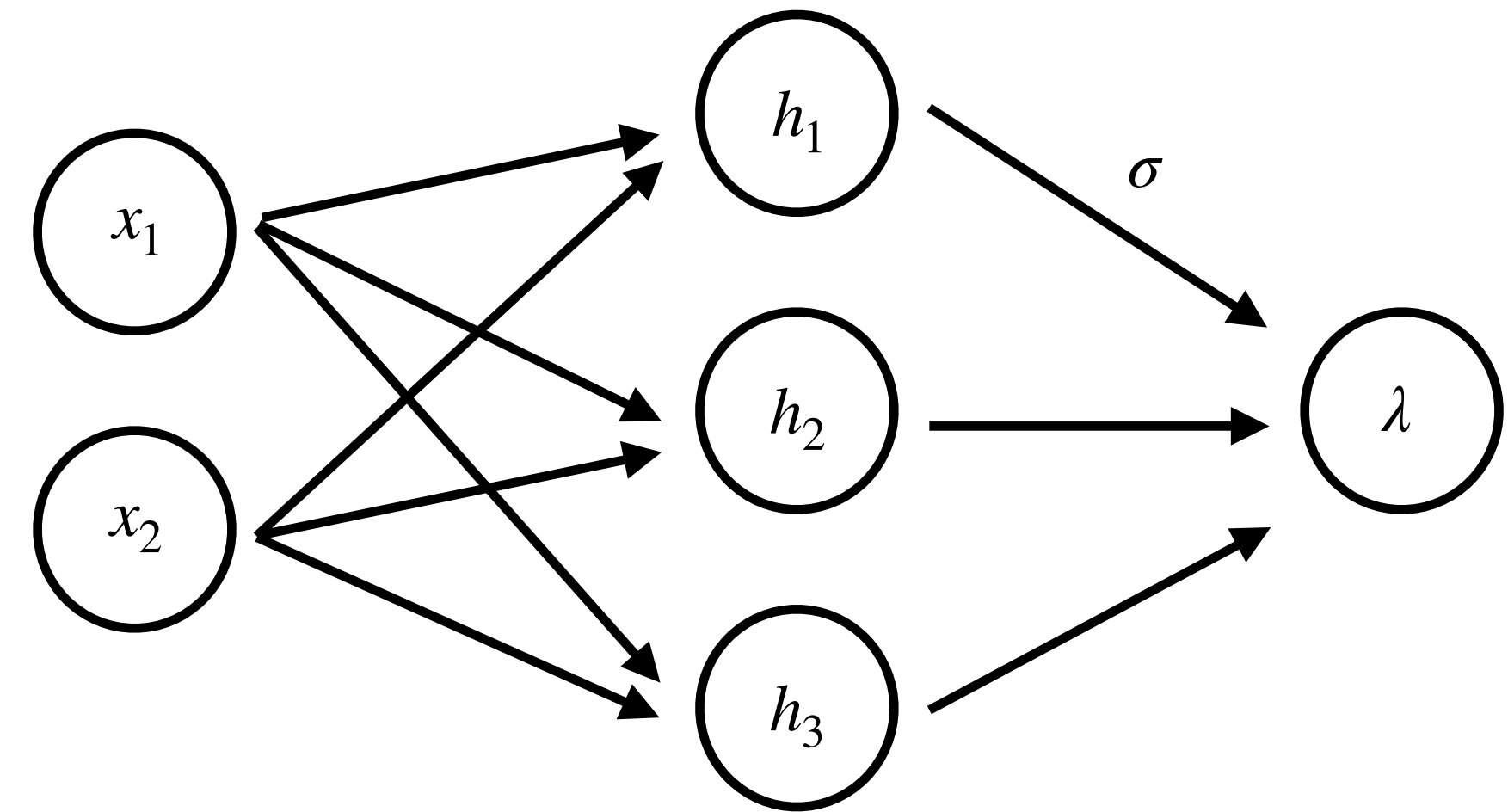2. Maximum Likelihood

3. Gradient Descent

4. Backpropagation

# Training loop for classification

Training set $\{x^{(i)}, y^{(i)}\}_{i=1,\dots,n}$ of input/output pairs

Iterate until convergence

1. Forward pass: for batch $x$ compute output $\hat{\lambda} = f(x; \theta)$

2. **Evaluate: compare the $\hat{\lambda}$ with the class label $y$**

3. Backward pass: update the parameters $\theta$
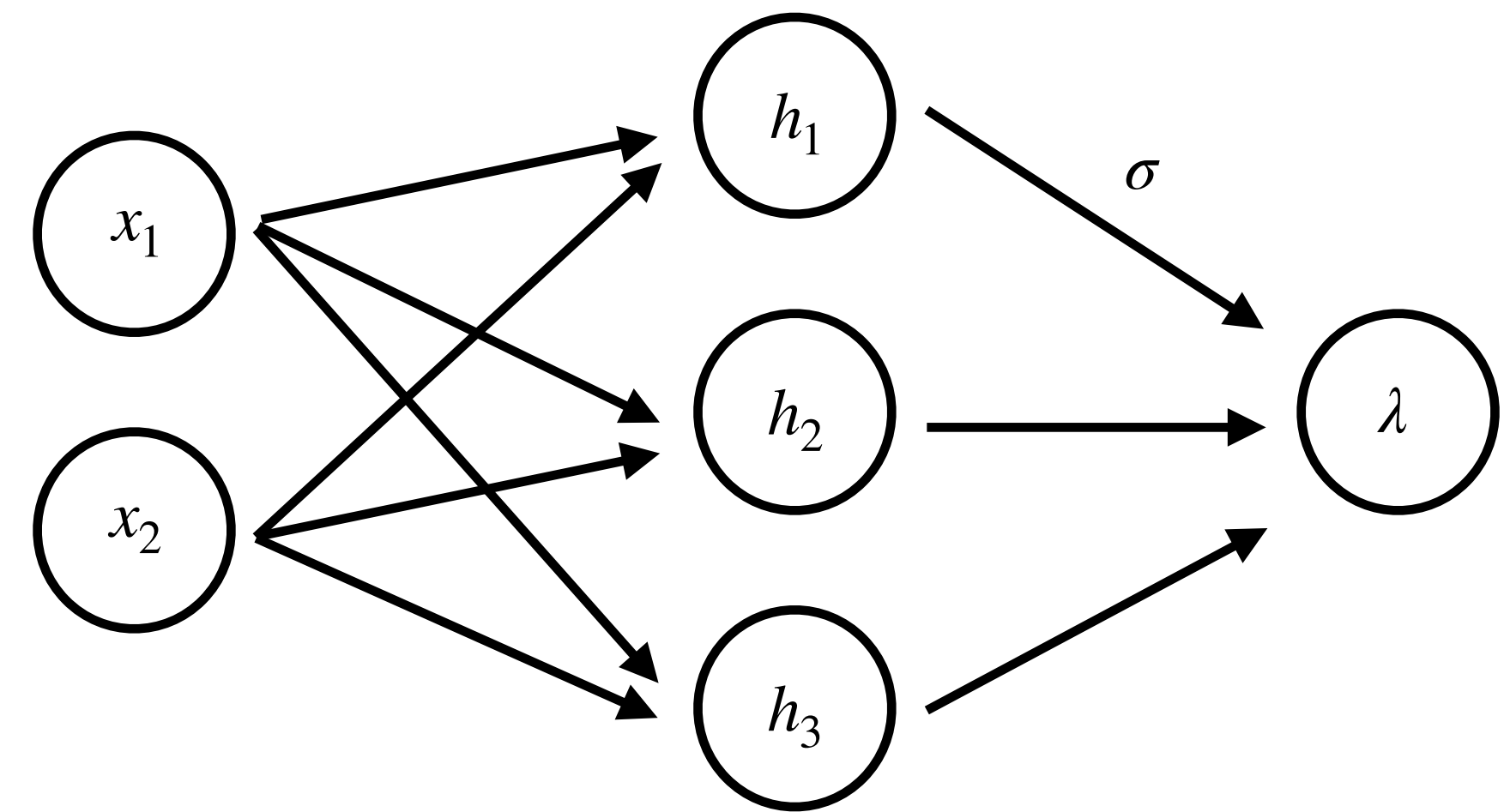
Scalar loss function $L(\theta)$

# Training loop for classification

Training set $\{x^{(i)}, y^{(i)}\}_{i=1,\ldots,n}$ of input/output pairs

Iterate until convergence

1. Forward pass: for batch $x$ compute output $\hat{\lambda} = f(x; \theta)$

2. **Evaluate: compare the $\hat{\lambda}$ with the class label $y$**

3. Backward pass: update the parameters $\theta$

Scalar loss function $L(\theta)$

- Error between ground truth $y$ and network output $\hat{\lambda}$
- Suitable for gradient learning ('usable' gradient)



4

# Binary Classification Loss

Which is a good candidate loss when the last activation returns $\hat{\lambda} = \sigma(\ldots)$?
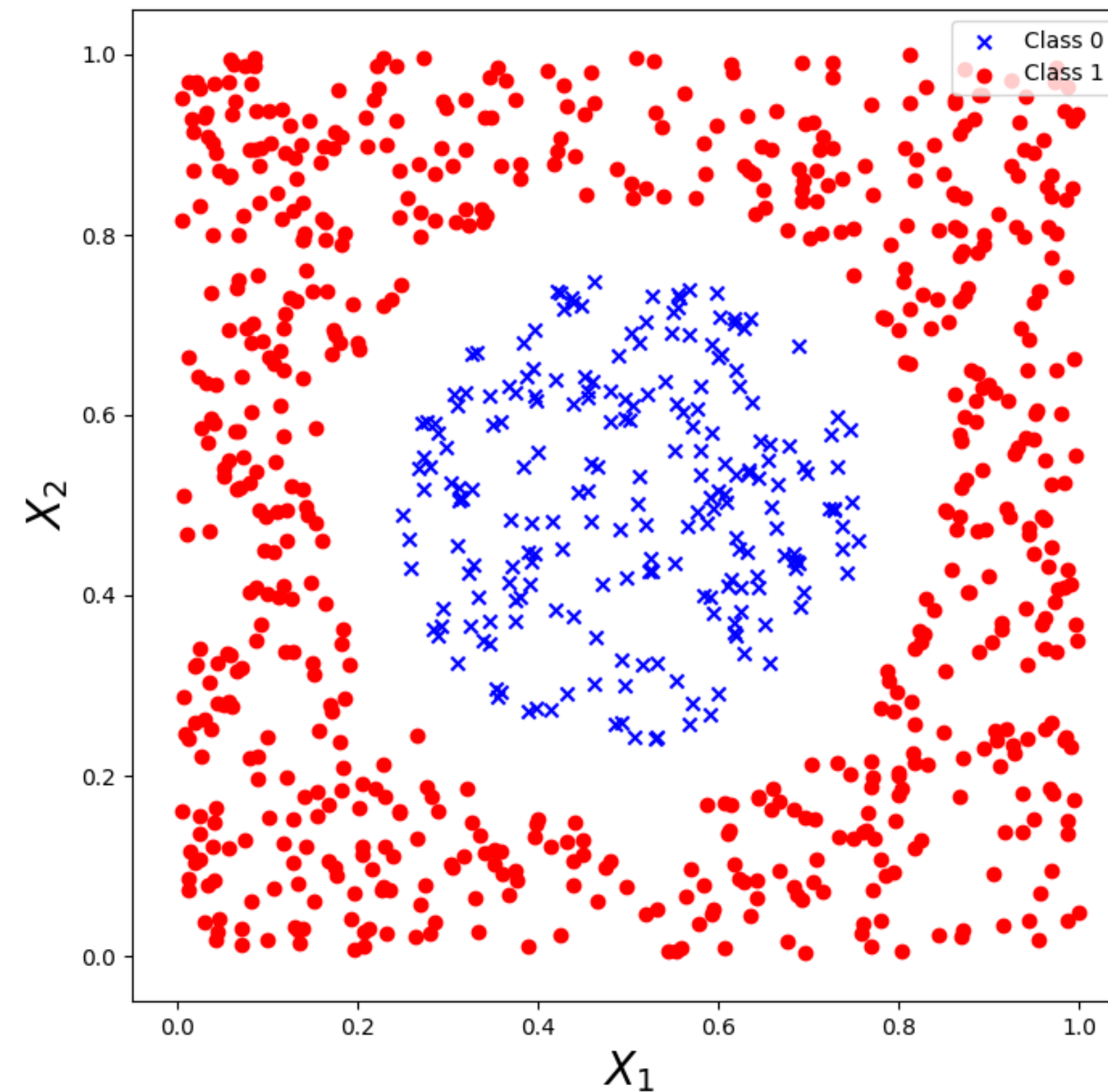
A. $L_{ratio} = \dfrac{y}{\hat{\lambda}}$

B. $L_{0-1} = \begin{cases} 0 & \text{if } y = \hat{\lambda} \\ 1 & \text{if } y \neq \hat{\lambda} \end{cases}$

C. $L_{SE} = \dfrac{1}{2}(y - \hat{\lambda})^2$

D. $L_1 = |y - \hat{\lambda}|$

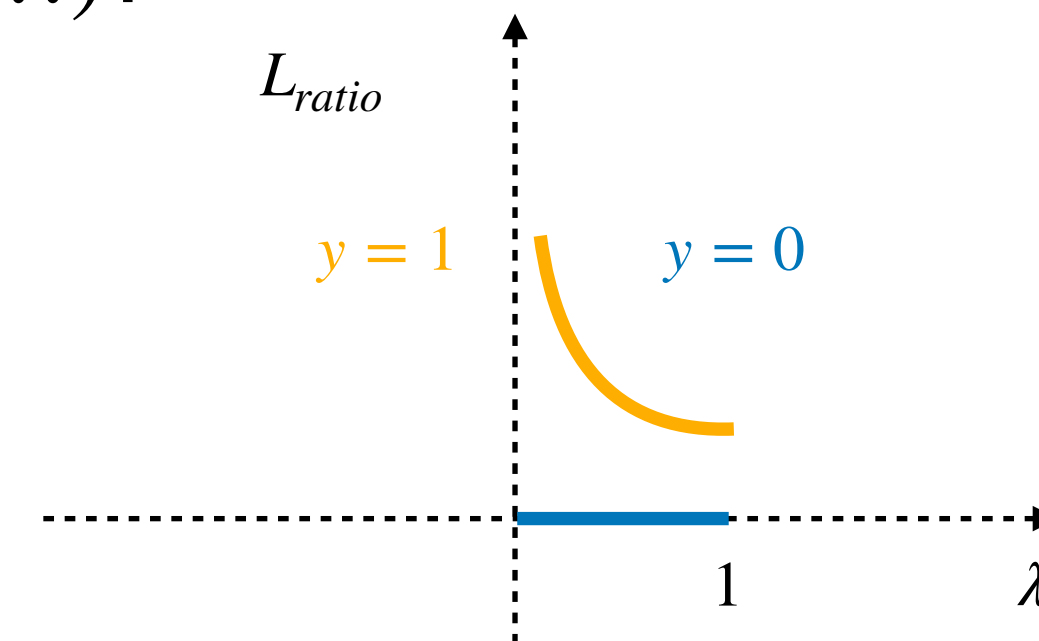E. $L_{CE} = -y\log(\hat{\lambda}) - (1-y)\log(1-\hat{\lambda})$

# Binary Classification Loss

Which is a good candidate loss when the last activation returns $\hat{\lambda} = \sigma(\ldots)$?
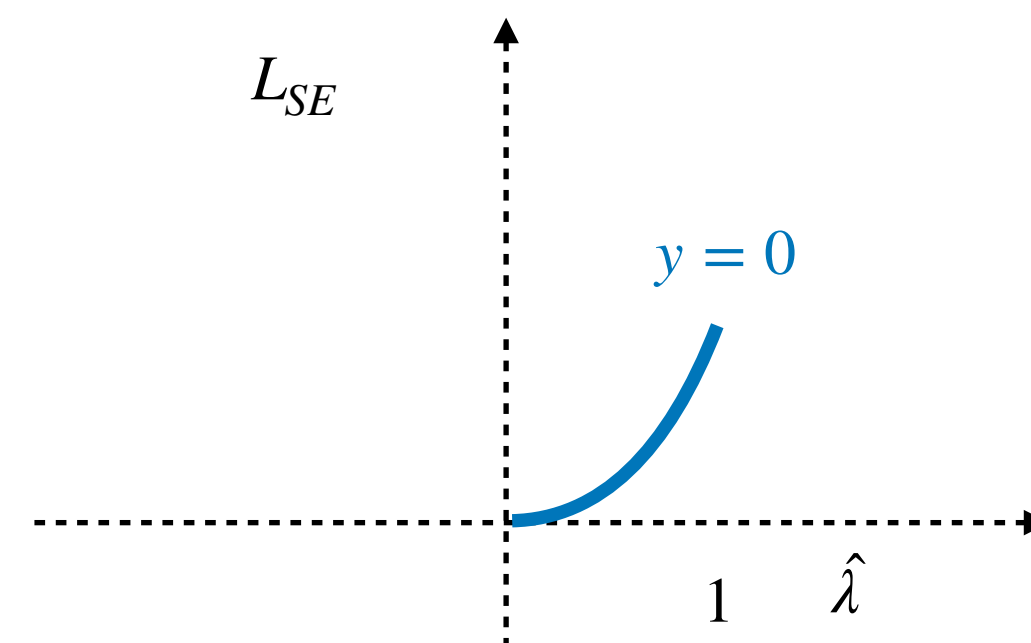
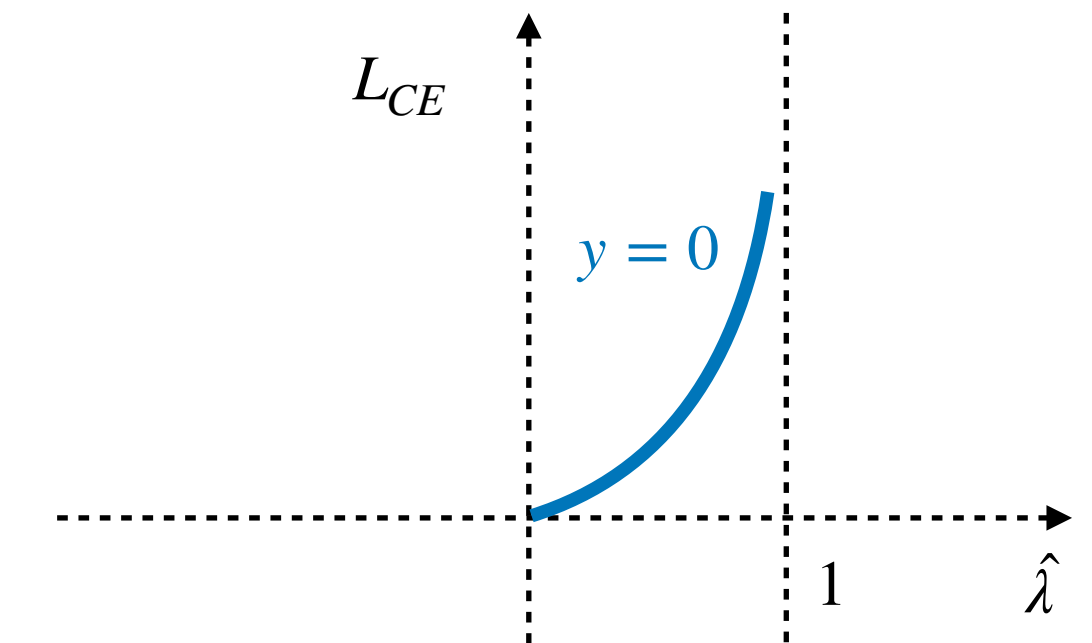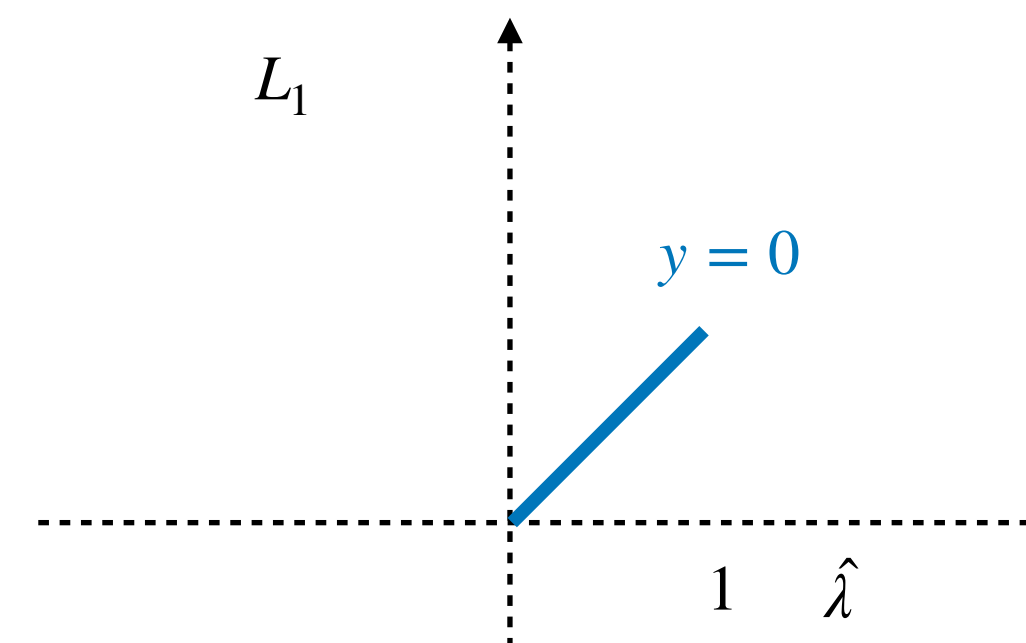A. $L_{ratio} = \dfrac{y}{\hat{\lambda}}$  $\longrightarrow$  No error measure for $y = 0$

B. $L_{0-1} = \begin{cases} 0 & \text{if } y = \hat{\lambda} \\ 1 & \text{if } y \neq \hat{\lambda} \end{cases}$  $\longrightarrow$  always 0
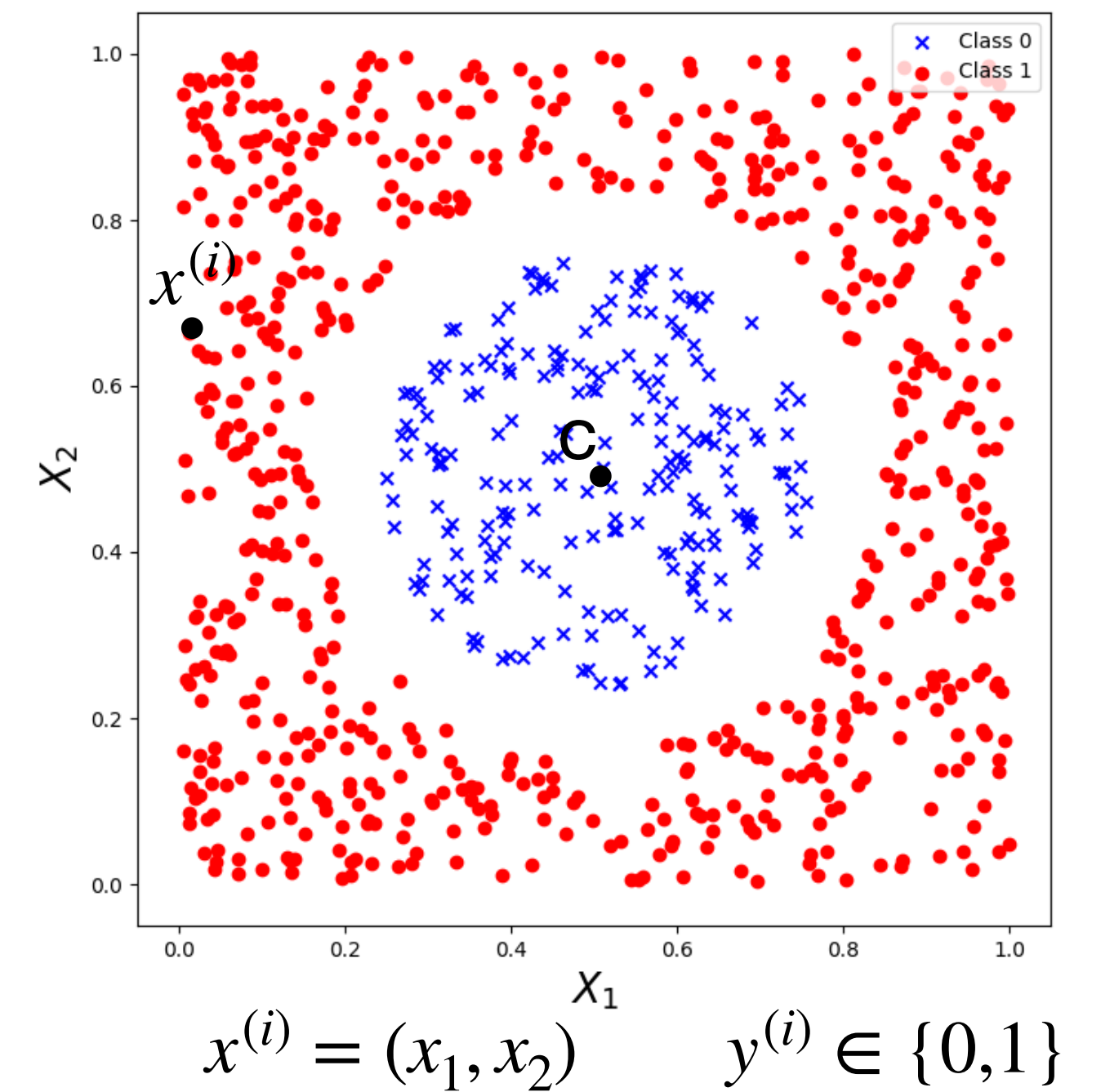
C. $L_{SE} = \dfrac{1}{2}(y - \hat{\lambda})^2$

D. $L_1 = |y - \hat{\lambda}|$

E. $L_{CE} = -y\log(\hat{\lambda}) - (1-y)\log(1-\hat{\lambda})$

# Maximum Likelihood

Training set $\{x^{(i)}, y^{(i)}\}_{i=1,\ldots,n}$ of input/output pairs



$$x^{(i)} = (x_1, x_2) \qquad y^{(i)} \in \{0,1\}$$

# Maximum Likelihood

Training set $\{x^{(i)}, y^{(i)}\}_{i=1,\ldots,n}$ of input/output pairs

1. $y^{(i)}$ drawn from an unknown distribution $p_{data}(y \mid x^{(i)})$



$$x^{(i)} = (x_1, x_2) \qquad y^{(i)} \in \{0,1\}$$

$$p_{data}(y \mid x^{(i)}) = \begin{cases} \delta_0(y) & dist(x^{(i)}, c) \leq r \\ \delta_1(y) & dist(x^{(i)}, c) > r \end{cases}$$

# Maximum Likelihood

Training set $\{x^{(i)}, y^{(i)}\}_{i=1,\dots,n}$ of input/output pairs

1. $y^{(i)}$ drawn from an unknown distribution $p_{data}(y \mid x^{(i)})$



$$x^{(i)} = (x_1, x_2) \qquad y^{(i)} \in \{0,1\}$$

$$p_{data}(y \mid x^{(i)}) = \begin{cases} \delta_0(y) & dist(x^{(i)}, c) \leq r \\ \delta_1(y) & dist(x^{(i)}, c) > r \end{cases}$$

Deterministic distribution
$$P(y = 0) = 1$$
$$P(y = 1) = 0$$

Euclidean distance

# Maximum Likelihood

Training set $\{x^{(i)}, y^{(i)}\}_{i=1,\ldots,n}$ of input/output pairs

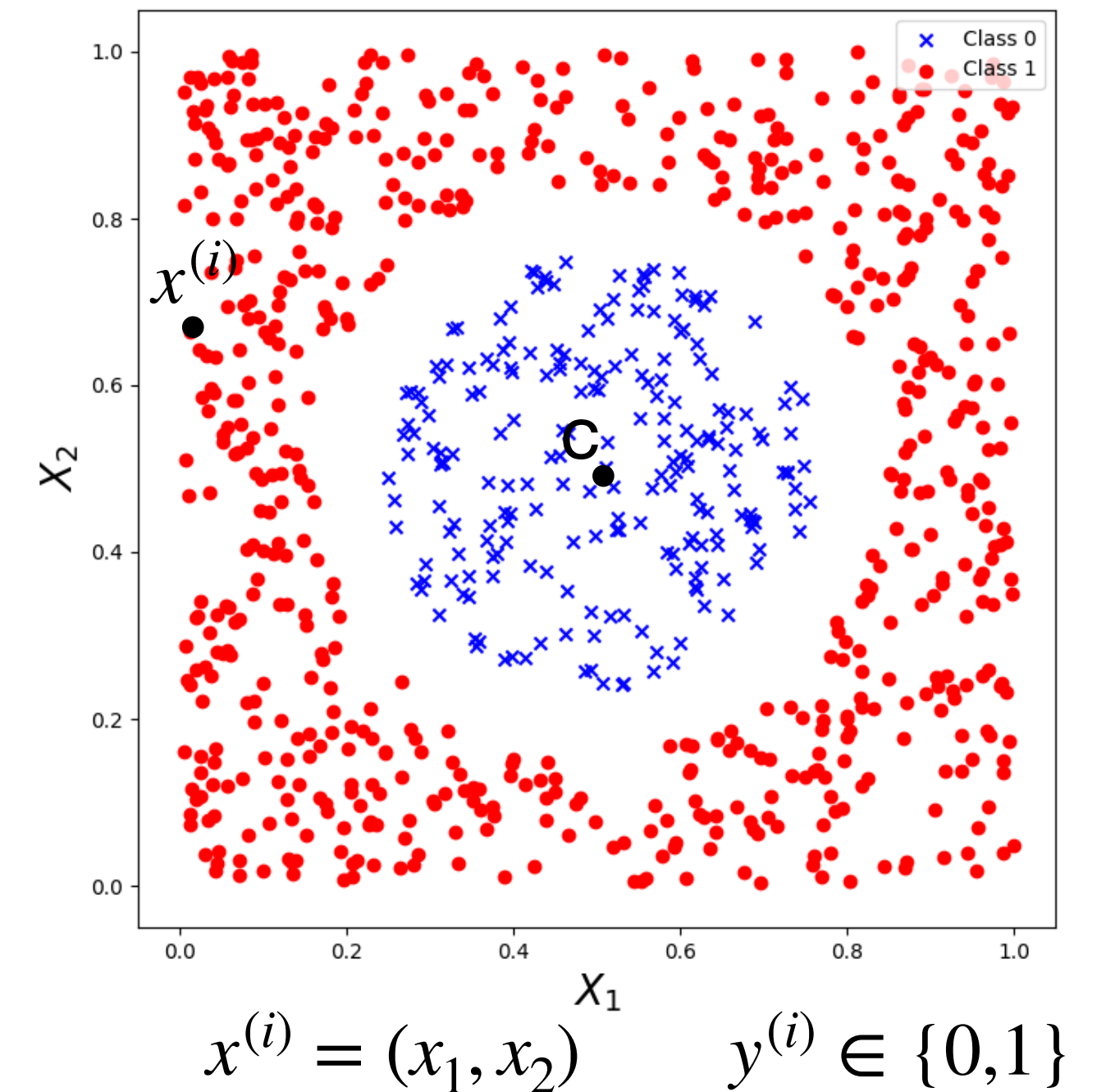1. $y^{(i)}$ drawn from an unknown distribution $p_{data}(y \mid x^{(i)})$



$$x^{(i)} = (x_1, x_2) \qquad y^{(i)} \in \{0,1\}$$

$$p_{data}(y \mid x^{(i)}) = \begin{cases} \delta_0(y) & dist(x^{(i)}, c) \leq r \\ \delta_1(y) & dist(x^{(i)}, c) > r \end{cases}$$

$\{x^{(1)}, \ldots, x^{(n)}\}$, classes $y^{(i)} \sim p_{data}(\cdot \mid x^{(i)})$

# Maximum Likelihood

Training set $\{x^{(i)}, y^{(i)}\}_{i=1,\ldots,n}$ of input/output pairs

1. $y^{(i)}$ drawn from an unknown distribution $p_{data}(y \mid x^{(i)})$

2. Learn an approximation $p_{model}$ of the data distribution



$$x^{(i)} = (x_1, x_2) \qquad y^{(i)} \in \{0,1\}$$

$$p_{model} \approx p_{data}(y \mid x^{(i)}) = \begin{cases} \delta_0(y) & dist(x^{(i)}, c) \leq r \\ \delta_1(y) & dist(x^{(i)}, c) > r \end{cases}$$
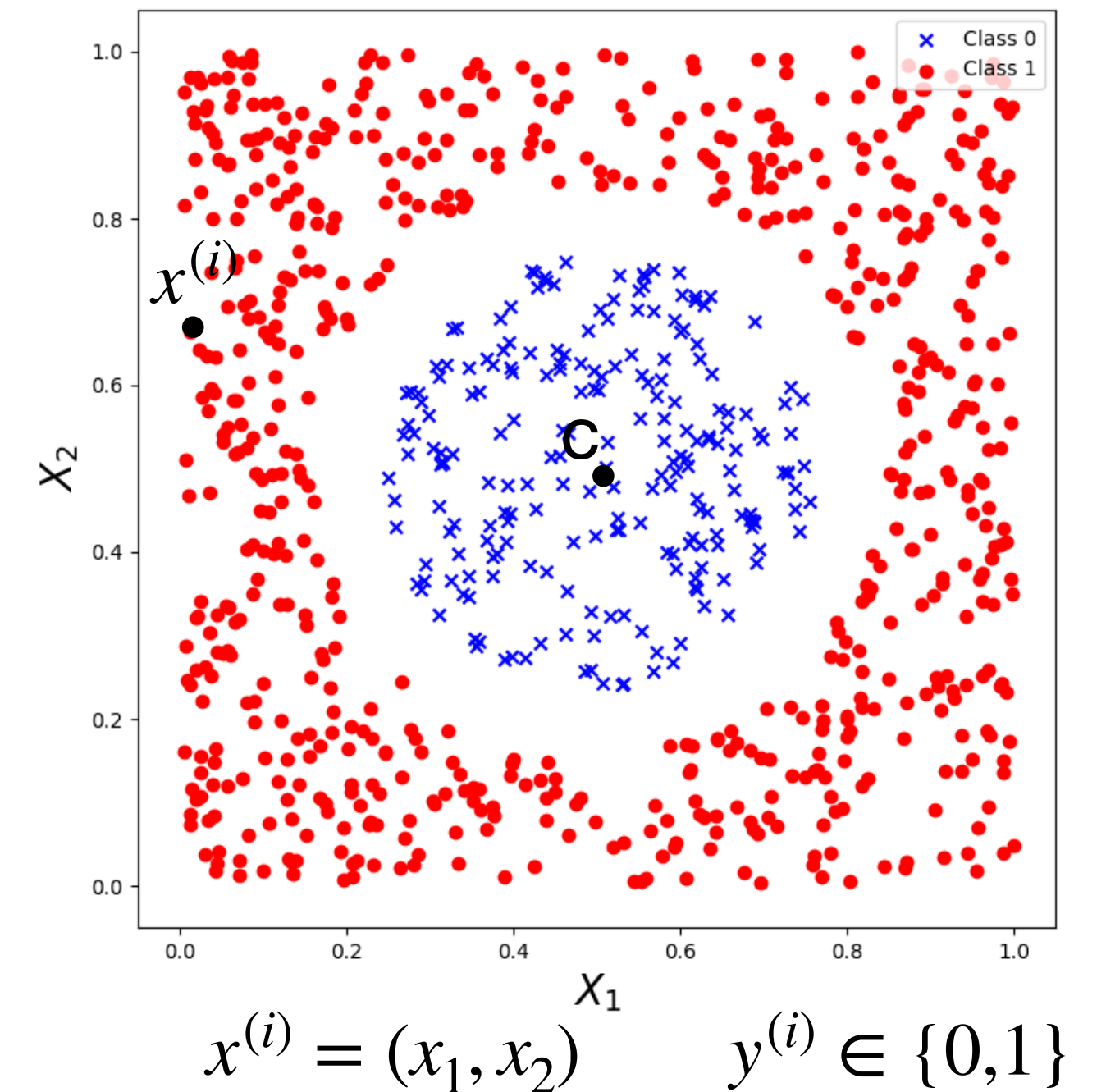
# Maximum Likelihood

Training set $\{x^{(i)}, y^{(i)}\}_{i=1,\ldots,n}$ of input/output pairs

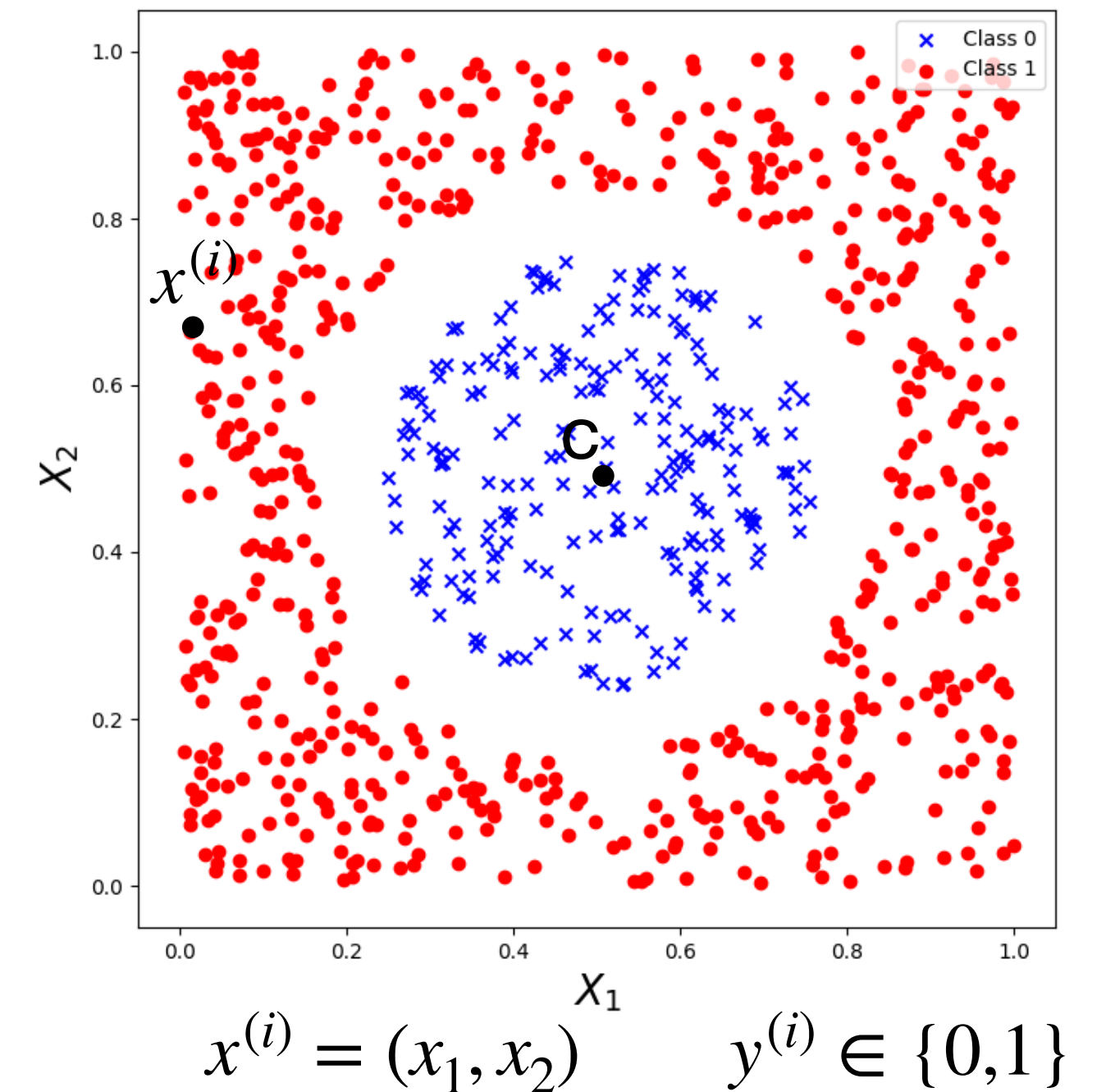1. $y^{(i)}$ drawn from an unknown distribution $p_{data}(y \,|\, x^{(i)})$

2. Learn an approximation $p_{model}$ of the data distribution

3. Distribution $p_{model}(y \,|\, \lambda)$ over the output space



$$x^{(i)} = (x_1, x_2) \qquad y^{(i)} \in \{0,1\}$$

$$p_{model} \approx p_{data}(y \,|\, x^{(i)}) = \begin{cases} \delta_0(y) & dist(x^{(i)}, c) \leq r \\ \delta_1(y) & dist(x^{(i)}, c) > r \end{cases}$$

$p_{model}(y \,|\, \lambda)$ over $\{0,1\}$ only one parameter $\lambda = P(y = 1 \,|\, x^{(i)})$

# Maximum Likelihood

Training set $\{x^{(i)}, y^{(i)}\}_{i=1,\dots,n}$ of input/output pairs

1. $y^{(i)}$ drawn from an unknown distribution $p_{data}(y \mid x^{(i)})$

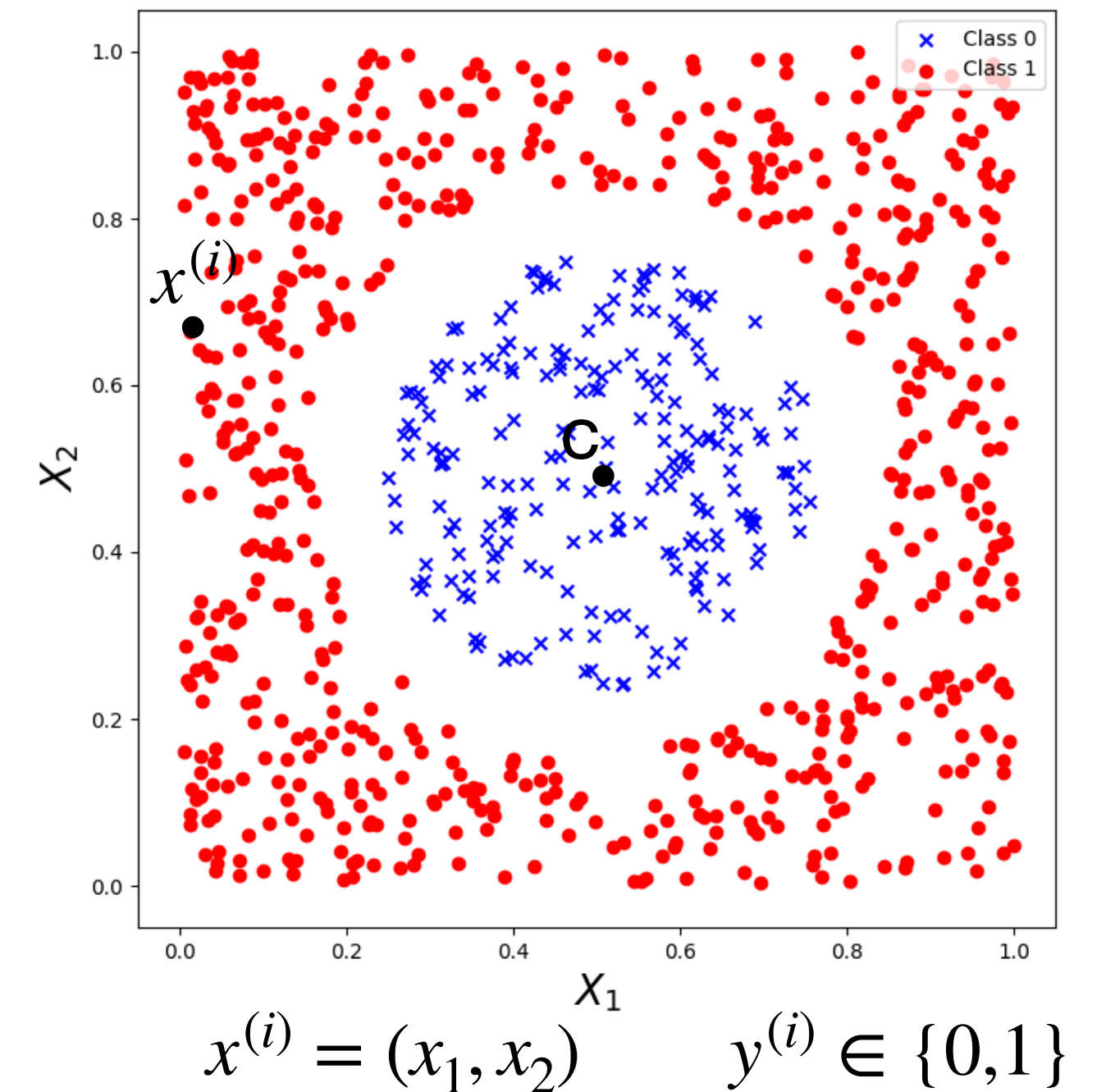2. Learn an approximation $p_{model}$ of the data distribution

3. Distribution $p_{model}(y \mid \lambda)$ over the output space
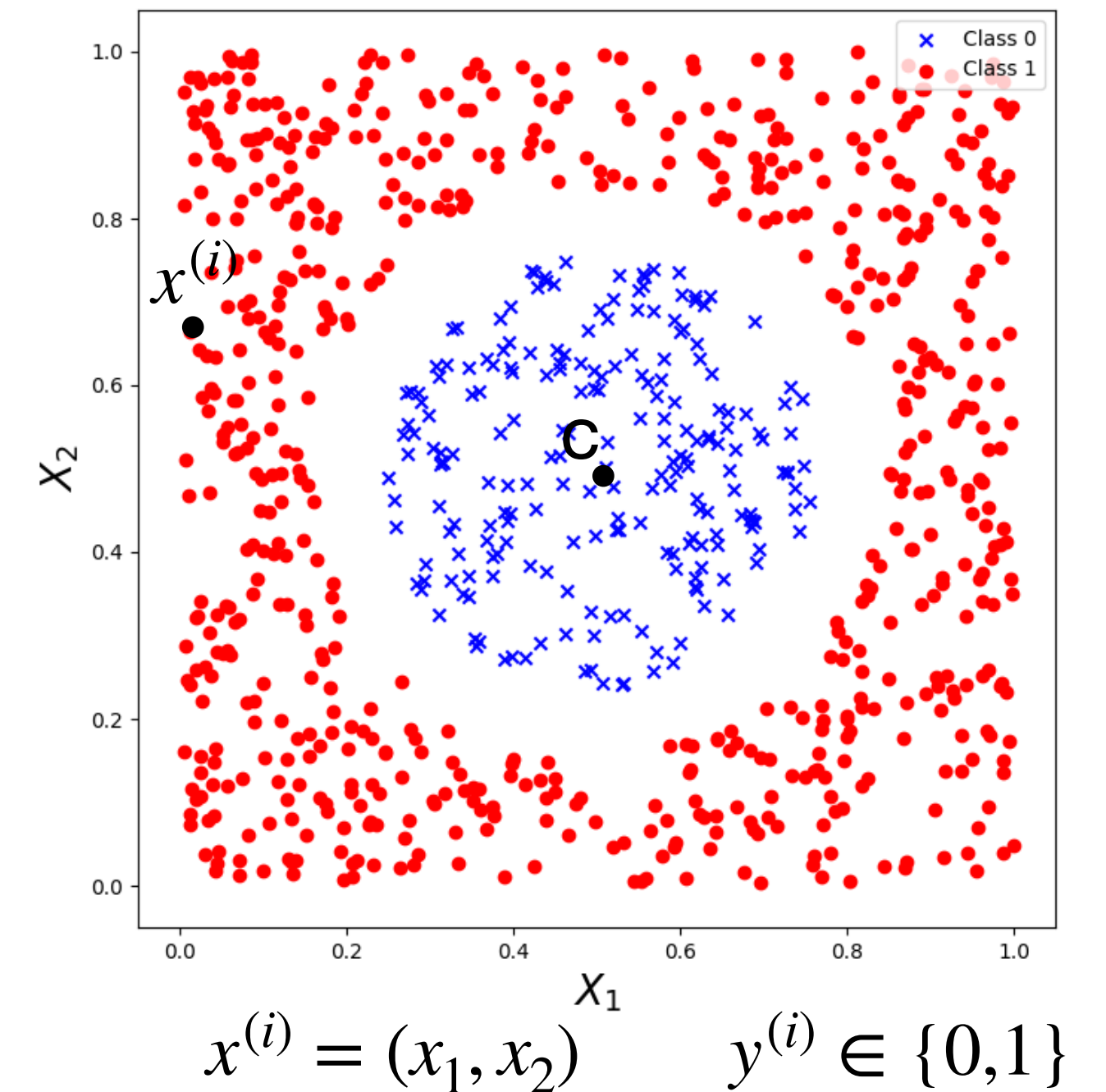


$$x^{(i)} = (x_1, x_2) \qquad y^{(i)} \in \{0,1\}$$

$$p_{model} \approx p_{data}(y \mid x^{(i)}) = \begin{cases} \delta_0(y) & dist(x^{(i)}, c) \leq r \\ \delta_1(y) & dist(x^{(i)}, c) > r \end{cases}$$

$p_{model}(y \mid \lambda)$ over $\{0,1\}$ only one parameter $\boxed{\lambda^{(i)}} = P(y = 1 \mid x^{(i)})$

$\lambda$ depends on the input x

13

# Maximum Likelihood

Training set $\{x^{(i)}, y^{(i)}\}_{i=1,\ldots,n}$ of input/output pairs

1. $y^{(i)}$ drawn from an unknown distribution $p_{data}(y \mid x^{(i)})$

2. Learn an approximation $p_{model}$ of the data distribution

3. Distribution $p_{model}(y \mid \lambda)$ over the output space

4. Set the network $f(x; \theta)$ to output (some of) the parameters $\lambda$



$$p_{model} \approx p_{data}(y \mid x^{(i)}) = \begin{cases} \delta_0(y) & dist(x^{(i)}, c) \leq r \\ \delta_1(y) & dist(x^{(i)}, c) > r \end{cases}$$
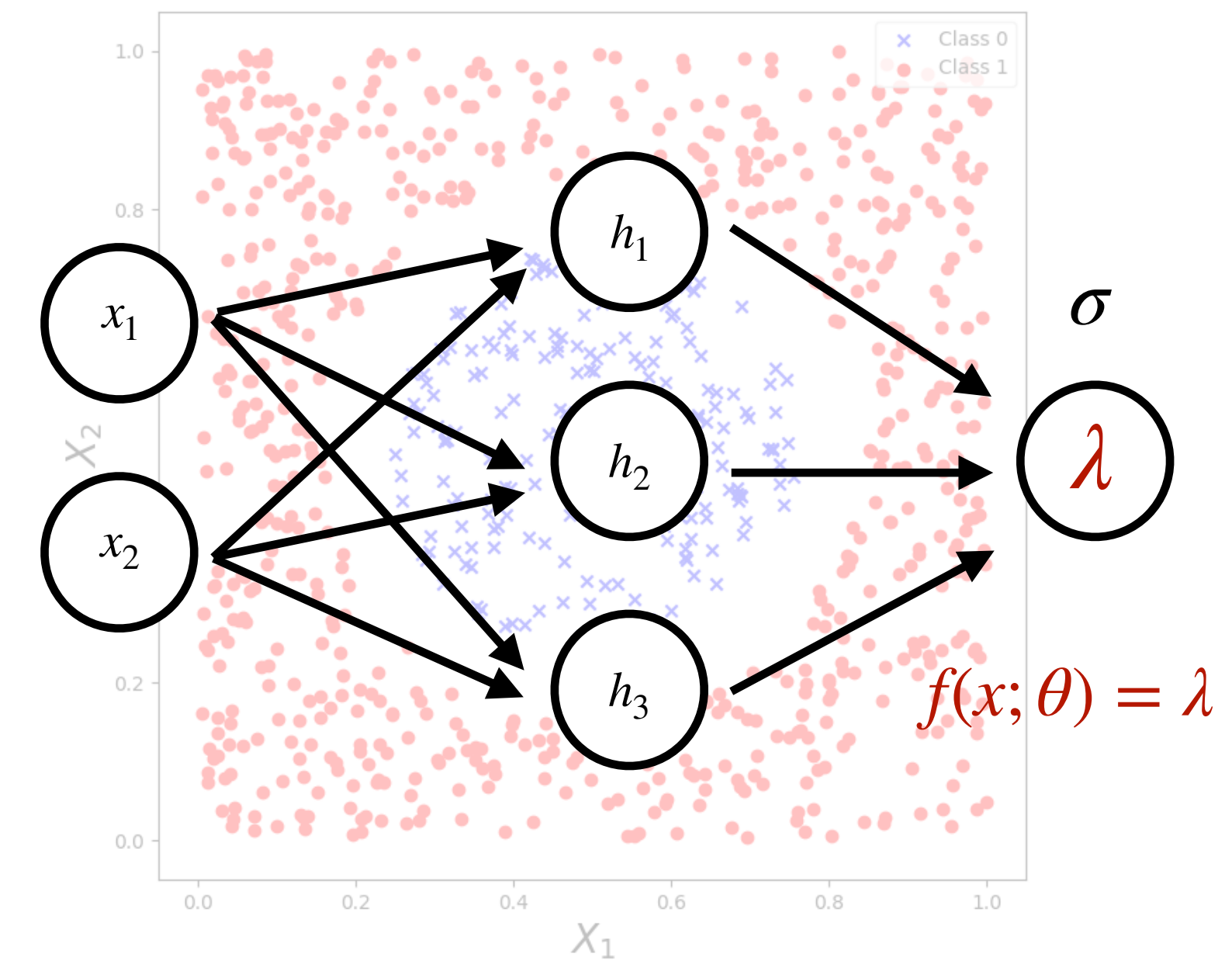
$p_{model}(y \mid \lambda)$ over $\{0,1\}$ only one parameter $\lambda^{(i)} = P(y = 1 \mid x^{(i)})$

# Maximum Likelihood

Training set $\{x^{(i)}, y^{(i)}\}_{i=1,\ldots,n}$ of input/output pairs

1. $y^{(i)}$ drawn from an unknown distribution $p_{data}(y \mid x^{(i)})$

2. Learn an approximation $p_{model}$ of the data distribution

3. Distribution $p_{model}(y \mid \lambda)$ over the output space

4. Set the network $f(x; \theta)$ to output (some of) the parameters $\lambda$

5. Look for parameters $\hat{\theta}$ to maximize the total likelihood of the training data

$$P(y^{(1)}, \ldots, y^{(n)} \mid x^{(1)}, \ldots, x^{(n)}; \theta)$$



$$p_{model} \approx p_{data}(y \mid x^{(i)}) = \begin{cases} \delta_0(y) & dist(x^{(i)}, c) \leq r \\ \delta_1(y) & dist(x^{(i)}, c) > r \end{cases}$$

$p_{model}(y \mid \lambda)$ over $\{0,1\}$ only one parameter $\lambda^{(i)} = P(y = 1 \mid x^{(i)})$
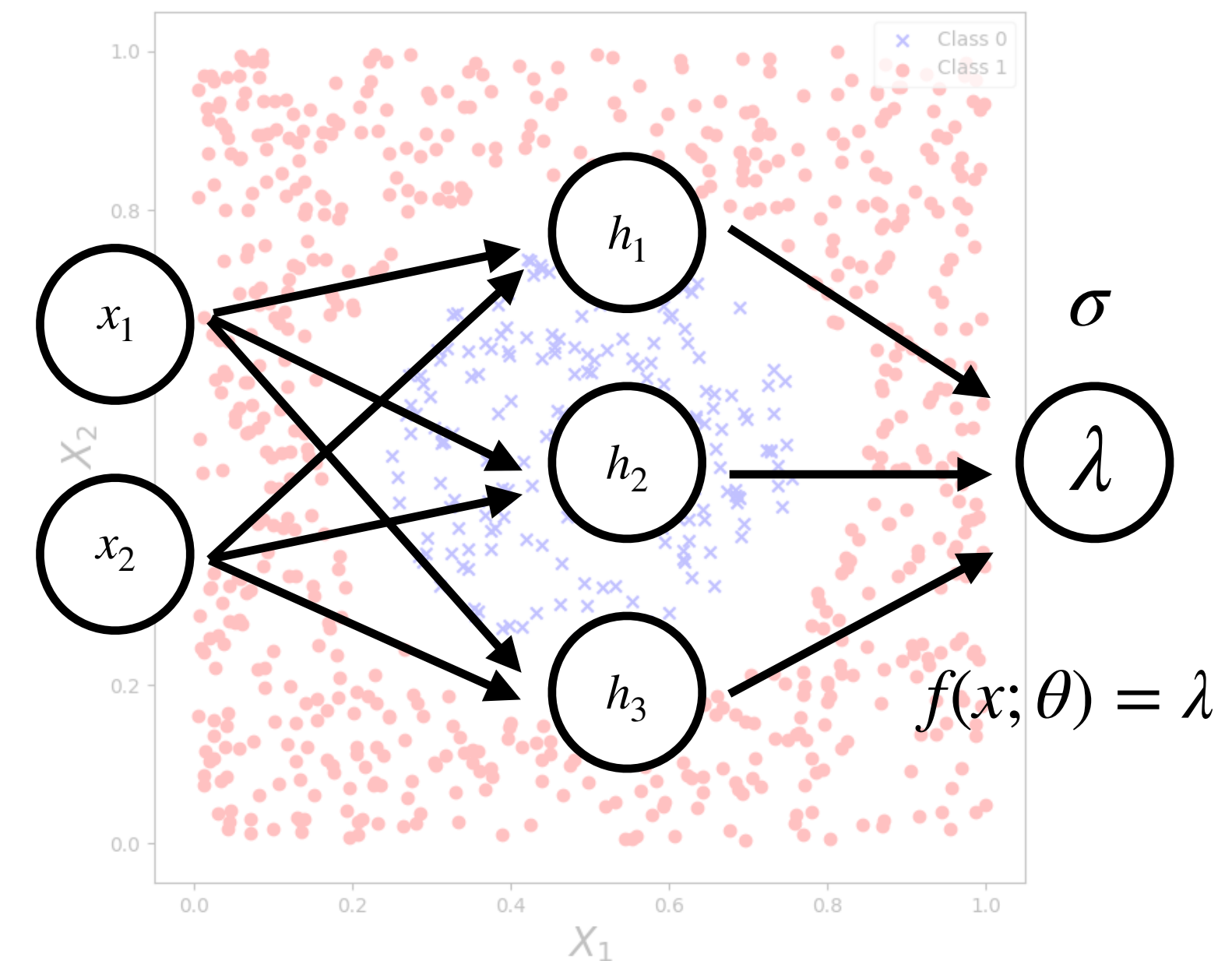
# Maximum Likelihood

Training set $\{x^{(i)}, y^{(i)}\}_{i=1,\ldots,n}$ of input/output pairs

1. $y^{(i)}$ drawn from an unknown distribution $p_{data}(y \mid x^{(i)})$

2. Learn an approximation $p_{model}$ of the data distribution

3. Distribution $p_{model}(y \mid \lambda)$ over the output space

4. Set the network $f(x; \theta)$ to output (some of) the parameters $\lambda$

5. Look for parameters $\hat{\theta}$ to maximize the total likelihood of the training data

$$P(y^{(1)}, \ldots, y^{(n)} \mid x^{(1)}, \ldots, x^{(n)}; \theta)$$

$$0, 1, 1, 1\ldots \quad (0.5,0.5), (0.2,0.2), (0.8,0.8), (0.8,0.1)\ldots$$



$$p_{model} \approx p_{data}(y \mid x^{(i)}) = \begin{cases} \delta_0(y) & dist(x^{(i)}, c) \leq r \\ \delta_1(y) & dist(x^{(i)}, c) > r \end{cases}$$

$p_{model}(y \mid \lambda)$ over $\{0,1\}$ only one parameter $\lambda^{(i)} = P(y = 1 \mid x^{(i)})$

# Maximum Likelihood

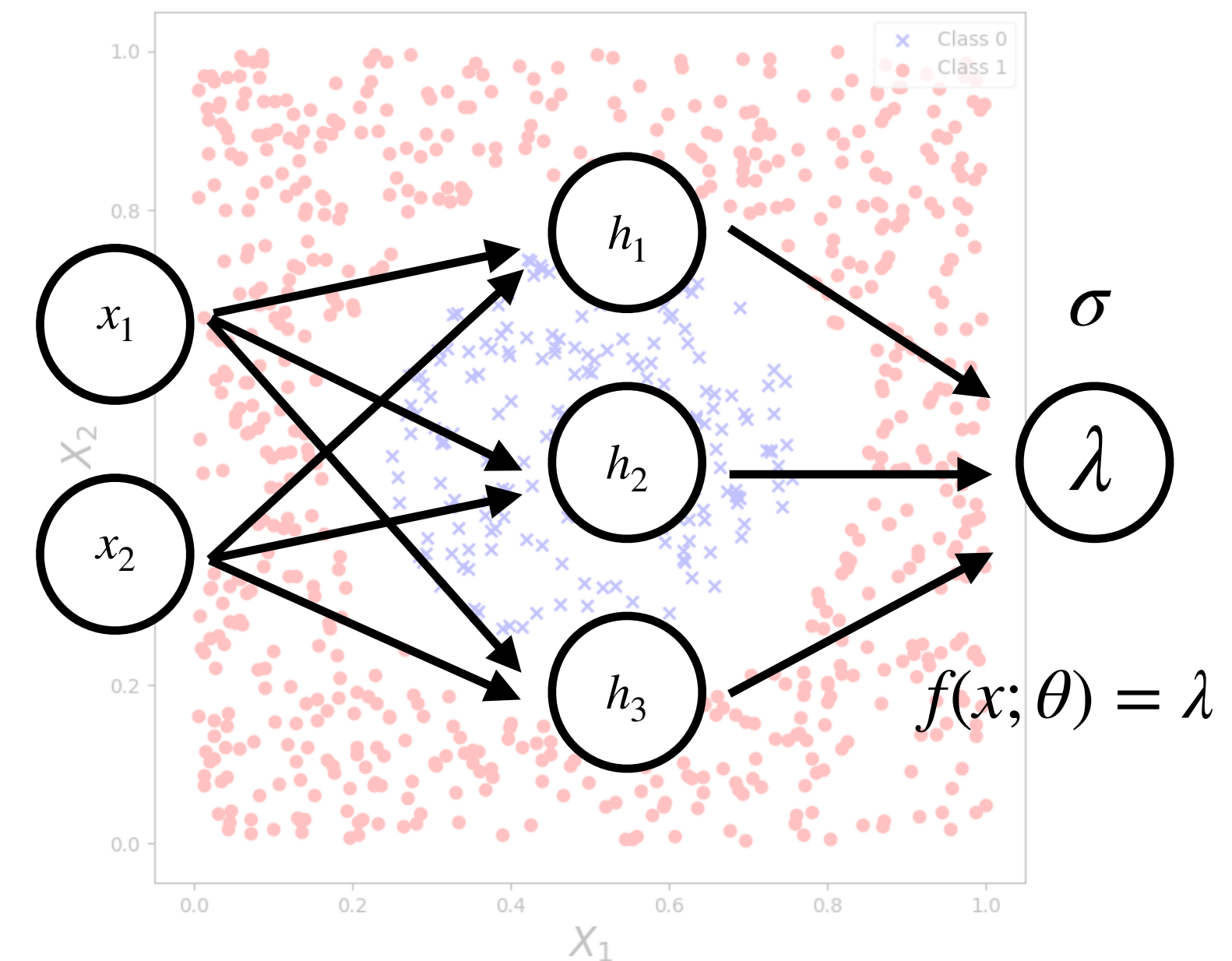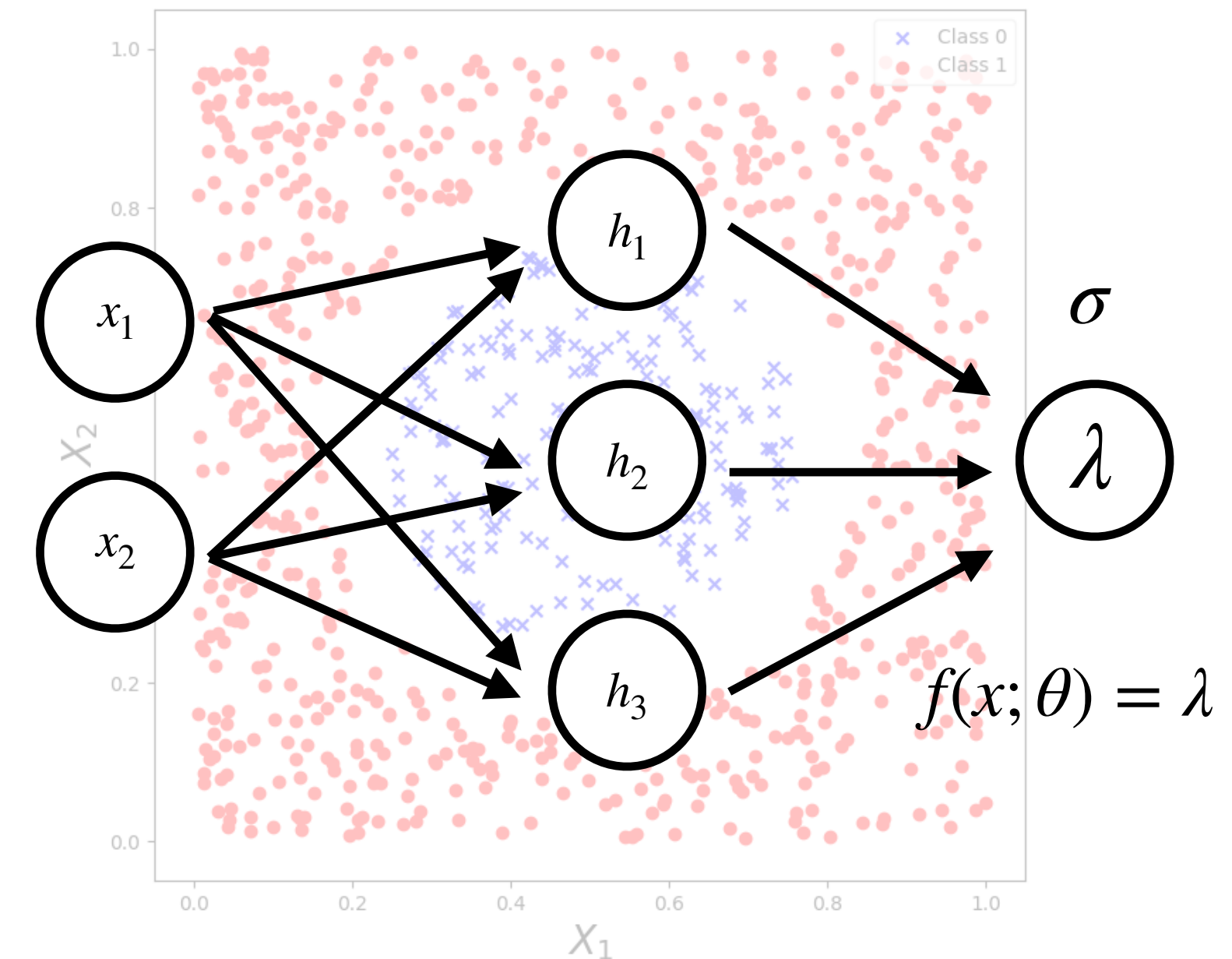Training set $\{x^{(i)}, y^{(i)}\}_{i=1,\ldots,n}$ of input/output pairs

1. $y^{(i)}$ drawn from an unknown distribution $p_{data}(y \mid x^{(i)})$

2. Learn an approximation $p_{model}$ of the data distribution

3. Distribution $p_{model}(y \mid \lambda)$ over the output space

4. Set the network $f(x; \theta)$ to output (some of) the parameters $\lambda$

5. Look for parameters $\hat{\theta}$ to maximize the total likelihood of the training data

$$P(y^{(1)}, \ldots, y^{(n)} \mid x^{(1)}, \ldots, x^{(n)}; \theta)$$

$$P(0, 1, 1, 1\ldots \mid (0.5,0.5), (0.2,0.2), (0.8,0.8), (0.8,0.1)\ldots; \theta)$$



$$p_{model} \approx p_{data}(y \mid x^{(i)}) = \begin{cases} \delta_0(y) & dist(x^{(i)}, c) \leq r \\ \delta_1(y) & dist(x^{(i)}, c) > r \end{cases}$$

$p_{model}(y \mid \lambda)$ over $\{0,1\}$ only one parameter $\lambda^{(i)} = P(y = 1 \mid x^{(i)})$

# Maximum Likelihood

Training set $\{x^{(i)}, y^{(i)}\}_{i=1,\ldots,n}$ of input/output pairs

1. $y^{(i)}$ drawn from an unknown distribution $p_{data}(y \,|\, x^{(i)})$

2. Learn an approximation $p_{model}$ of the data distribution

3. Distribution $p_{model}(y \,|\, \lambda)$ over the output space

4. Set the network $f(x; \theta)$ to output (some of) the parameters $\lambda$

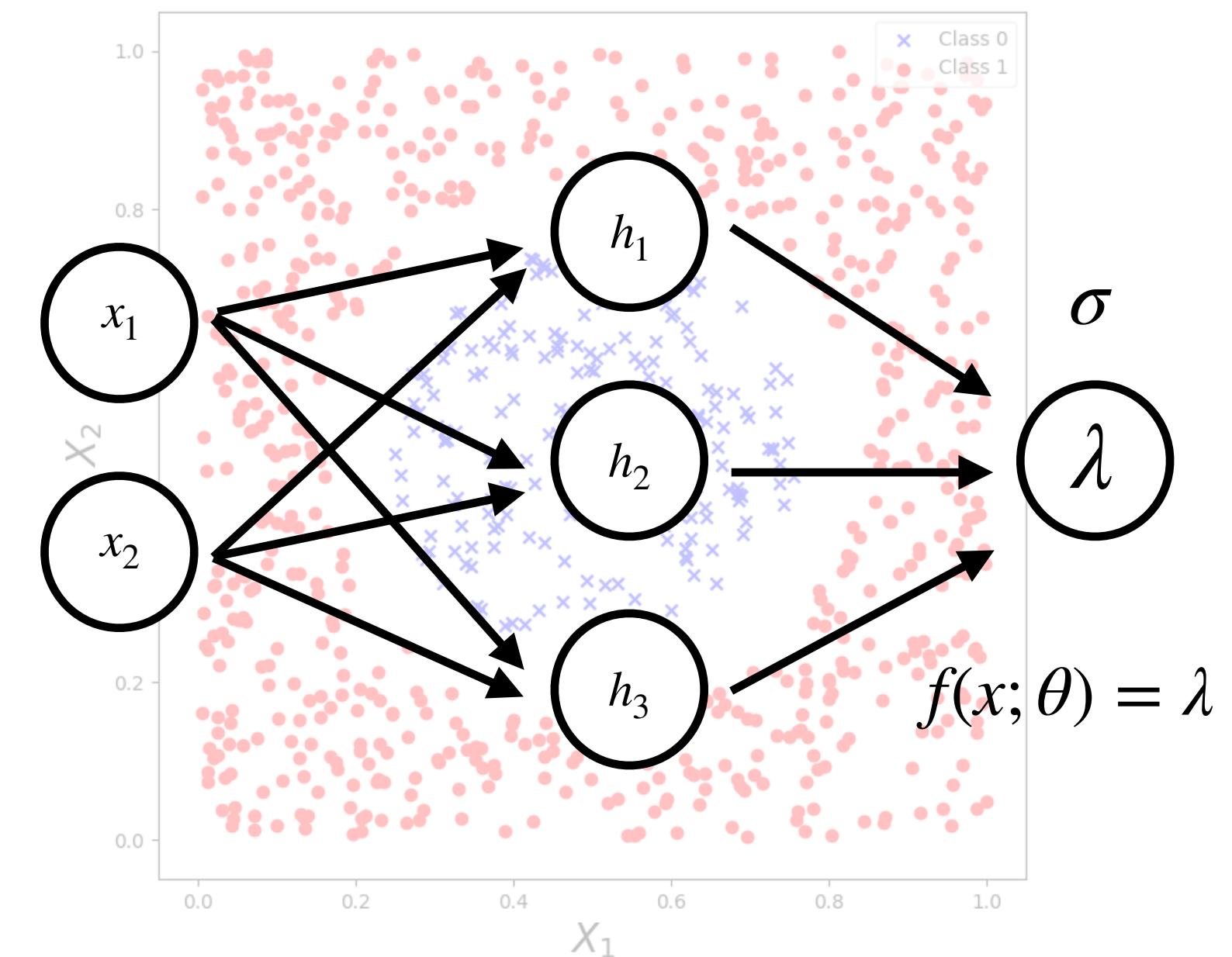5. Look for parameters $\hat{\theta}$ to maximize the total likelihood of the training data

$$P(y^{(1)}, \ldots, y^{(n)} \,|\, x^{(1)}, \ldots, x^{(n)}; \theta) = \prod_{i=1}^{n} P(y^{(i)} \,|\, x^{(i)}; \theta)$$

$\downarrow$

Multiplication $P(y^{(1)} | x^{(1)}; \theta) \times P(y^{(2)} | x^{(2)}; \theta) \times \ldots \times P(y^{(n)} | x^{(n)}; \theta)$



$$p_{model} \approx p_{data}(y \,|\, x^{(i)}) = \begin{cases} \delta_0(y) & dist(x^{(i)}, c) \leq r \\ \delta_1(y) & dist(x^{(i)}, c) > r \end{cases}$$

$p_{model}(y \,|\, \lambda)$ over $\{0,1\}$ only one parameter $\lambda^{(i)} = P(y = 1 \,|\, x^{(i)})$
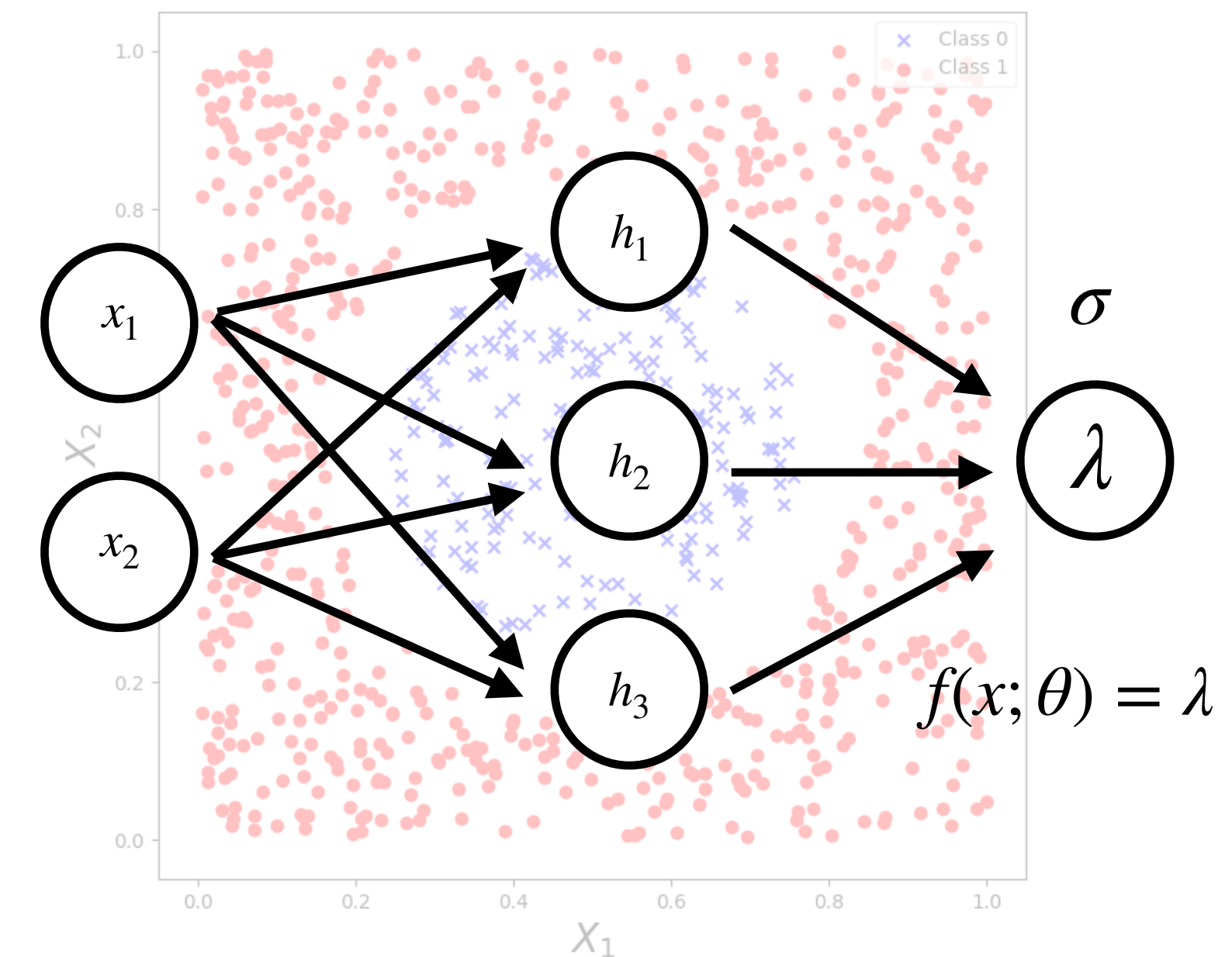
18

# Maximum Likelihood

Training set $\{x^{(i)}, y^{(i)}\}_{i=1,\ldots,n}$ of input/output pairs

1. $y^{(i)}$ drawn from an unknown distribution $p_{data}(y \mid x^{(i)})$

2. Learn an approximation $p_{model}$ of the data distribution

3. Distribution $p_{model}(y \mid \lambda)$ over the output space

4. Set the network $f(x; \theta)$ to output (some of) the parameters $\lambda$

5. Look for parameters $\hat{\theta}$ to maximize the total likelihood of the training data

$$P(y^{(1)}, \ldots, y^{(n)} \mid x^{(1)}, \ldots, x^{(n)}; \theta) = \prod_{i=1}^{n} P(y^{(i)} \mid x^{(i)}; \theta)$$

Can we do that?



$$p_{model} \approx p_{data}(y \mid x^{(i)}) = \begin{cases} \delta_0(y) & dist(x^{(i)}, c) \leq r \\ \delta_1(y) & dist(x^{(i)}, c) > r \end{cases}$$

$p_{model}(y \mid \lambda)$ over $\{0,1\}$ only one parameter $\lambda^{(i)} = P(y = 1 \mid x^{(i)})$

# Maximum Likelihood

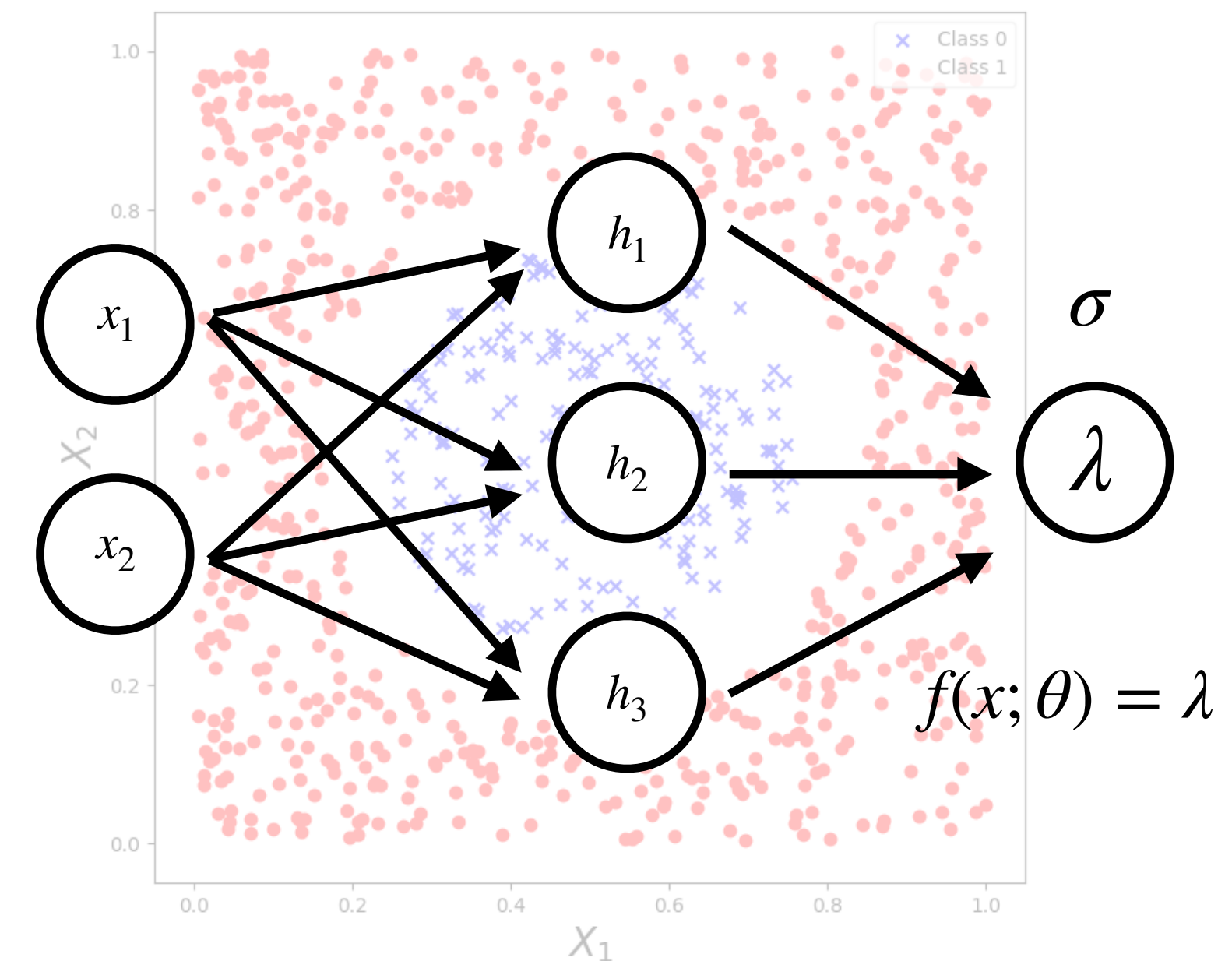Training set $\{x^{(i)}, y^{(i)}\}_{i=1,\ldots,n}$ of i.i.d. input/output pairs

1. $y^{(i)}$ drawn from an unknown distribution $p_{data}(y \mid x^{(i)})$

2. Learn an approximation $p_{model}$ of the data distribution

3. Distribution $p_{model}(y \mid \lambda)$ over the output space

4. Set the network $f(x; \theta)$ to output (some of) the parameters $\lambda$

5. Look for parameters $\hat{\theta}$ to maximize the total likelihood of the training data

$$P(y^{(1)}, \ldots, y^{(n)} \mid x^{(1)}, \ldots, x^{(n)}; \theta) = \prod_{i=1}^{n} P(y^{(i)} \mid x^{(i)}; \theta)$$



$$p_{model} \approx p_{data}(y \mid x^{(i)}) = \begin{cases} \delta_0(y) & dist(x^{(i)}, c) \leq r \\ \delta_1(y) & dist(x^{(i)}, c) > r \end{cases}$$

$p_{model}(y \mid \lambda)$ over $\{0,1\}$ only one parameter $\lambda^{(i)} = P(y = 1 \mid x^{(i)})$

Assuming $y$ conditional independent on $x$ and identically distributed
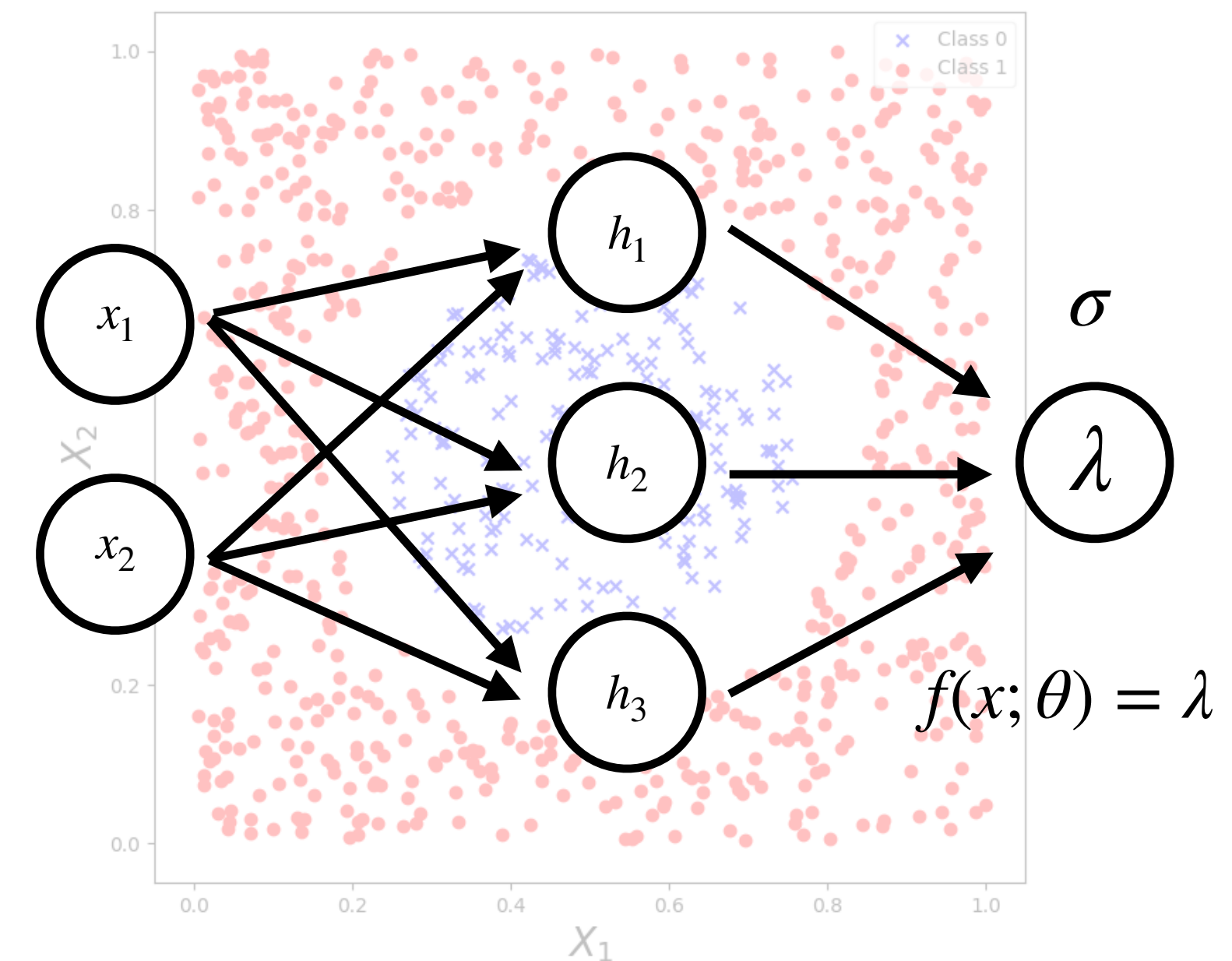
20

# Maximum Likelihood

Training set $\{x^{(i)}, y^{(i)}\}_{i=1,\ldots,n}$ of i.i.d. input/output pairs

1. $y^{(i)}$ drawn from an unknown distribution $p_{data}(y \mid x^{(i)})$

2. Learn an approximation $p_{model}$ of the data distribution

3. Distribution $p_{model}(y \mid \lambda)$ over the output space

4. Set the network $f(x; \theta)$ to output (some of) the parameters $\lambda$

5. Look for parameters $\hat{\theta}$ to maximize the total likelihood of the training data

$$P(y^{(1)}, \ldots, y^{(n)} \mid x^{(1)}, \ldots, x^{(n)}; \theta) = \prod_{i=1}^{n} p_{model}(y^{(i)} \mid f(x^{(i)}; \theta))$$
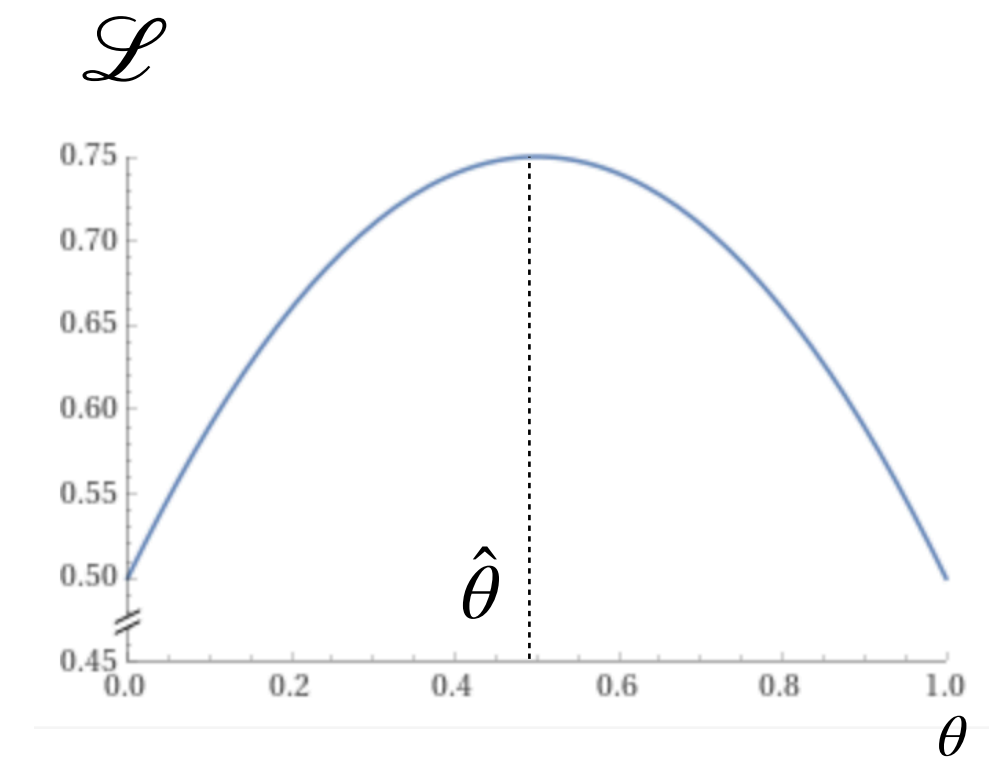


$$p_{model} \approx p_{data}(y \mid x^{(i)}) = \begin{cases} \delta_0(y) & dist(x^{(i)}, c) \leq r \\ \delta_1(y) & dist(x^{(i)}, c) > r \end{cases}$$

$p_{model}(y \mid \lambda)$ over $\{0,1\}$ only one parameter $\lambda^{(i)} = P(y = 1 \mid x^{(i)})$

# Maximum Likelihood

$$\text{likelihood}$$

**Maximum Likelihood criterion:** $\hat{\theta} = argmax_\theta \left( \overbrace{\prod_{i=1}^{n} p_{model}(y^{(i)} \mid f(x^{(i)}; \theta))}^{likelihood} \right)$

# Maximum Likelihood

$$\text{likelihood}$$

**Maximum Likelihood criterion:** $\hat{\theta} = argmax_\theta \left( \overbrace{\prod_{i=1}^{n} p_{model}(y^{(i)} \,|\, f(x^{(i)}; \theta))}^{likelihood} \right)$



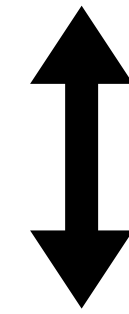Multiplying many probabilities (numbers lower than 1): the product can become very small!

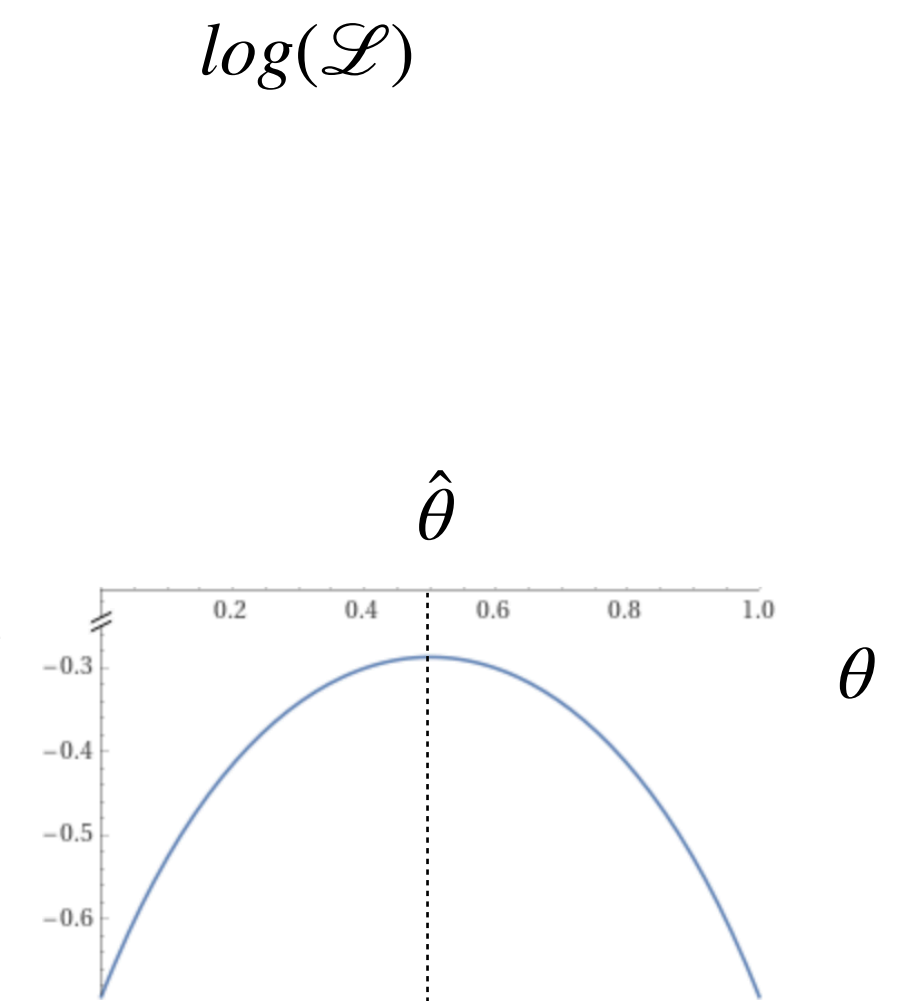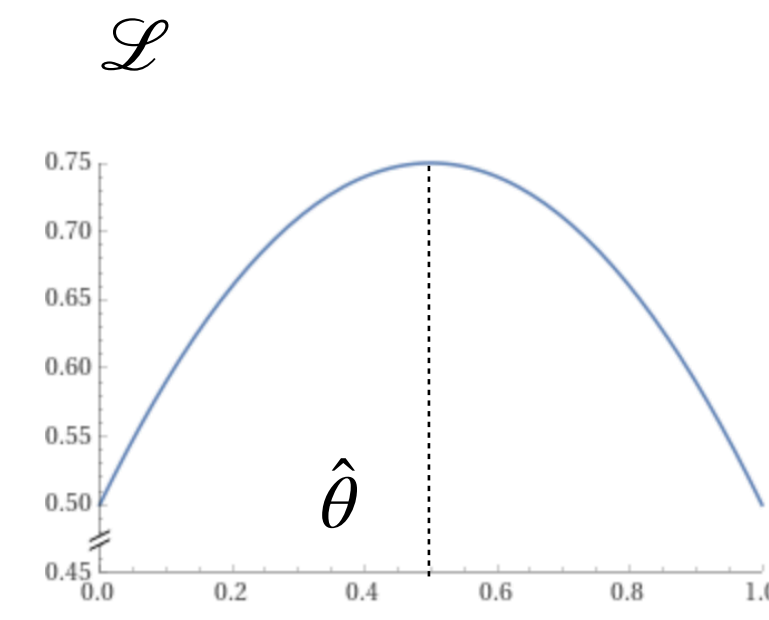$$0.01 \times 0.1 \times \ldots \approx 0$$

# Maximum Likelihood

**Maximum Likelihood criterion:** $\hat{\theta} = argmax_\theta \left( \prod_{i=1}^{n} p_{model}(y^{(i)} \mid f(x^{(i)}; \theta)) \right)$

*log- likelihood*

**Maximum log-Likelihood criterion:** $\hat{\theta} = argmax_\theta \left( \log \left( \prod_{i=1}^{n} p_{model}(y^{(i)} \mid f(x^{(i)}; \theta)) \right) \right)$



$\mathcal{L}$

$log(\mathcal{L})$

$\hat{\theta}$

$\hat{\theta}$

$\theta$

Log is an increasing function

# Maximum Likelihood

**Maximum Likelihood criterion:** $\hat{\theta} = argmax_\theta \left( \prod_{i=1}^{n} p_{model}(y^{(i)} \mid f(x^{(i)}; \theta)) \right)$

$\updownarrow$

*log- likelihood*

**Maximum log-Likelihood criterion:** $\hat{\theta} = argmax_\theta \left( \log \left( \overbrace{\prod_{i=1}^{n} p_{model}(y^{(i)} \mid f(x^{(i)}; \theta))} \right) \right) = argmax_\theta \left( \sum_{i=1}^{n} \log(p_{model}(y^{(i)} \mid f(x^{(i)}; \theta))) \right)$

✅ Sum of log probabilities!

# Maximum Likelihood

**Maximum Likelihood criterion:** $\hat{\theta} = argmax_\theta \left( \prod_{i=1}^{n} p_{model}(y^{(i)} \mid f(x^{(i)}; \theta)) \right)$

$\updownarrow$

**Maximum log-Likelihood criterion:** $\hat{\theta} = argmax_\theta \left( \sum_{i=1}^{n} \log(p_{model}(y^{(i)} \mid f(x^{(i)}; \theta))) \right)$

$\updownarrow$

*negative log-likelihood*

**Minimum negative log-Likelihood criterion:** $\hat{\theta} = argmin_\theta \left( \overbrace{\sum_{i=1}^{n} - \log(p_{model}(y^{(i)} \mid f(x^{(i)}; \theta)))} \right)$

$log(\mathscr{L})$

$-log(\mathscr{L})$

$\hat{\theta}$

$\hat{\theta}$

$\theta$

$argmin(-\mathscr{L}) = argmax(\mathscr{L})$

# Maximum Likelihood

**Maximum Likelihood criterion:** $\hat{\theta} = argmax_{\theta} \left( \prod_{i=1}^{n} p_{model}(y^{(i)} \,|\, f(x^{(i)}; \theta)) \right)$

$\updownarrow$

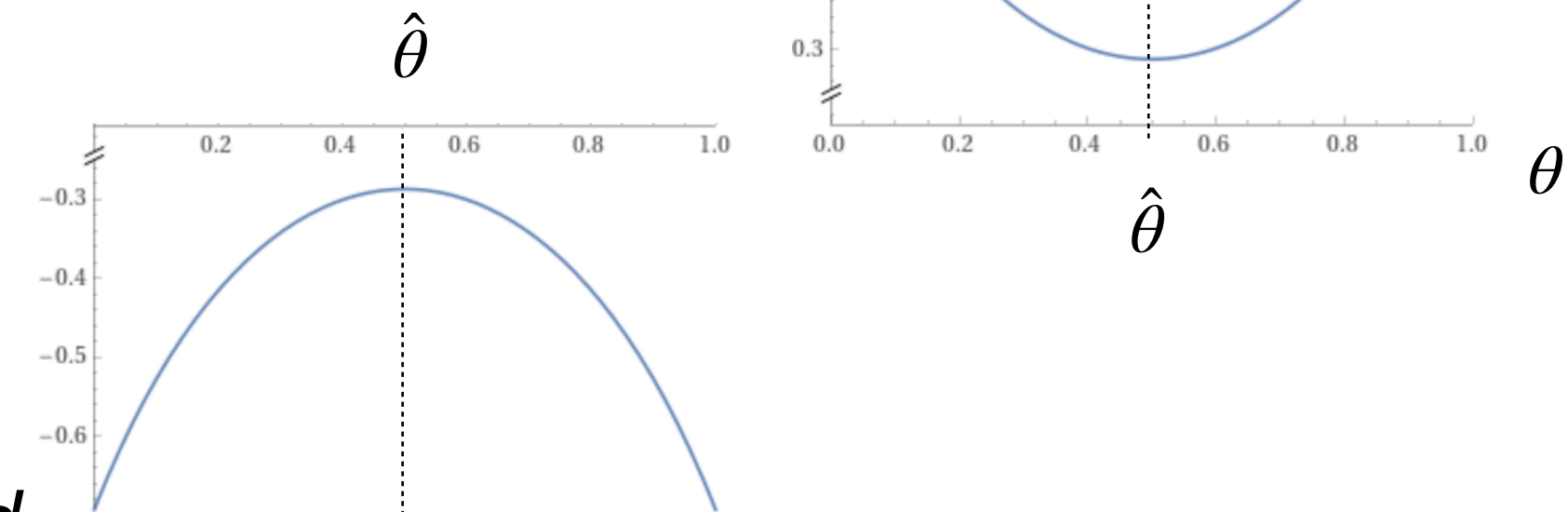**Maximum log-Likelihood criterion:** $\hat{\theta} = argmax_{\theta} \left( \sum_{i=1}^{n} \log(p_{model}(y^{(i)} \,|\, f(x^{(i)}; \theta))) \right)$

$\updownarrow$

**Minimum negative log-Likelihood criterion:** $\hat{\theta} = argmin_{\theta} \left( \sum_{i=1}^{n} -\log(p_{model}(y^{(i)} \,|\, f(x^{(i)}; \theta))) \right)$

$\updownarrow$
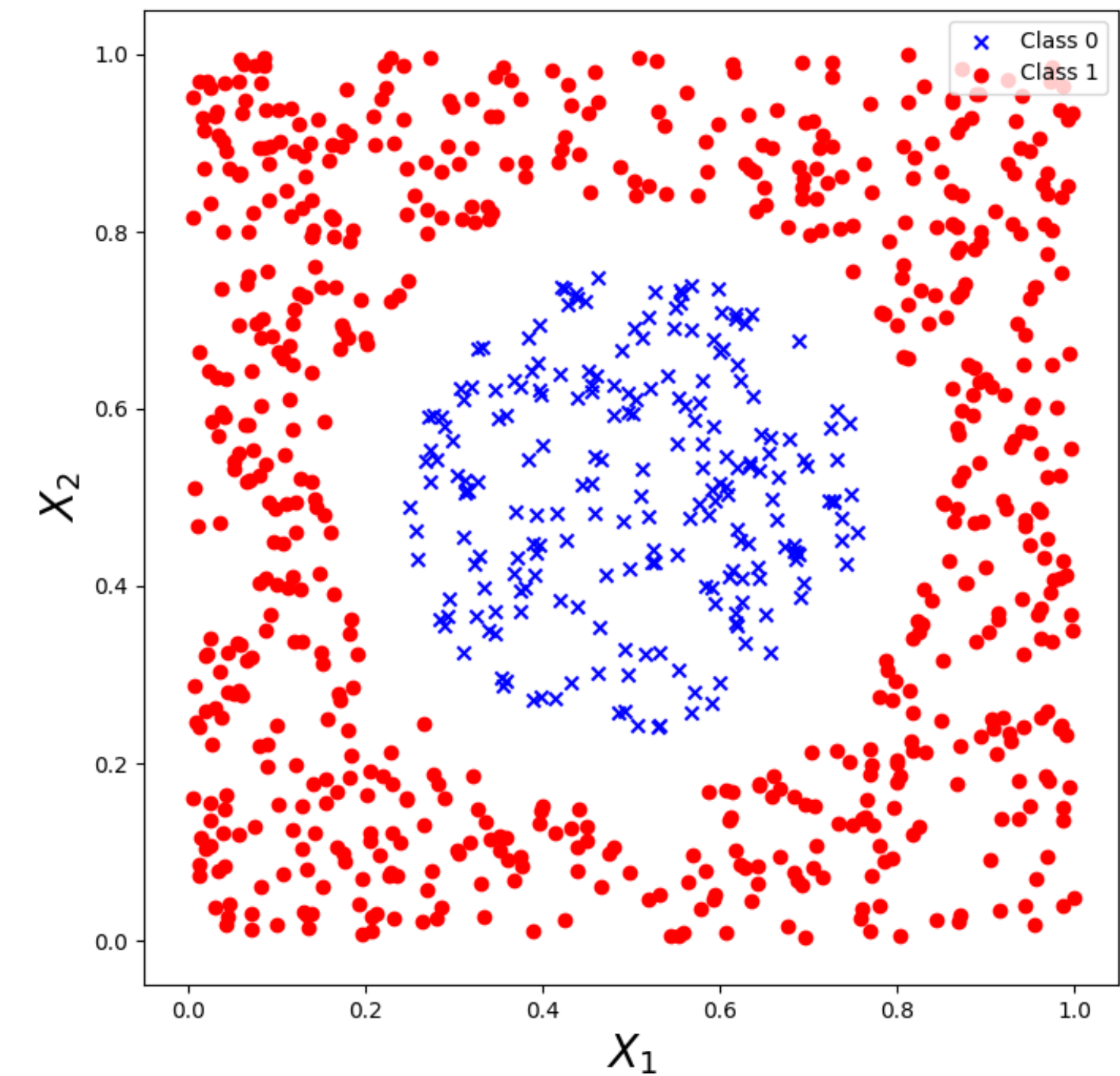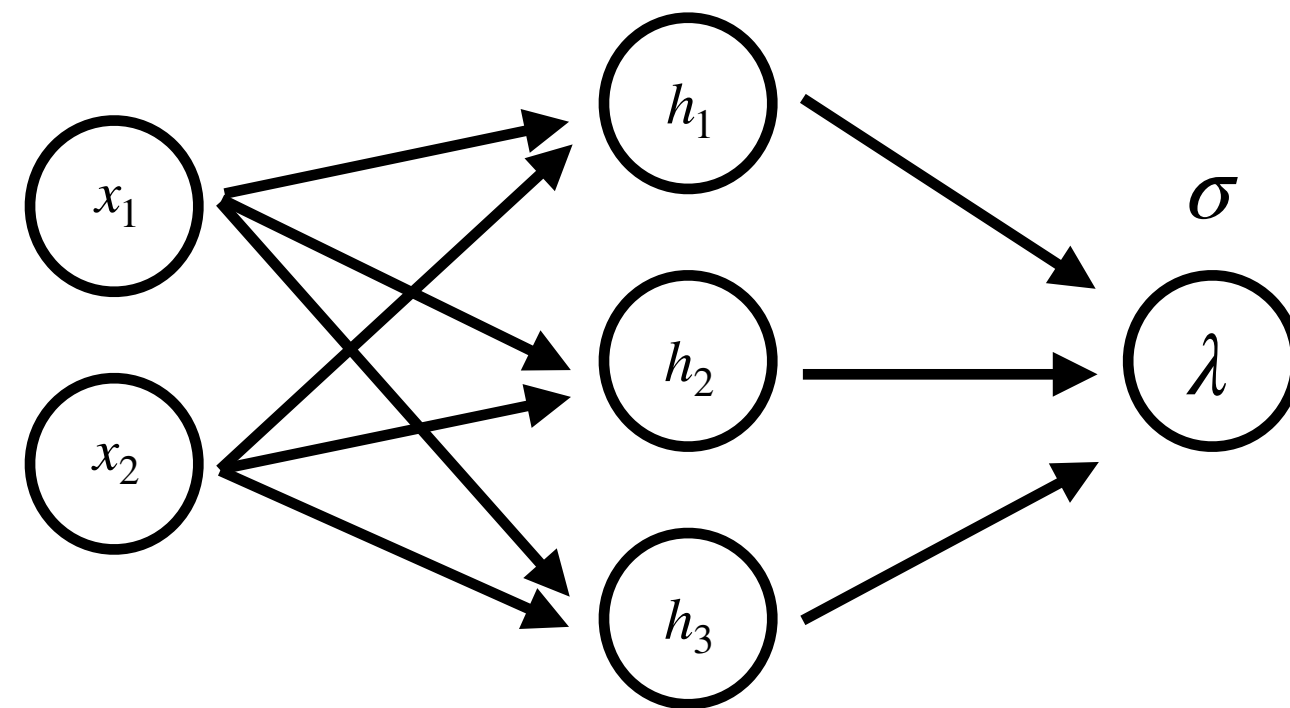
*Loss*

**Minimum loss function:** $\hat{\theta} = argmin_{\theta} \left( \frac{1}{n} \sum_{i=1}^{n} -\log(p_{model}(y^{(i)} \,|\, f(x^{(i)}; \theta))) \right) = argmin_{\theta} L(\theta)$

# Building the Loss Function: Binary Classification

Training dataset $\{x^{(i)}, y^{(i)}\}_{i=1}^{n}$

$f(x^{(i)}; \theta) = \lambda^{(i)}$ probability that an input $x^{(i)}$ belongs to class 1



**Loss function:** $L(\theta) = \dfrac{1}{n} \displaystyle\sum_{i=1}^{n} -\log(p_{model}(y^{(i)} \,|\, f(x^{(i)}; \theta)))$

# Building the Loss Function: Binary Classification

Training dataset $\{x^{(i)}, y^{(i)}\}_{i=1}^{n}$

$f(x^{(i)}; \theta) = \lambda^{(i)}$ probability that an input $x^{(i)}$ belongs to class 1



**Loss function:** $L(\theta) = \dfrac{1}{n} \displaystyle\sum_{i=1}^{n} - \log \boxed{p_{model}(y^{(i)} \mid f(x^{(i)}; \theta))}$ ⟶ $p_{model}(y^{(i)} \mid f(x^{(i)}; \theta)) = \begin{cases} \lambda^{(i)} & \text{if } y^{(i)} = 1 \\ 1 - \lambda^{(i)} & \text{if } y^{(i)} = 0 \end{cases}$

29

# Building the Loss Function: Binary Classification

Training dataset $\{x^{(i)}, y^{(i)}\}_{i=1}^{n}$

$f(x^{(i)}; \theta) = \lambda^{(i)}$ probability that an input $x^{(i)}$ belongs to class 1
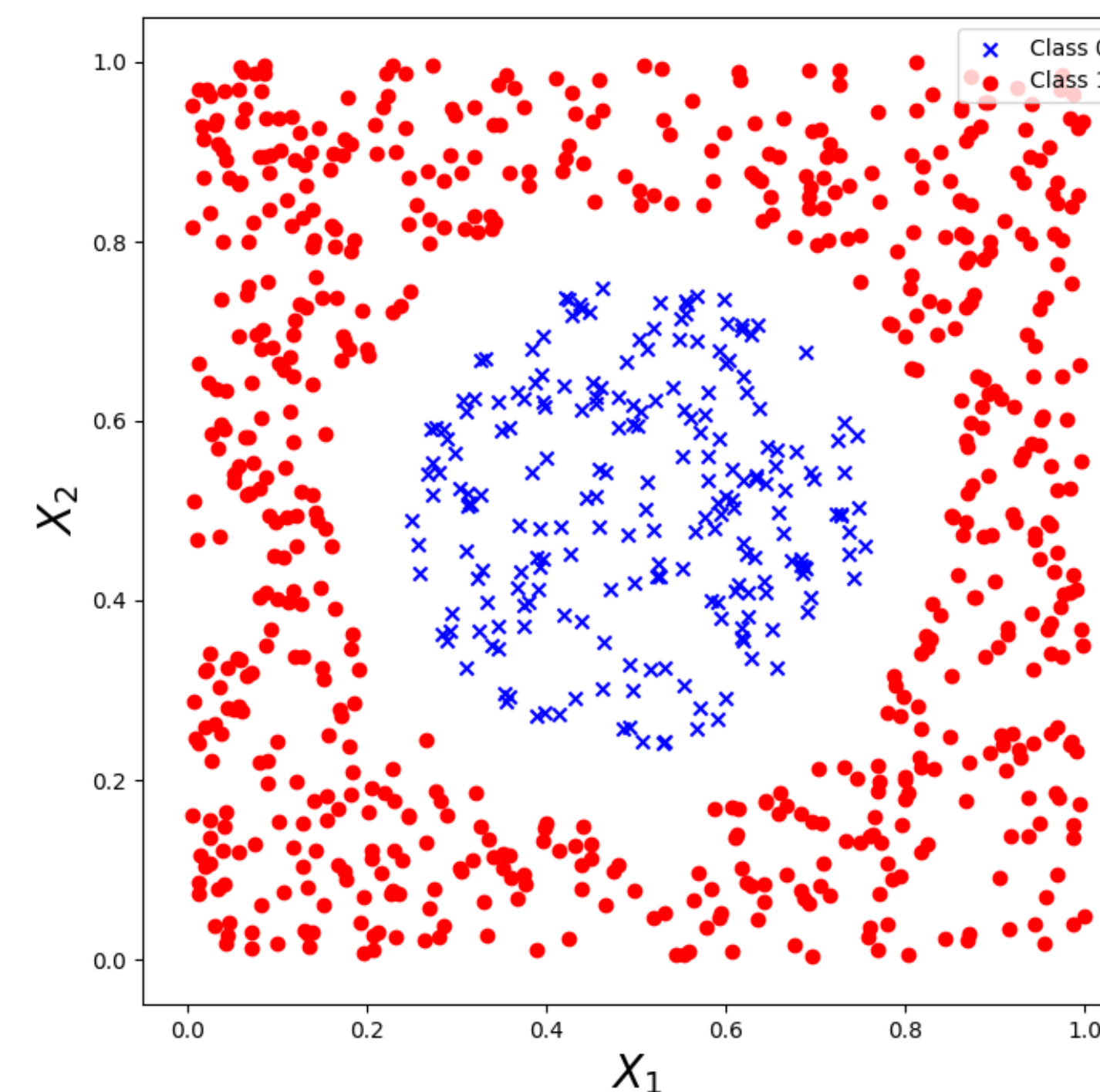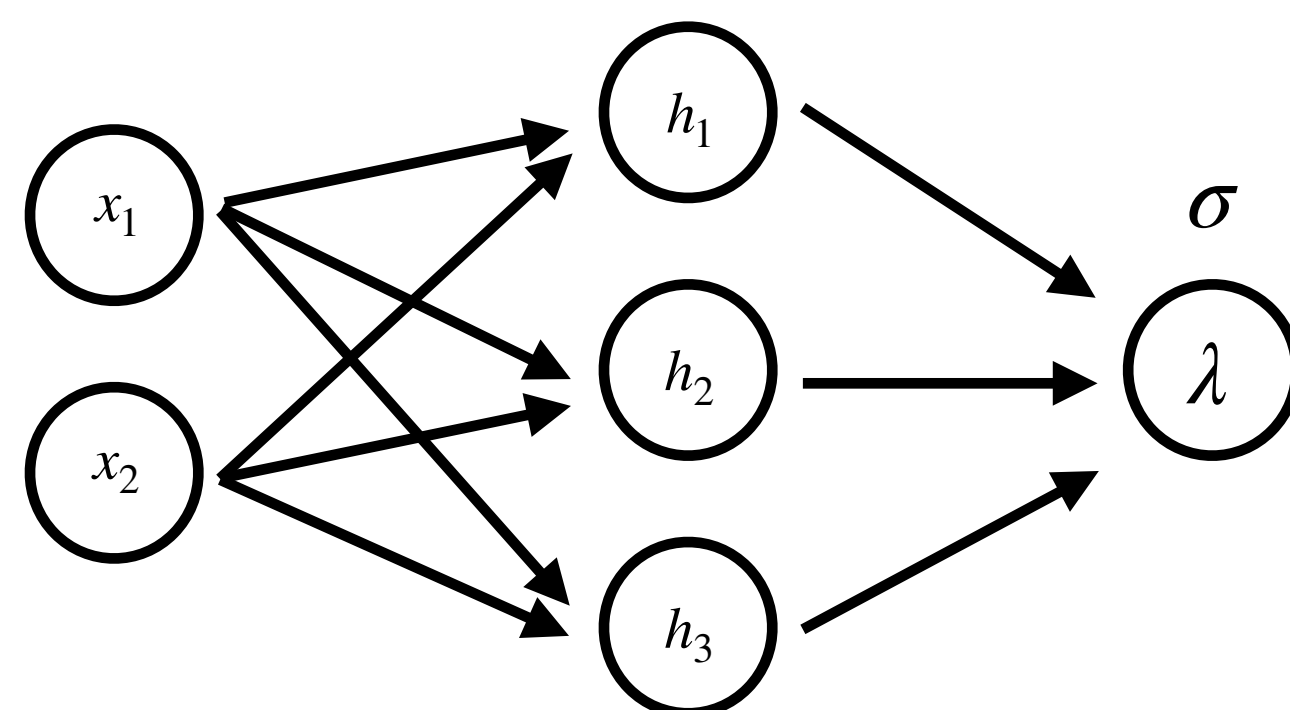


**Loss function:** $L(\theta) = \dfrac{1}{n} \sum\limits_{i=1}^{n} -\log(p_{model}(y^{(i)} \mid f(x^{(i)}; \theta)))$

| Input x | Class y | Output $\lambda$ | Error $l$ | Predicted $\hat{y}$ |
|---------|---------|--------|---------|-----------|
| (0.5, 0.5) | 0 | **0.1** | **?** | **?** |

# Building the Loss Function: Binary Classification

Training dataset $\{x^{(i)}, y^{(i)}\}_{i=1}^{n}$

$f(x^{(i)}; \theta) = \lambda^{(i)}$ probability that an input $x^{(i)}$ belongs to class 1
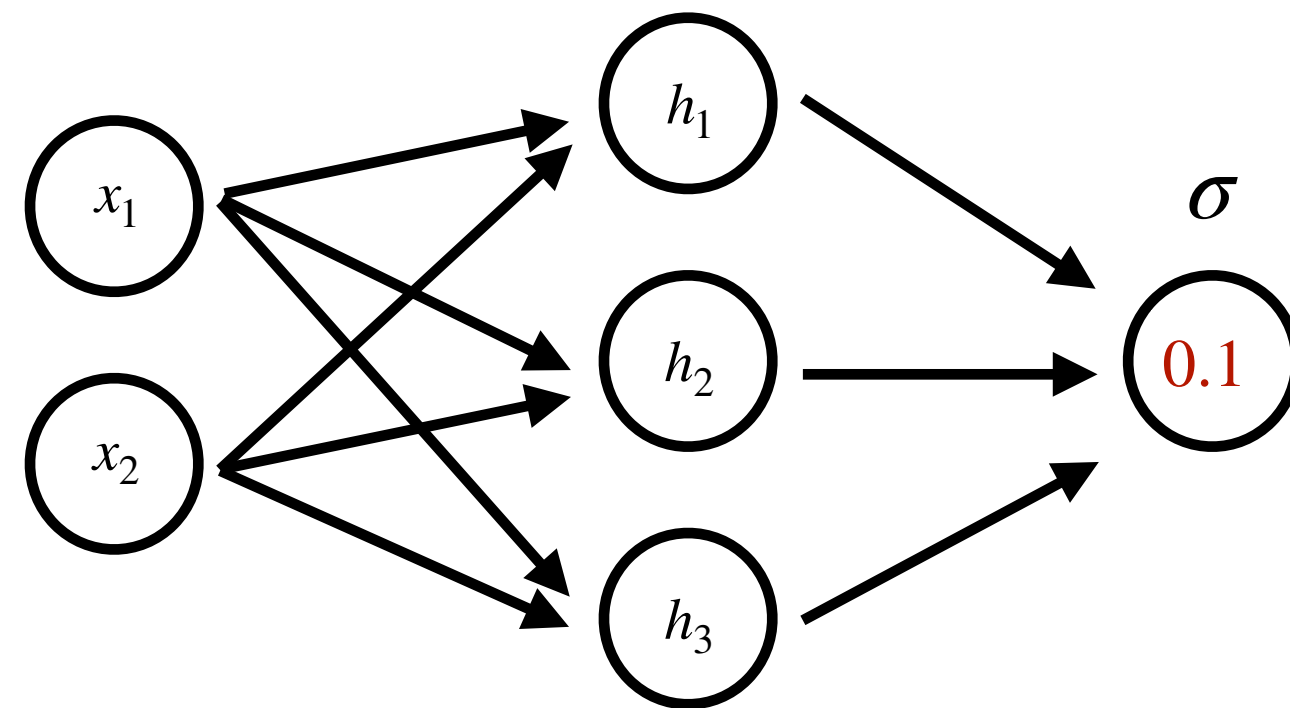


**Loss function:** $L(\theta) = \dfrac{1}{n} \sum_{i=1}^{n} \boxed{-\log(p_{model}(y^{(i)} \mid f(x^{(i)}; \theta)))}$

$l^{(i)} = -\log(p_{model}(y^{(i)} = 0 \mid \lambda^{(i)} = 0.1)) = -\log(0.9)$

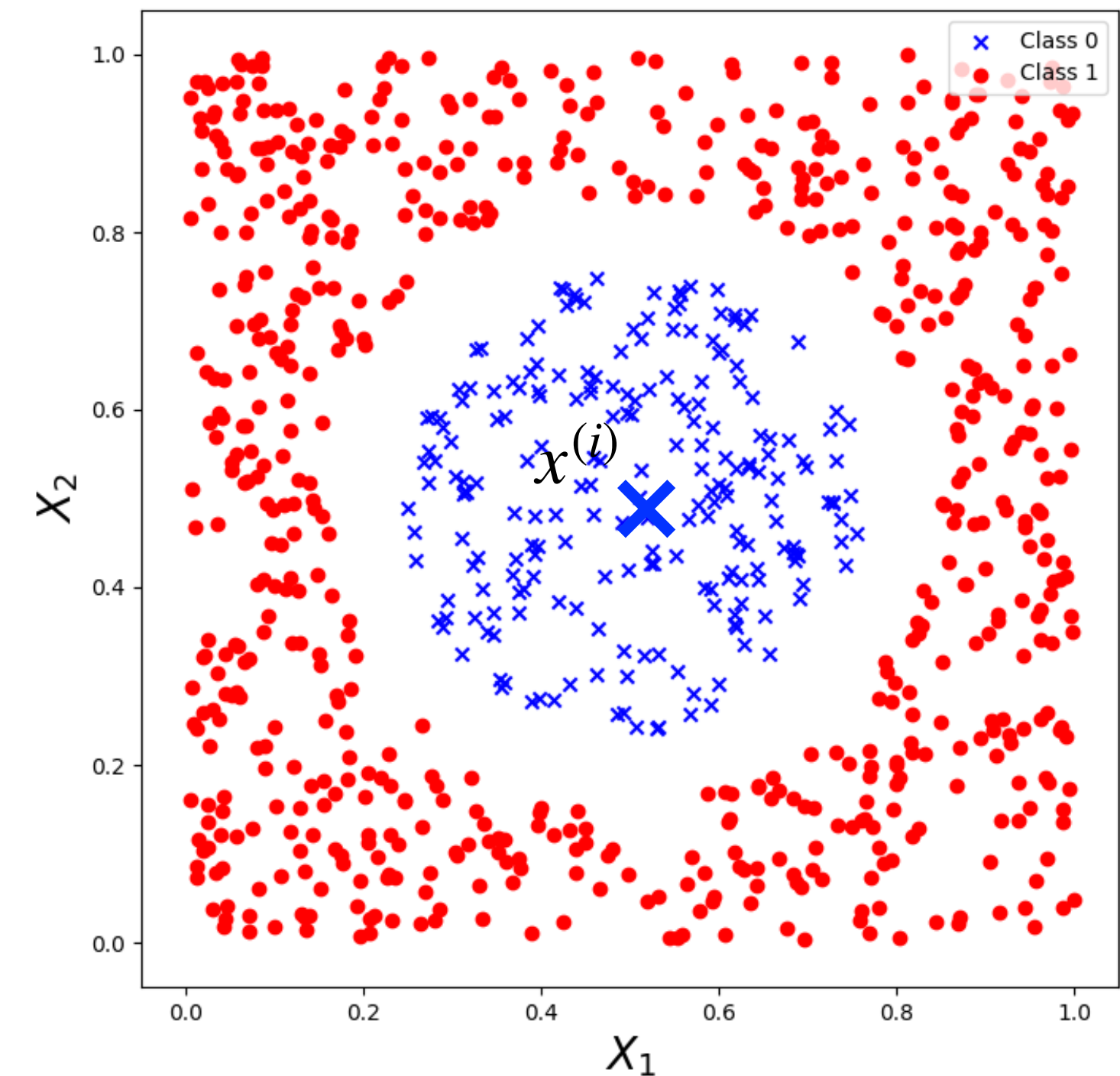| Input x | Class y | Output $\lambda$ | Error $l$ | Predicted $\hat{y}$ |
|---------|---------|------------------|-----------|---------------------|
| (0.5, 0.5) | 0 | 0.1 | **-log(0.9)** | **?** |

# Building the Loss Function: Binary Classification

Training dataset $\{x^{(i)}, y^{(i)}\}_{i=1}^{n}$

$f(x^{(i)}; \theta) = \lambda^{(i)}$ probability that an input $x^{(i)}$ belongs to class 1



**Loss function:** $L(\theta) = \dfrac{1}{n} \displaystyle\sum_{i=1}^{n} -\log(p_{model}(y^{(i)} \mid f(x^{(i)}; \theta)))$

$$l^{(i)} = -y^{(i)} \log(\lambda^{(i)}) - (1 - y^{(i)})\log(1 - \lambda^{(i)})$$

| Input x | Class y | Output $\lambda$ | Error $l$ | Predicted $\hat{y}$ |
|---------|---------|------------------|-----------|---------------------|
| (0.5, 0.5) | 0 | 0.1 | **-log(0.9)** | **?** |

32

# Building the Loss Function: Binary Classification

Training dataset $\{x^{(i)}, y^{(i)}\}_{i=1}^{n}$

$f(x^{(i)}; \theta) = \lambda^{(i)}$ probability that an input $x^{(i)}$ belongs to class 1



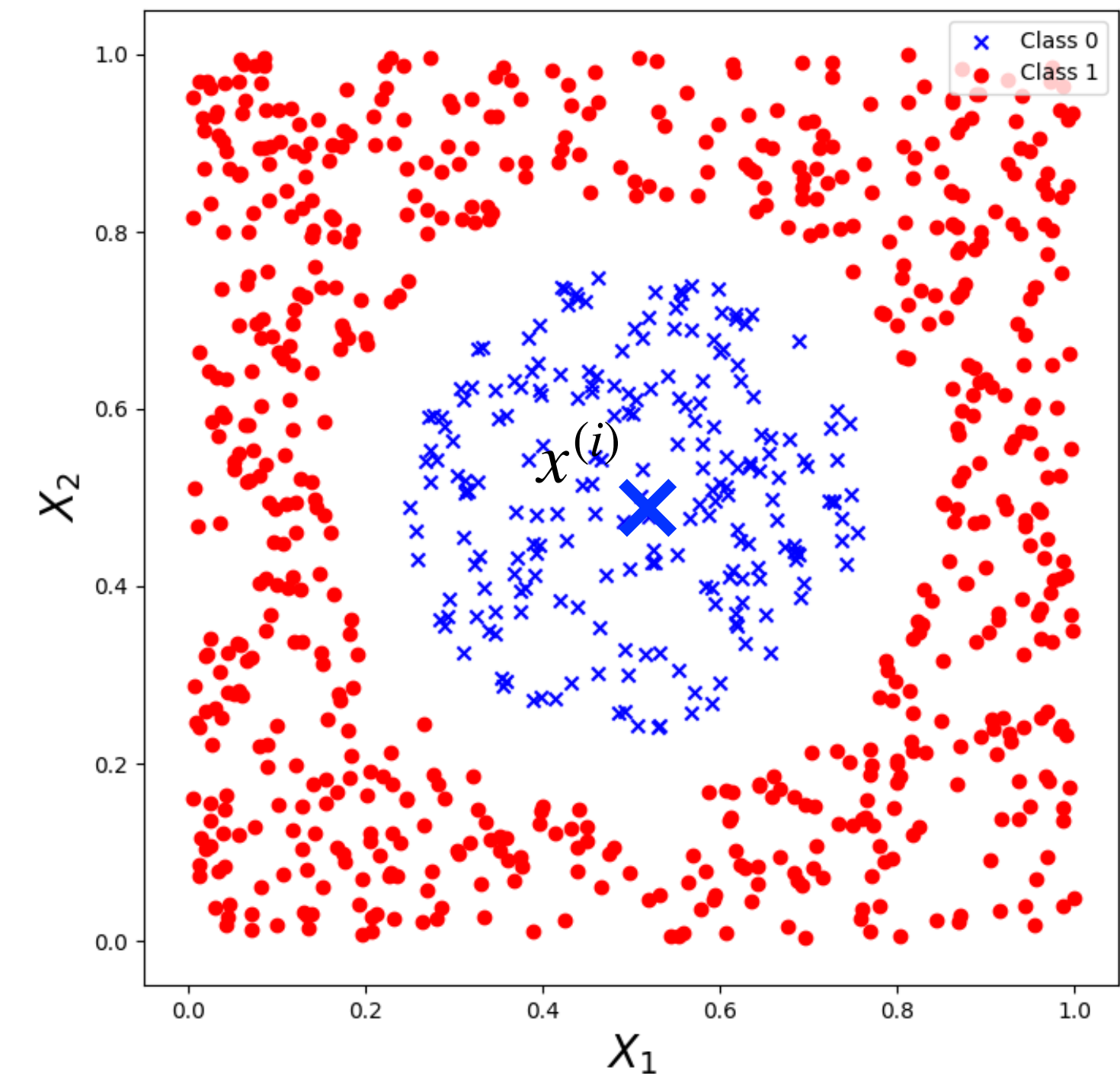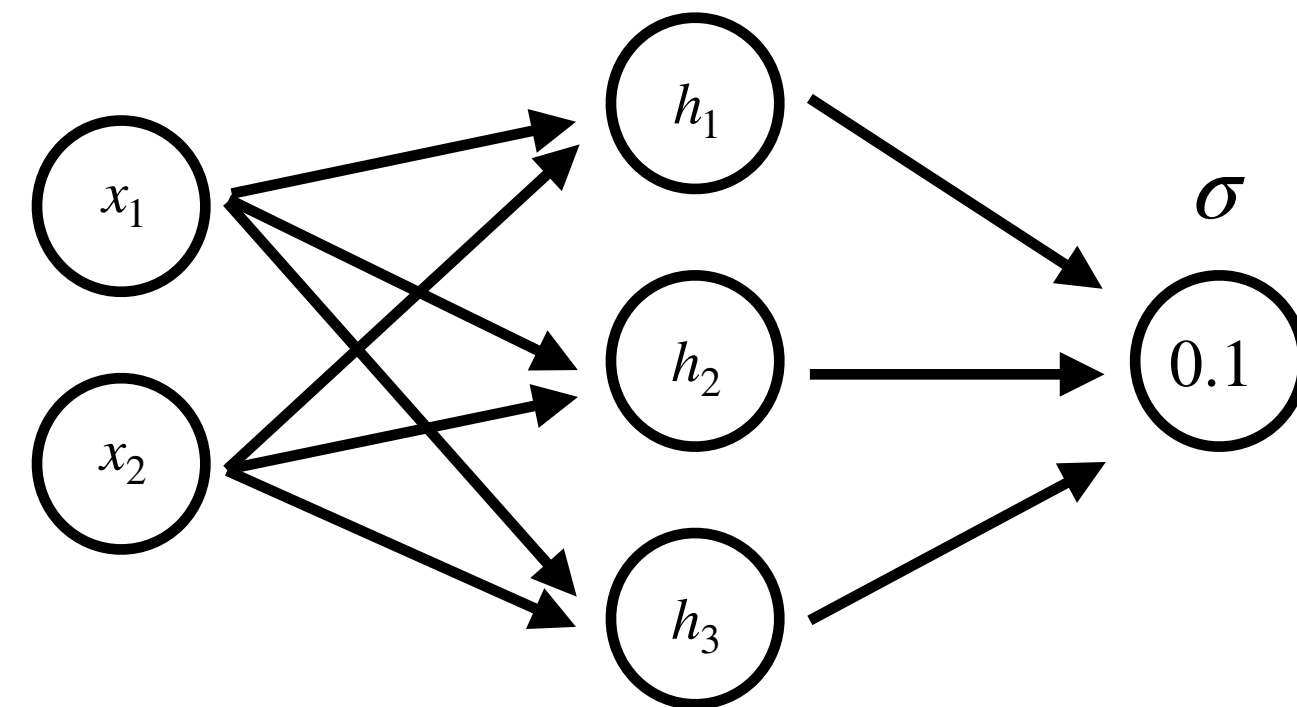**Loss function:** $L(\theta) = \dfrac{1}{n} \sum_{i=1}^{n} -\log(p_{model}(y^{(i)} \mid f(x^{(i)}; \theta)))$
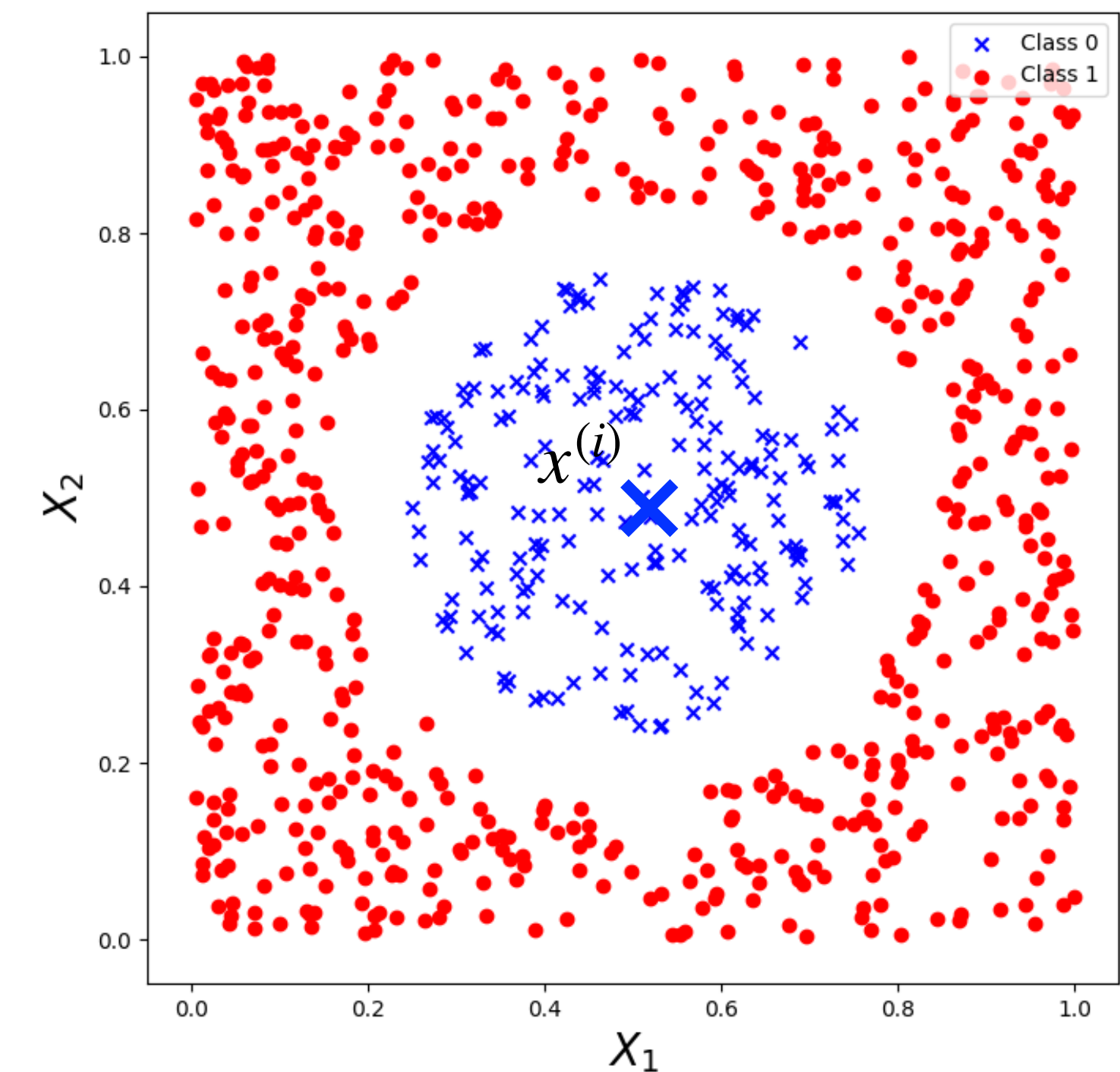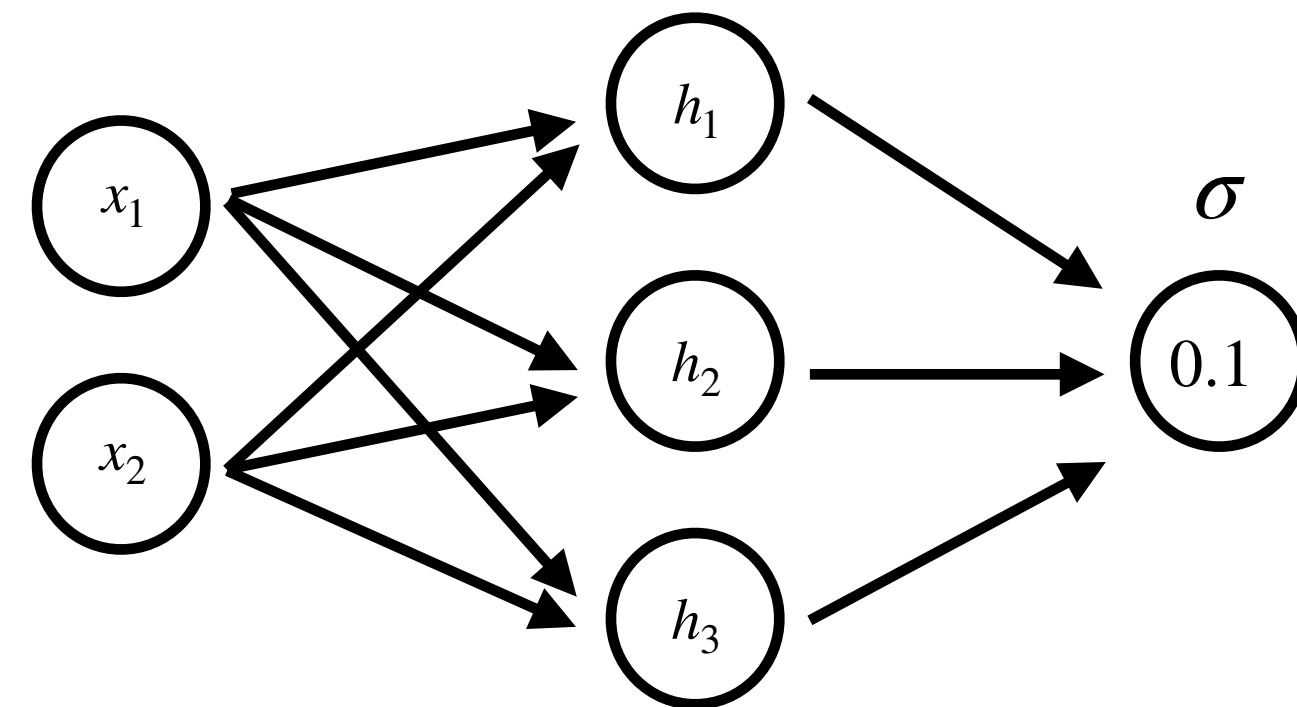
$l^{(i)} = -\boxed{y^{(i)}}\log(\lambda^{(i)}) - (1 - \boxed{y^{(i)}})\log(1 - \lambda^{(i)})$

$\qquad\qquad\quad 0 \qquad\qquad\qquad\quad 0$

| Input x | Class y | Output $\lambda$ | Error $l$ | Predicted $\hat{y}$ |
|---------|---------|------------------|-----------|---------------------|
| (0.5, 0.5) | 0 | 0.1 | **-log(0.9)** | **?** |

# Building the Loss Function: Binary Classification

Training dataset $\{x^{(i)}, y^{(i)}\}_{i=1}^{n}$

$f(x^{(i)}; \theta) = \lambda^{(i)}$ probability that an input $x^{(i)}$ belongs to class 1



**Loss function:** $L(\theta) = \dfrac{1}{n} \sum_{i=1}^{n} -\log(p_{model}(y^{(i)} | f(x^{(i)}; \theta)))$
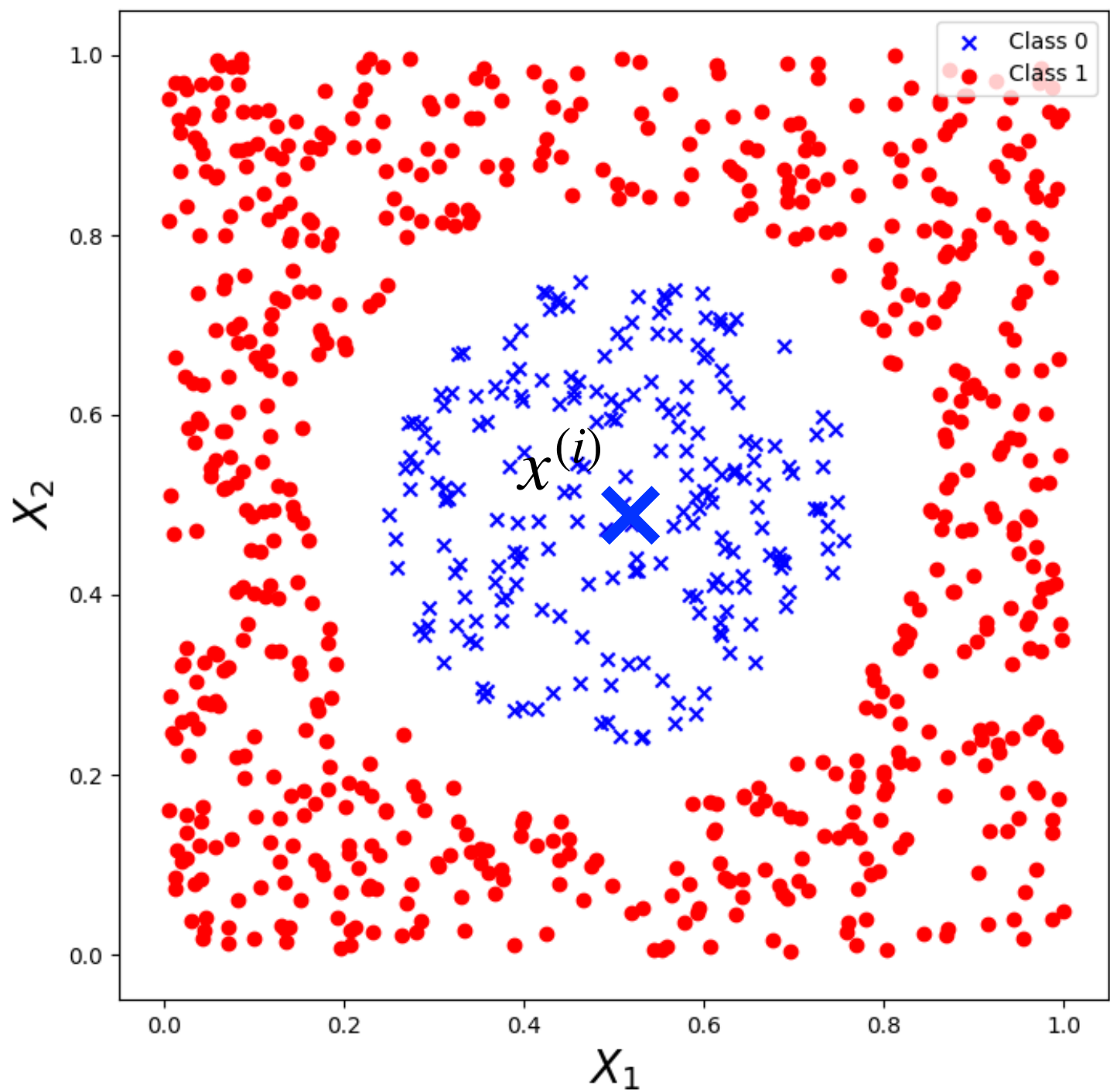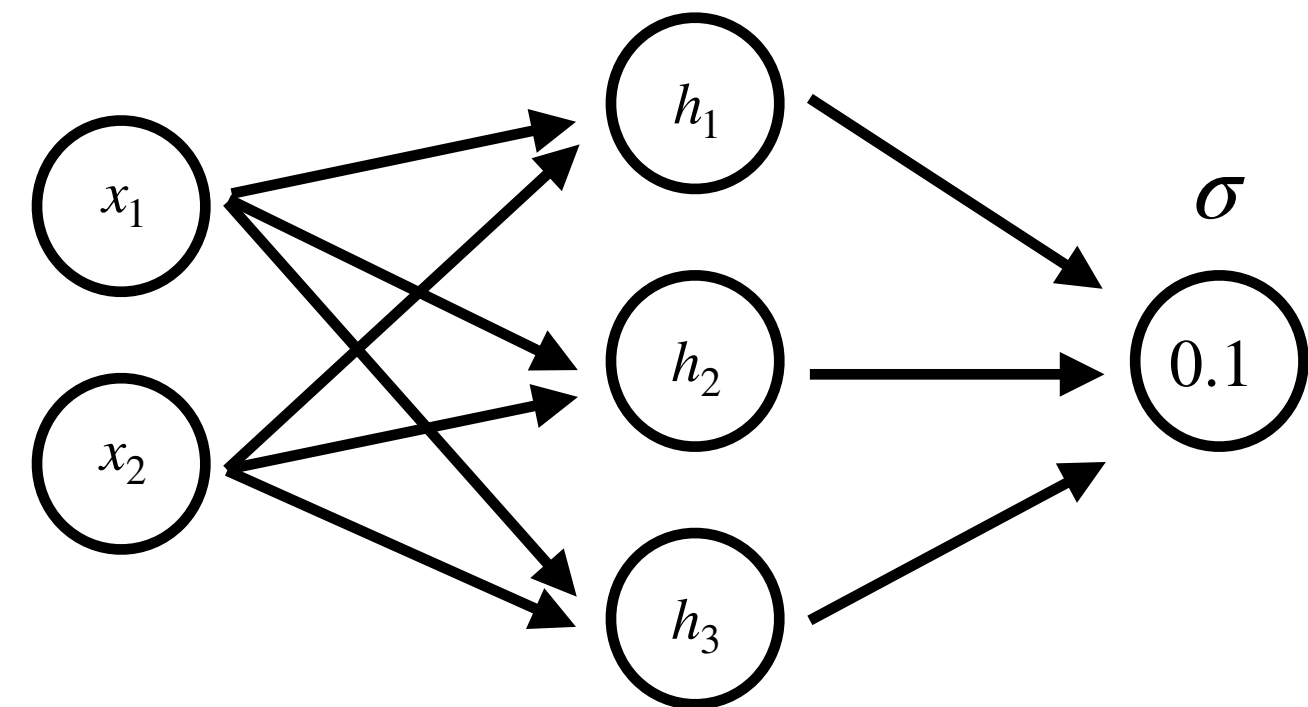
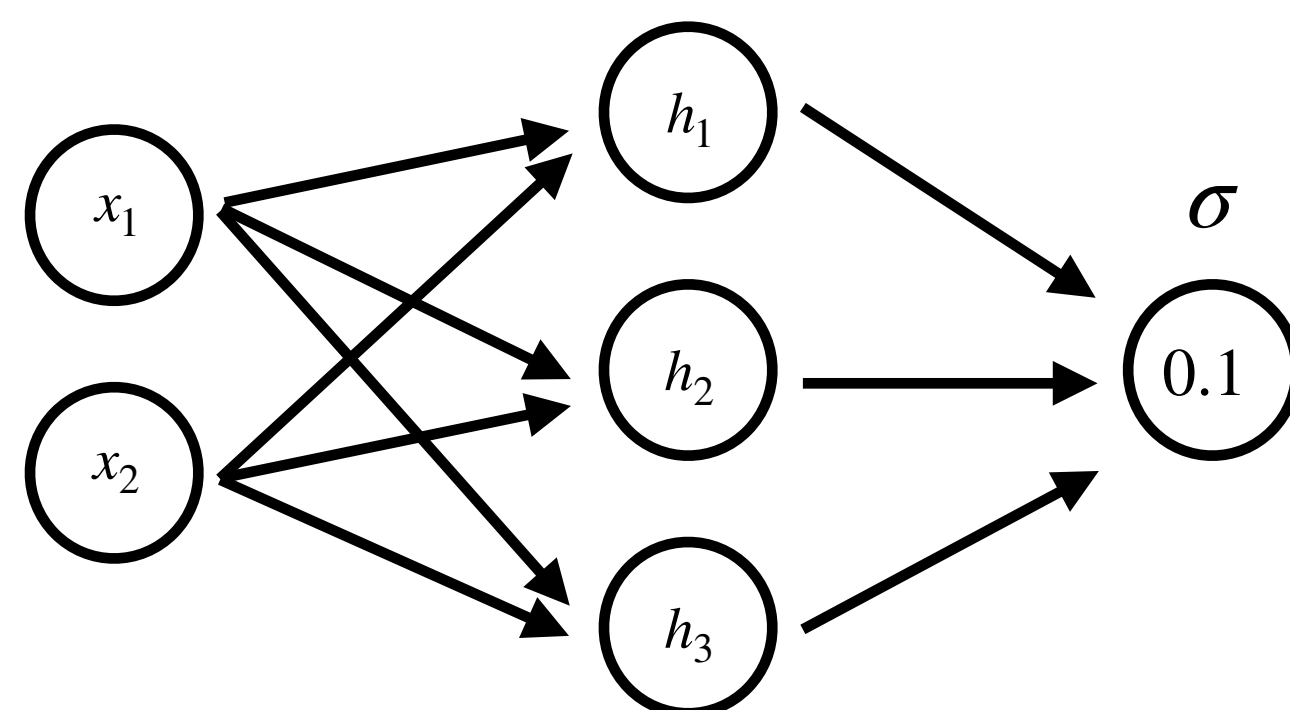$$l^{(i)} = -\log(1 - \lambda^{(i)}) = -\log(0.9)$$

| Input x | Class y | Output $\lambda$ | Error $l$ | Predicted $\hat{y}$ |
|---------|---------|------------------|-----------|---------------------|
| (0.5, 0.5) | 0 | 0.1 | **-log(0.9)** | **?** |

34

# Building the Loss Function: Binary Classification

Training dataset $\{x^{(i)}, y^{(i)}\}_{i=1}^{n}$

$f(x^{(i)}; \theta) = \lambda^{(i)}$ probability that an input $x^{(i)}$ belongs to class 1



**Loss function:** $L_{CE}(\theta) = \dfrac{1}{n} \sum\limits_{i=1}^{n} - y^{(i)} \log(\lambda^{(i)}) - (1 - y^{(i)})\log(1 - \lambda^{(i)})$

| Input x | Class y | Output $\lambda$ | Error $l$ | Predicted $\hat{y}$ |
|---------|---------|------------------|-----------|---------------------|
| (0.5, 0.5) | 0 | 0.1 | **-log(0.9)** | **?** |

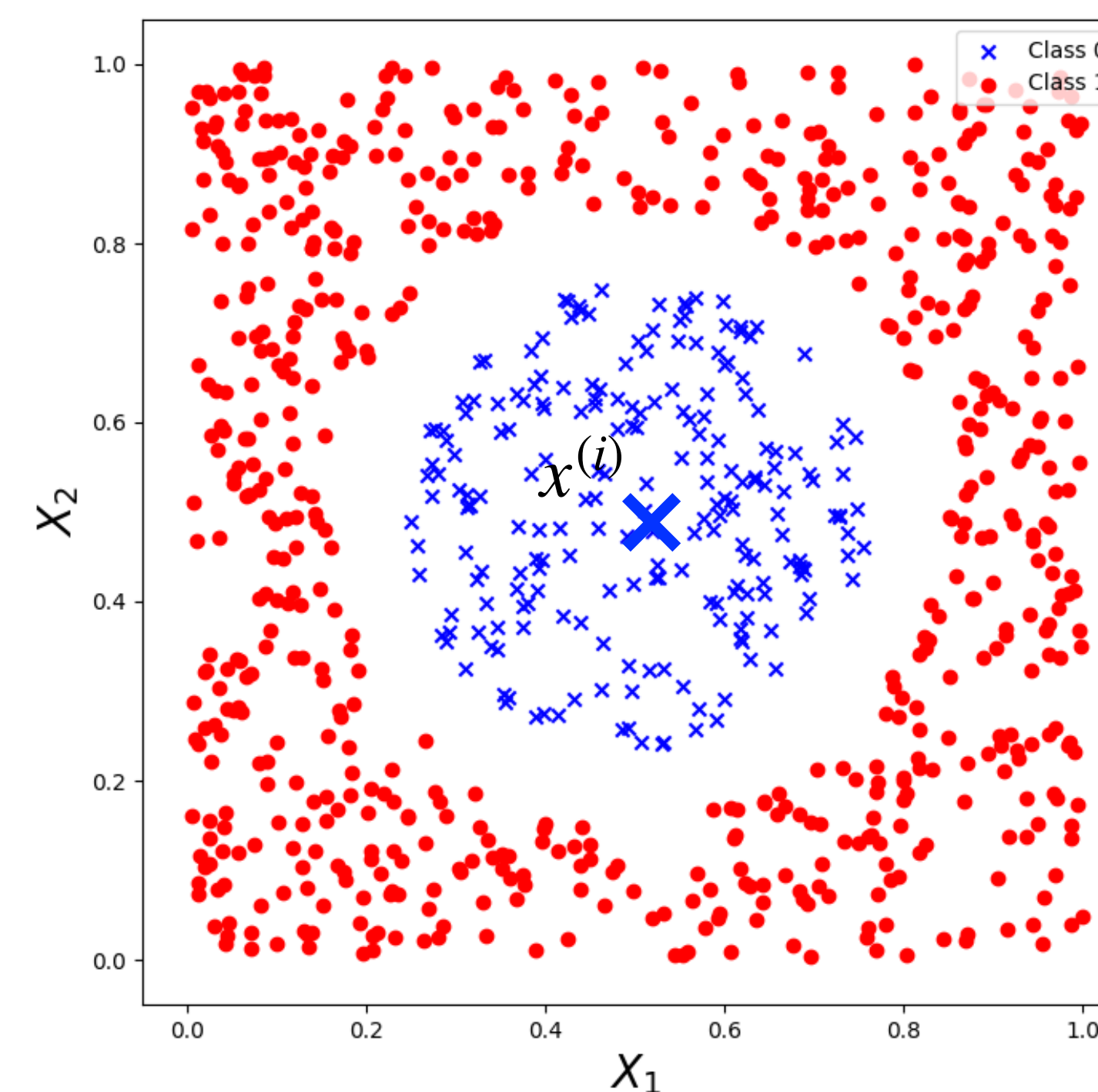# Building the Loss Function: Binary Classification

Training dataset $\{x^{(i)}, y^{(i)}\}_{i=1}^{n}$

$f(x^{(i)}; \theta) = \lambda^{(i)}$ probability that an input $x^{(i)}$ belongs to class 1



**Loss function:** $L_{CE}(\theta) = \dfrac{1}{n} \sum\limits_{i=1}^{n} -y^{(i)} \log(\lambda^{(i)}) - (1 - y^{(i)})\log(1 - \lambda^{(i)})$

$$\hat{y}^{(i)} = argmax_y \left( p_{model}(y \,|\, \lambda^{(i)}) \right)$$

| Input x | Class y | Output $\lambda$ | Error $l$ | Predicted $\hat{y}$ |
|---------|---------|------------------|-----------|---------------------|
| (0.5, 0.5) | 0 | 0.1 | -log(0.9) | **?** |

36

# Building the Loss Function: Binary Classification

Training dataset $\{x^{(i)}, y^{(i)}\}_{i=1}^n$

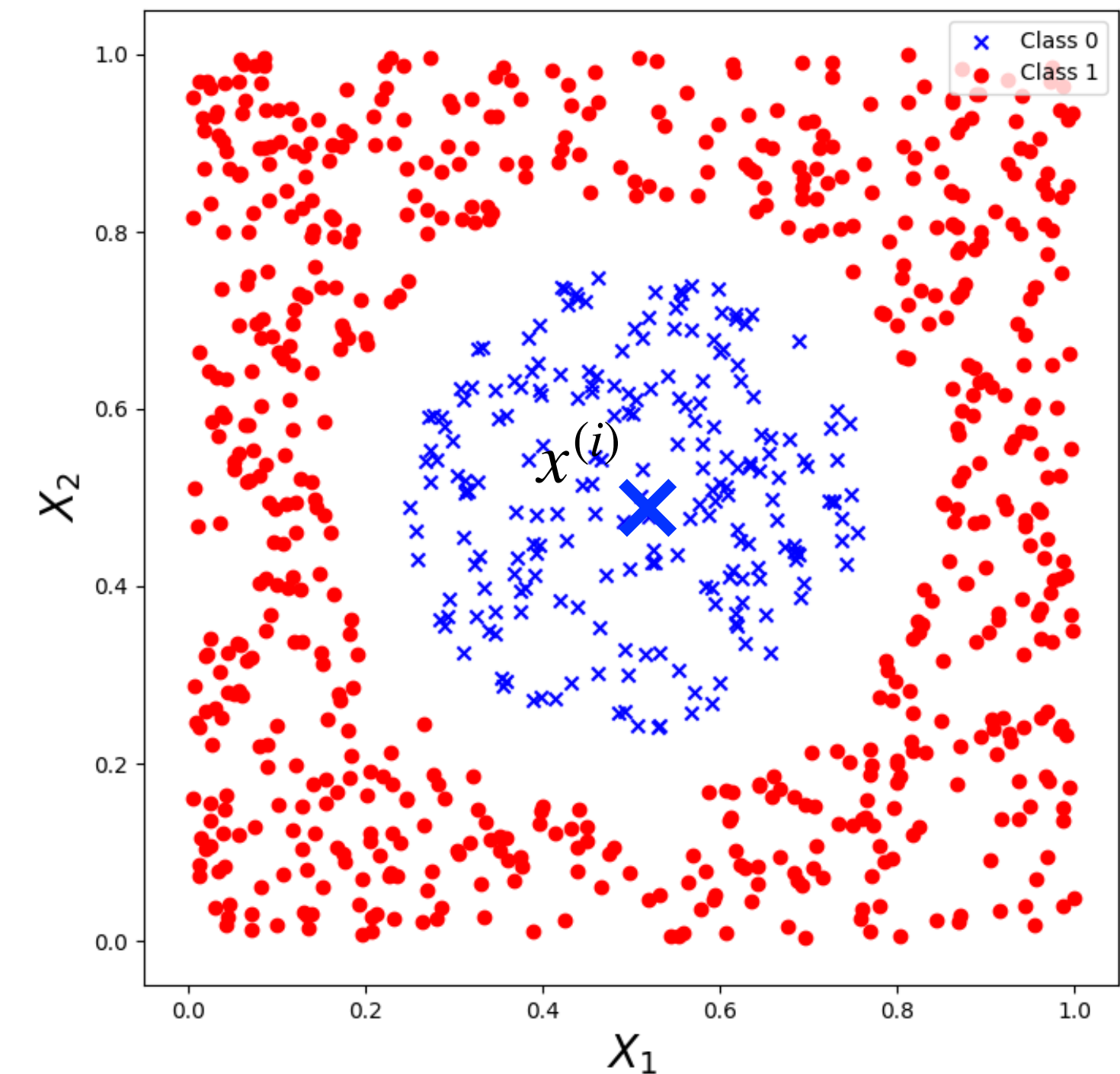$f(x^{(i)}; \theta) = \lambda^{(i)}$ probability that an input $x^{(i)}$ belongs to class 1



**Loss function:** $L_{CE}(\theta) = \dfrac{1}{n} \displaystyle\sum_{i=1}^n - y^{(i)} \log(\lambda^{(i)}) - (1 - y^{(i)}) \log(1 - \lambda^{(i)})$

$$\hat{y}^{(i)} = argmax_y \left( p_{model}(y \,|\, \lambda^{(i)}) \right)$$

$y = 0 \quad y = 1$

$$p_{model}(y \,|\, \lambda^{(i)} = 0.1) = (0.9, \ 0.1)$$

| Input x | Class y | Output $\lambda$ | Error $l$ | Predicted $\hat{y}$ |
|---------|---------|---------|---------|---------|
| (0.5, 0.5) | 0 | 0.1 | -log(0.9) | **?** |

# Building the Loss Function: Binary Classification

Training dataset $\{x^{(i)}, y^{(i)}\}_{i=1}^{n}$

$f(x^{(i)}; \theta) = \lambda^{(i)}$ probability that an input $x^{(i)}$ belongs to class 1
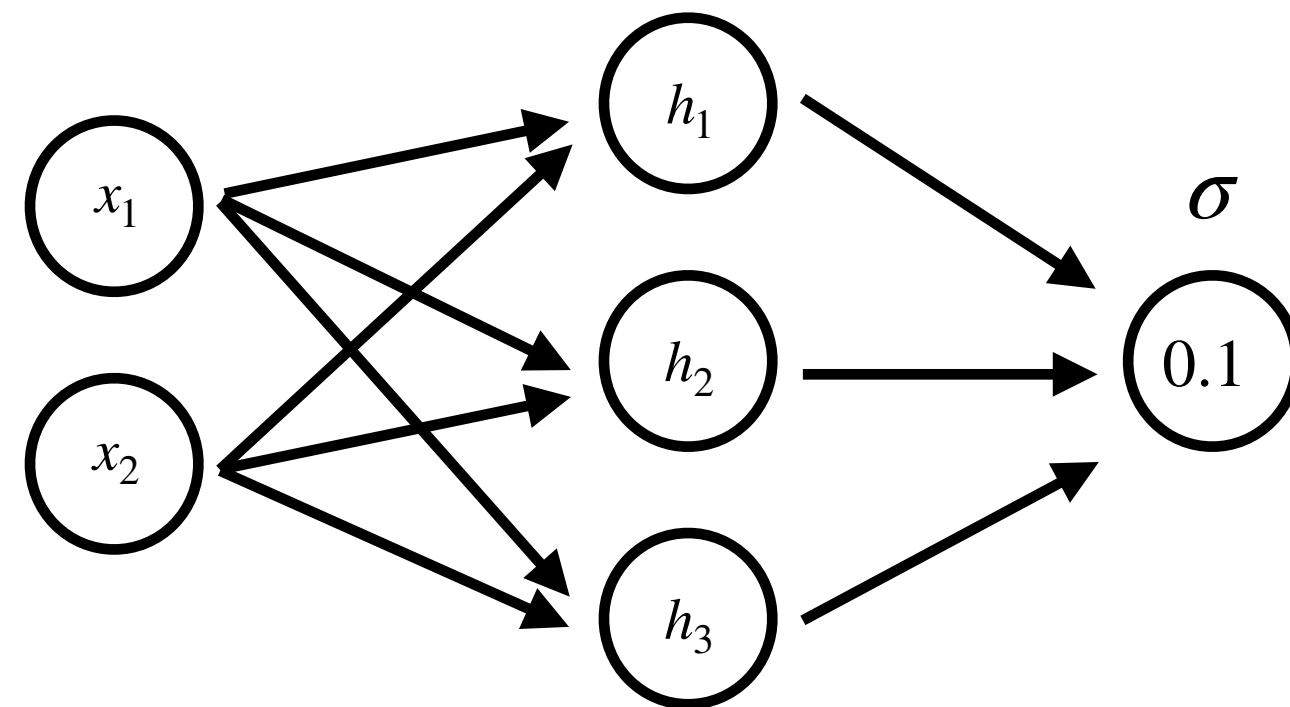


**Loss function:** $L_{CE}(\theta) = \dfrac{1}{n} \sum\limits_{i=1}^{n} - y^{(i)} \log(\lambda^{(i)}) - (1 - y^{(i)}) \log(1 - \lambda^{(i)})$
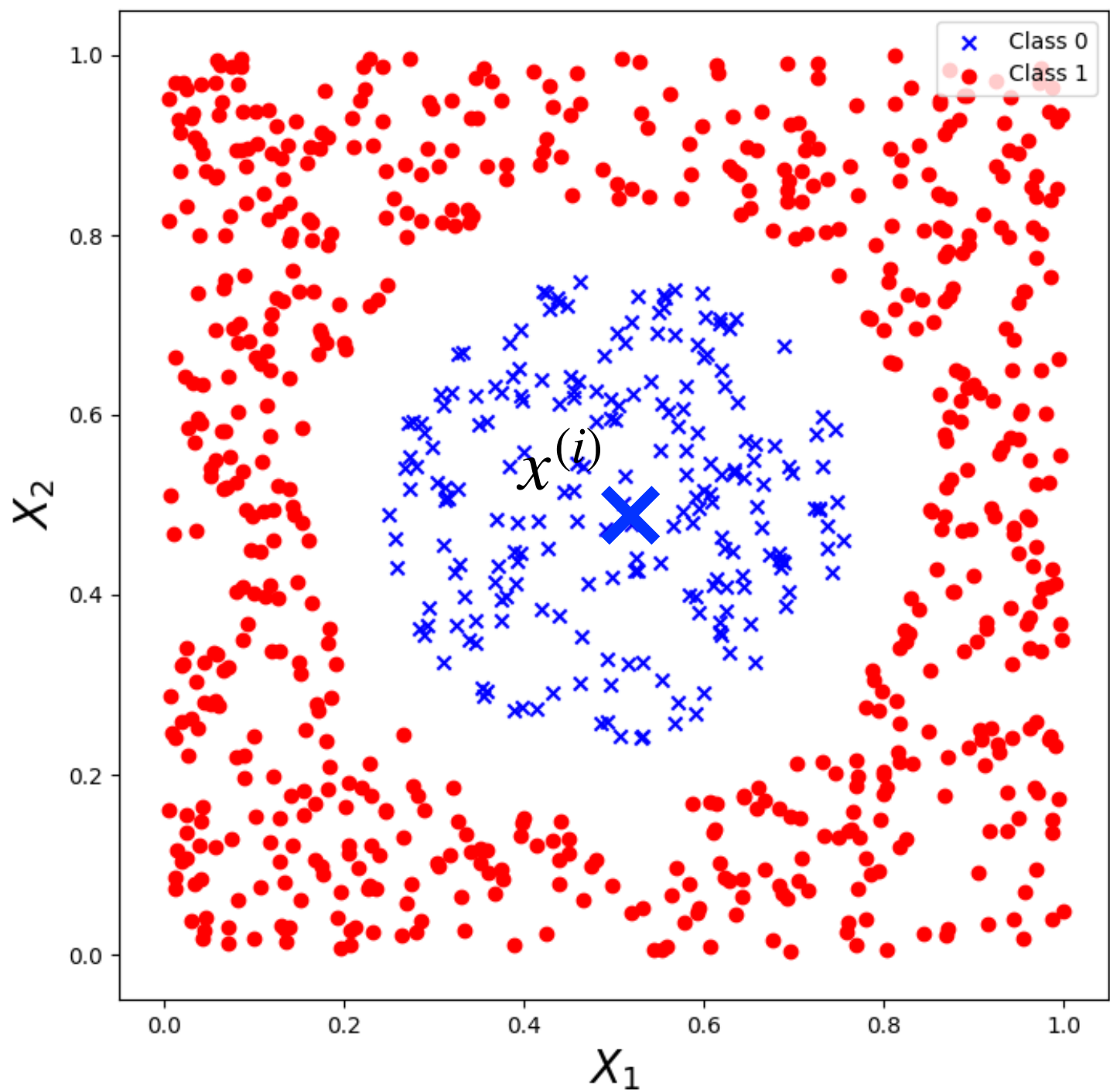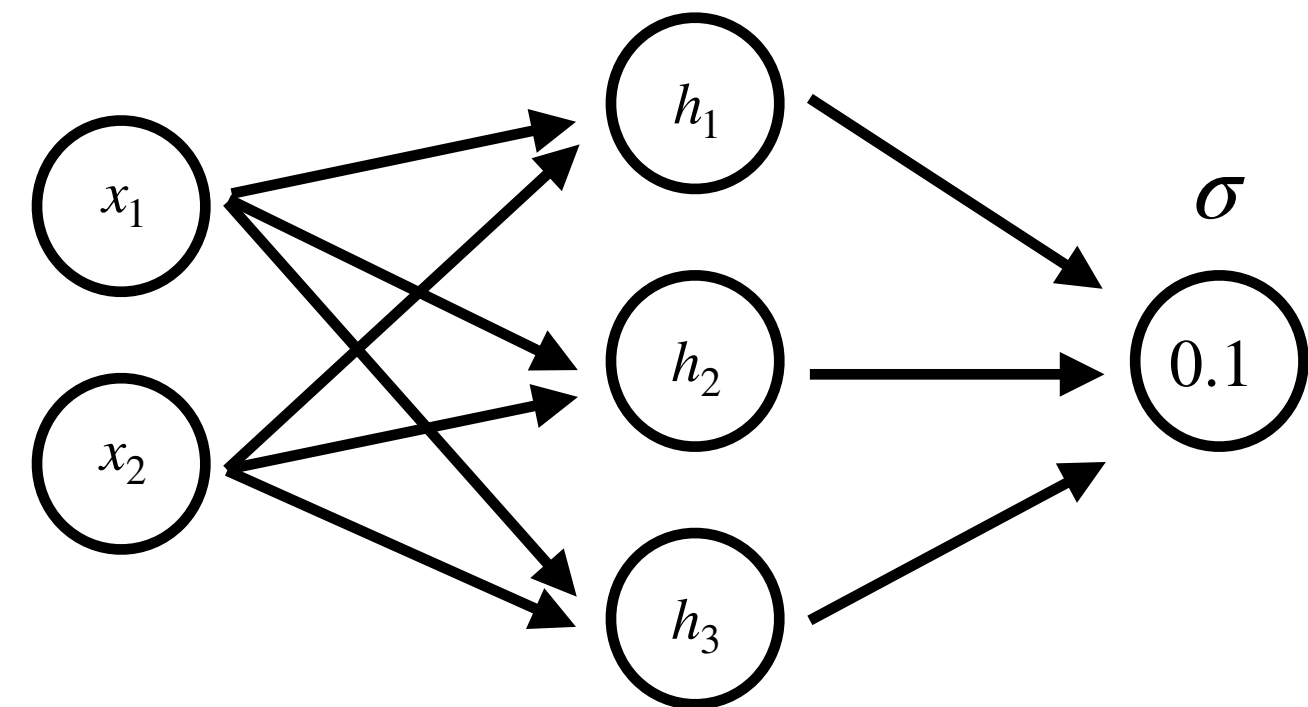
$$\hat{y}^{(i)} = argmax_y \left( (0.9, 0.1) \right) = 0$$

| Input x | Class y | Output $\lambda$ | Error $l$ | Predicted $\hat{y}$ |
|---------|---------|--------|---------|-------------|
| (0.5, 0.5) | 0 | 0.1 | -log(0.9) | **0** |

# Cross Entropy

Given the ground-truth distribution $p_{data}$ and an approximation $p_{model}$

$$CE(p_{data}, p_{model}) = -\sum_{y} p_{data}(y) \, \log(p_{model}(y))^*$$

*Colah's blog: Visual Information Theory

# Cross Entropy

Given the ground-truth distribution $p_{data}$ and an approximation $p_{model}$

$$CE(p_{data}, p_{model}) = -\sum_y p_{data}(y) \log(p_{model}(y))^*$$

The loss $L$ estimates $CE(p_{data}, p_{model})$



$CE(\delta_0, (1-\lambda, \lambda))$

$p_{data} = \delta_0$

$p_{model} = (1-\lambda, \lambda)$

$CE$

$\lambda$

# Cross Entropy

Given the ground-truth distribution $p_{data}$ and an approximation $p_{model}$

$$CE(p_{data}, p_{model}) = -\sum_{y} p_{data}(y) \log(p_{model}(y))^*$$

The loss $L$ estimates $CE(p_{data}, p_{model})$



$CE(\delta_0, (1-\lambda, \lambda))$
$p_{data} = \delta_0$
$p_{model} = (1-\lambda, \lambda)$

$CE$

$\hat{\lambda} = 0$

$\lambda$

*Colah's blog: Visual Information Theory

Optimal output

# Cross Entropy

Given the ground-truth distribution $p_{data}$ and an approximation $p_{model}$

$$CE(p_{data}, p_{model}) = -\sum_y p_{data}(y) \log(p_{model}(y))^*$$

The loss $L$ estimates $CE(p_{data}, p_{model})$



$CE(\delta_0, (1-\lambda, \lambda))$
$p_{data} = \delta_0$
$p_{model} = (1-\lambda, \lambda)$

$CE$

$CE(\delta_0, (1,0)) = 0$

$\hat{\lambda} = 0$

$\lambda$

$\hat{p}_{model} = (1,0) = \delta_0 = p_{data}$
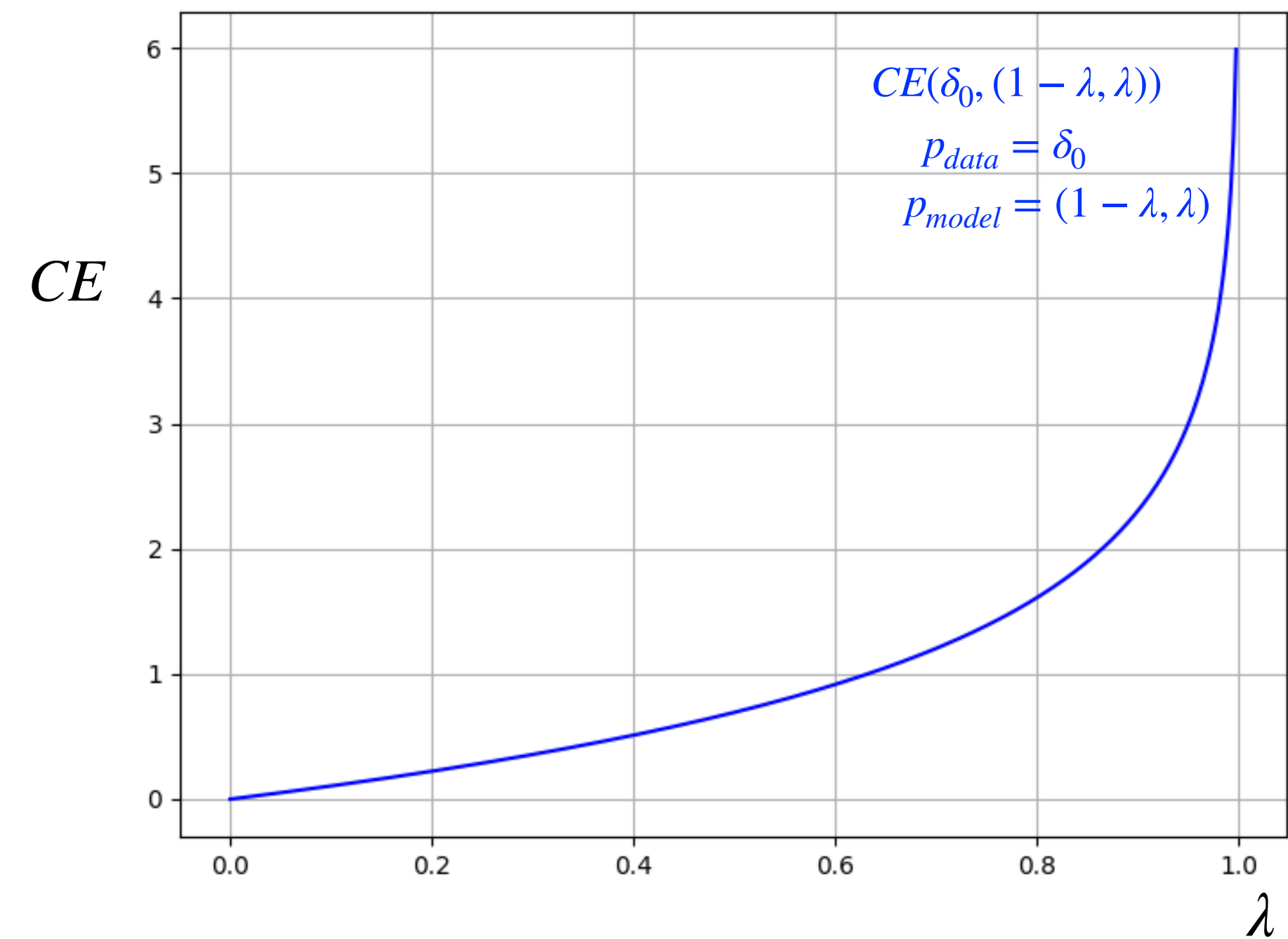
Optimal output

*Colah's blog: Visual Information Theory

# Cross Entropy

Given the ground-truth distribution $p_{data}$ and an approximation $p_{model}$

$$CE(p_{data}, p_{model}) = -\sum_{y} p_{data}(y) \, \log(p_{model}(y))*$$

The loss $L$ estimates $CE(p_{data}, p_{model})$



*CE*

$CE((0.5,0.5), (1-\lambda, \lambda))$
$p_{data} = (0.5,0.5)$
$p_{model} = (1-\lambda, \lambda)$
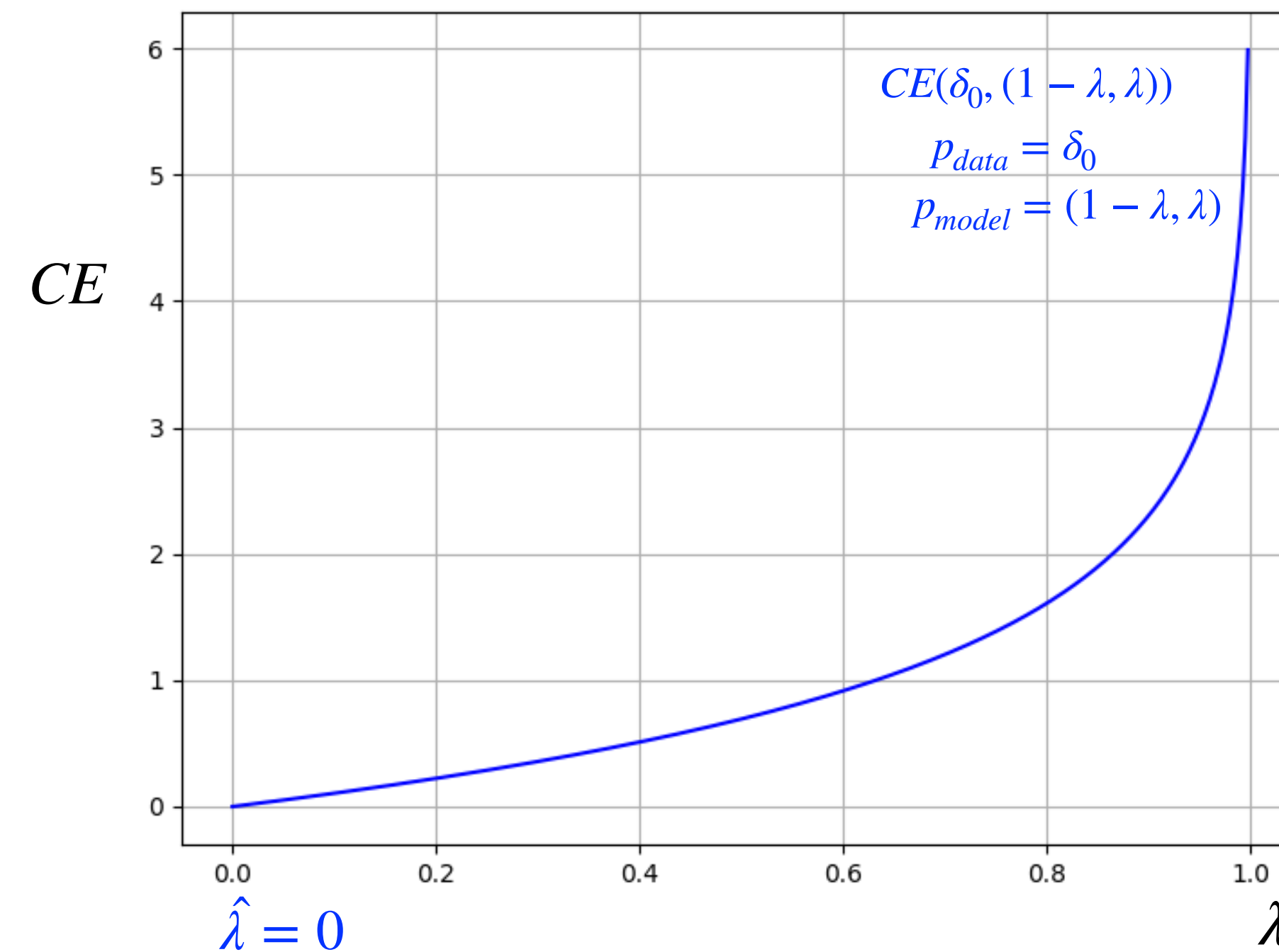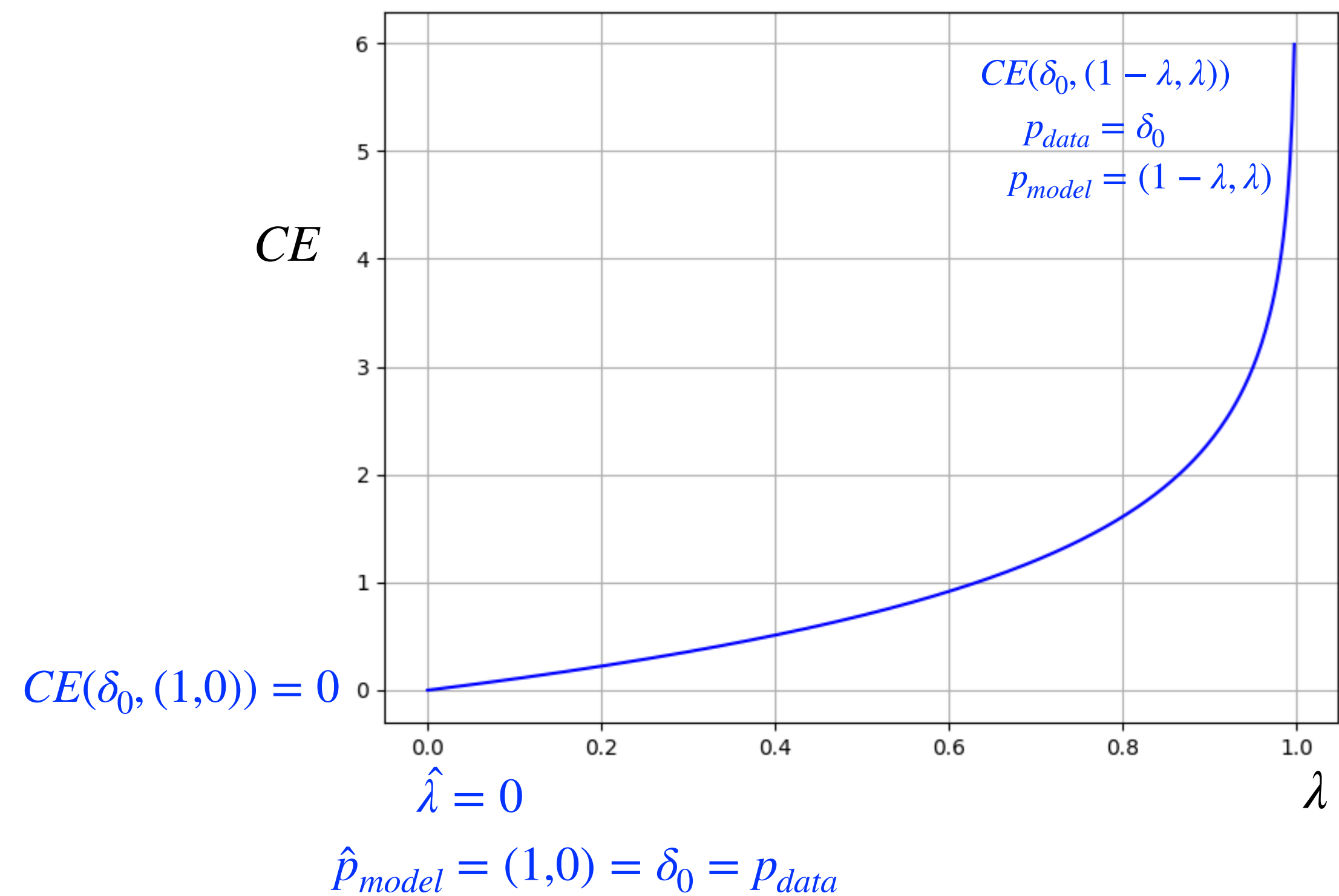
$\lambda$

*Colah's blog: Visual Information Theory

# Cross Entropy

Given the ground-truth distribution $p_{data}$ and an approximation $p_{model}$

$$CE(p_{data}, p_{model}) = -\sum_{y} p_{data}(y) \, \log(p_{model}(y))^*$$

The loss $L$ estimates $CE(p_{data}, p_{model})$



$CE((0.5,0.5),(1-\lambda,\lambda))$
$p_{data} = (0.5,0.5)$
$p_{model} = (1-\lambda,\lambda)$

$CE$

$\hat{\lambda} = 0.5$

$\lambda$

$\hat{p}_{model} = (0.5, 0.5) = p_{data}$

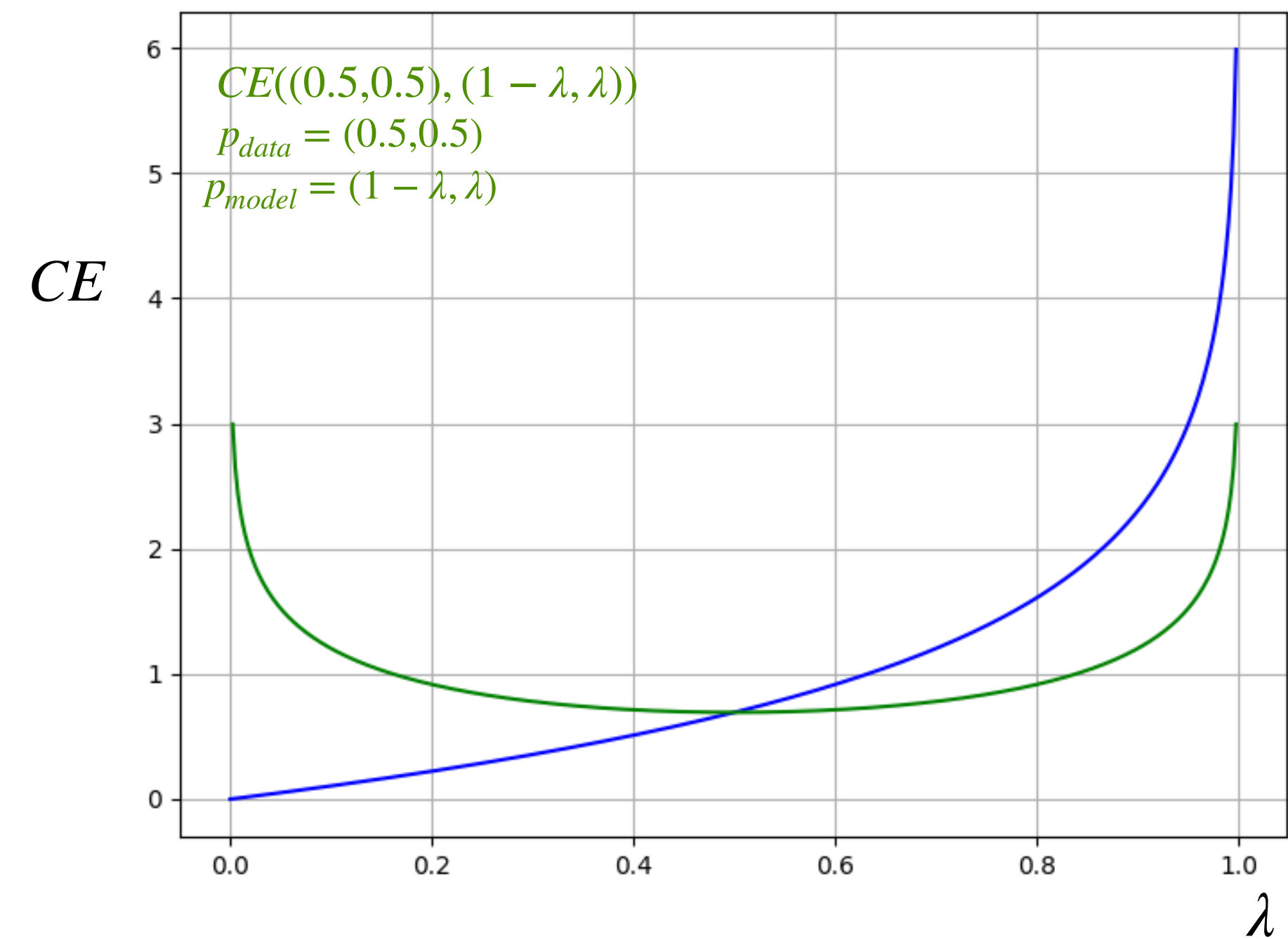Optimal output

44

*Colah's blog: Visual Information Theory

# Cross Entropy

Given the ground-truth distribution $p_{data}$ and an approximation $p_{model}$

$$CE(p_{data}, p_{model}) = -\sum_y p_{data}(y) \log(p_{model}(y))^*$$

The loss $L$ estimates $CE(p_{data}, p_{model})$



$CE((0.5,0.5), (1-\lambda, \lambda))$
$p_{data} = (0.5, 0.5)$
$p_{model} = (1-\lambda, \lambda)$

$CE$

$CE((0.5,0.5), (0.5,0.5)) = 0.69$

$\hat{\lambda} = 0.5$

$\hat{p}_{model} = (0.5, 0.5) = p_{data}$

$\lambda$

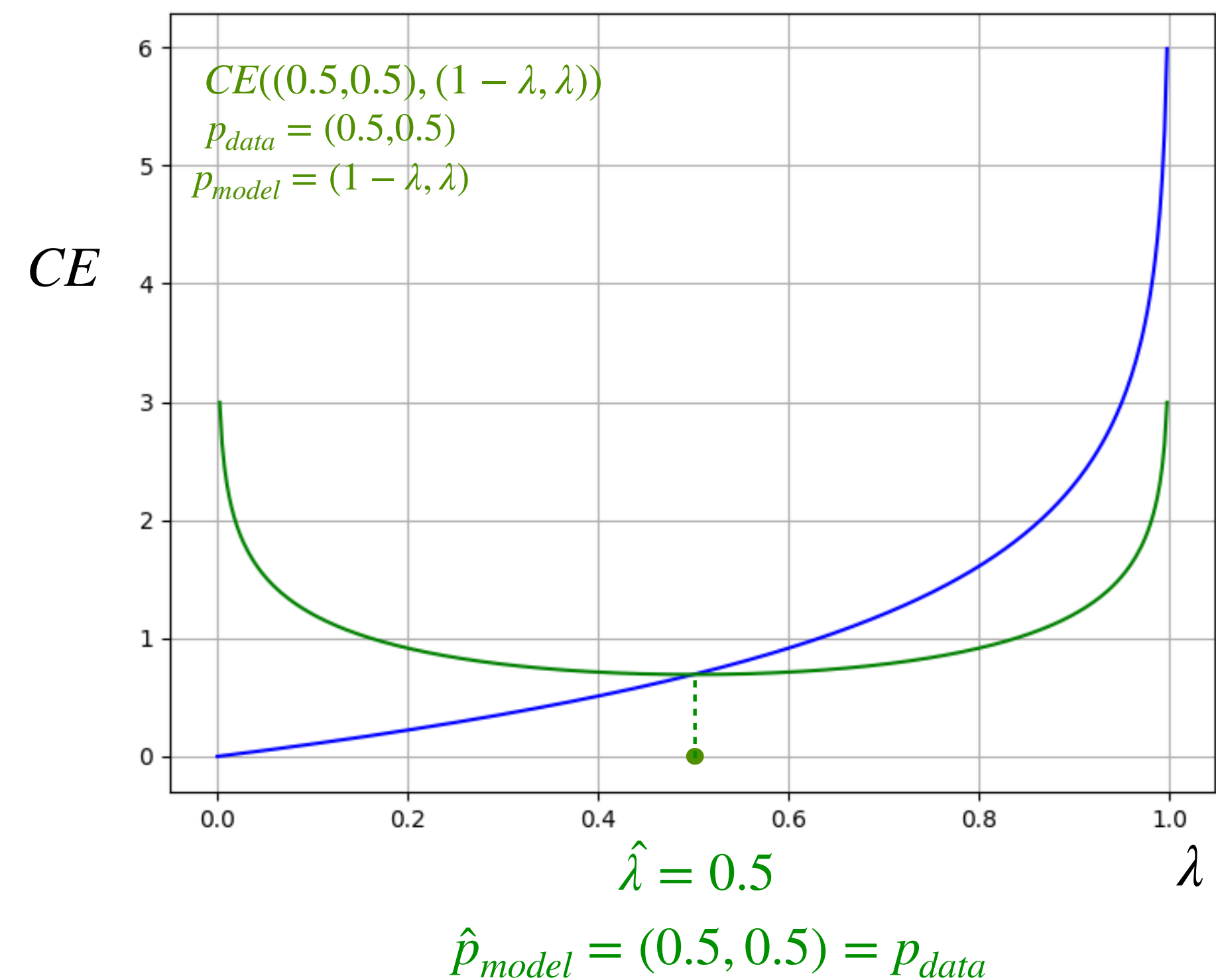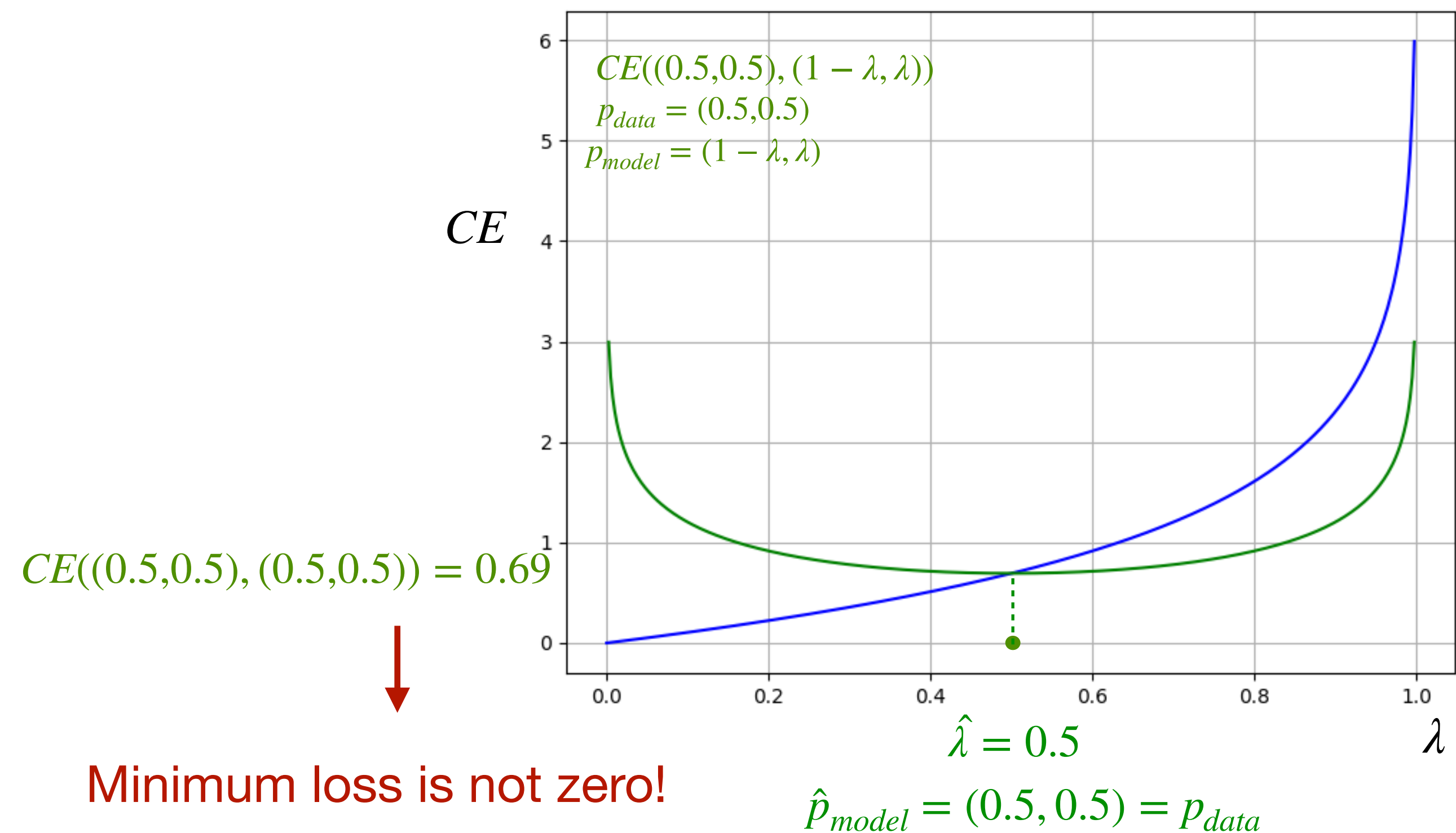**Minimum loss is not zero!**

Optimal output

# Cross Entropy

Given the ground-truth distribution $p_{data}$ and an approximation $p_{model}$

$$CE(p_{data}, p_{model}) = - \sum_y p_{data}(y) \, \log(p_{model}(y))^*$$

The loss $L$ estimates $CE(p_{data}, p_{model})$
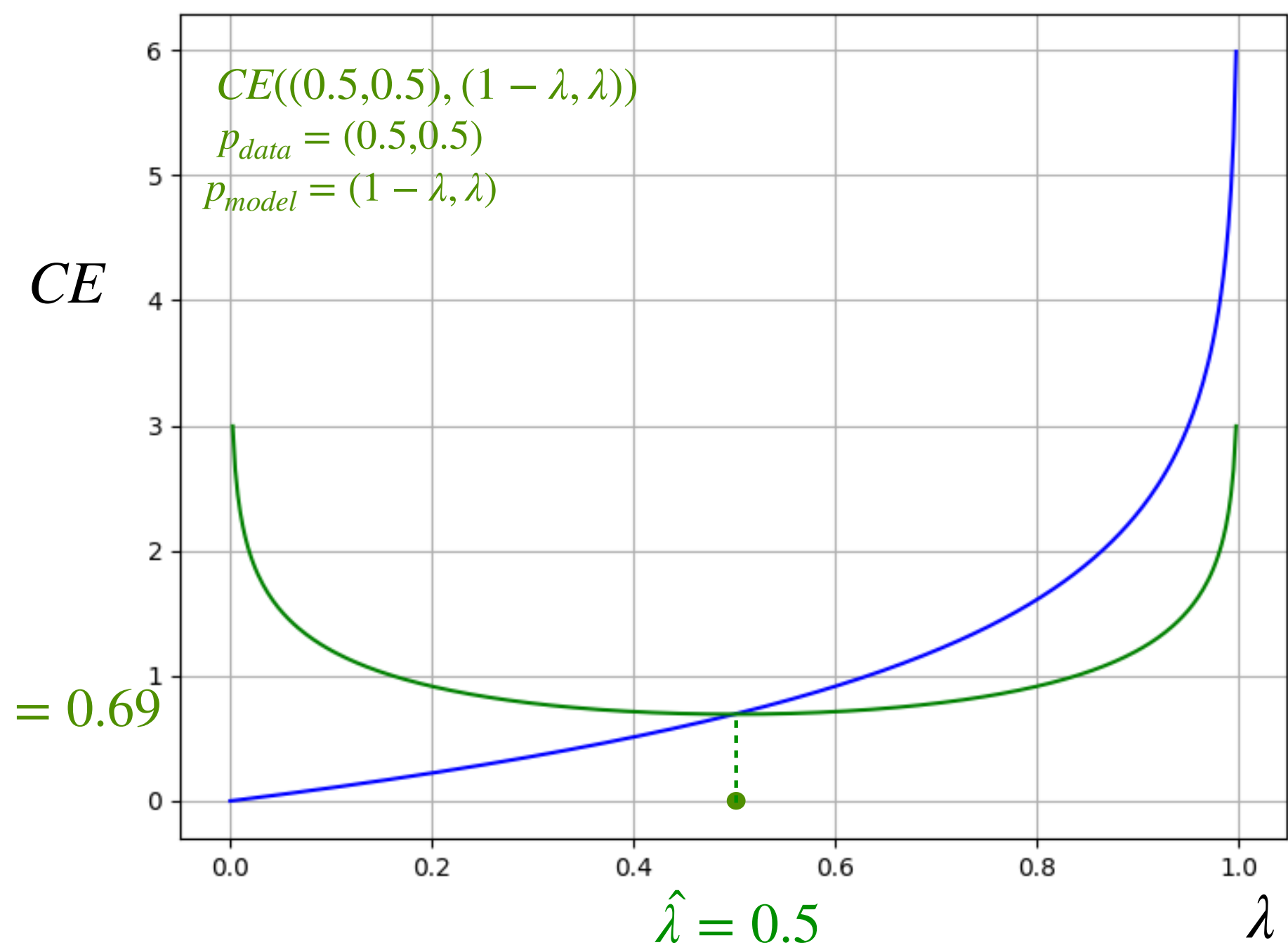
A prefect model $p_{data} = p_{model}$

$$CE(p_{data}, p_{model}) = H(p_{data}) \text{ entropy}$$



$CE((0.5,0.5),(1-\lambda,\lambda))$
$p_{data} = (0.5,0.5)$
$p_{model} = (1-\lambda,\lambda)$

$CE$

$CE((0.5,0.5),(0.5,0.5)) = 0.69$

**Minimum loss is not zero!**

$\hat{\lambda} = 0.5$

$\hat{p}_{model} = (0.5, 0.5) = p_{data}$

$\lambda$

Optimal output

# Maximum Likelihood

- Open the notebook in Brightspace $\longrightarrow$ Gradient-based Optimization $\longrightarrow$ Maximum Likelihood - Bandits

- Work in group of 2 or 3

- Go through the notebook and answer the questions with your peers

- If you have doubts, raise your hand or come to me

- Around 15 minutes

Discussion next Monday

# Loss Function Minimization

Solve the minimization problem $\hat{\theta} = argmin_\theta L(\theta)$ through gradient descent.

$$\theta = \{w, b\}$$

# Loss Function Minimization

Solve the minimization problem $\hat{\theta} = argmin_\theta L(\theta)$ through gradient descent.

$$\theta = \{w, b\}$$



$L(\theta)$

$w$

$b$



$b$

$w$
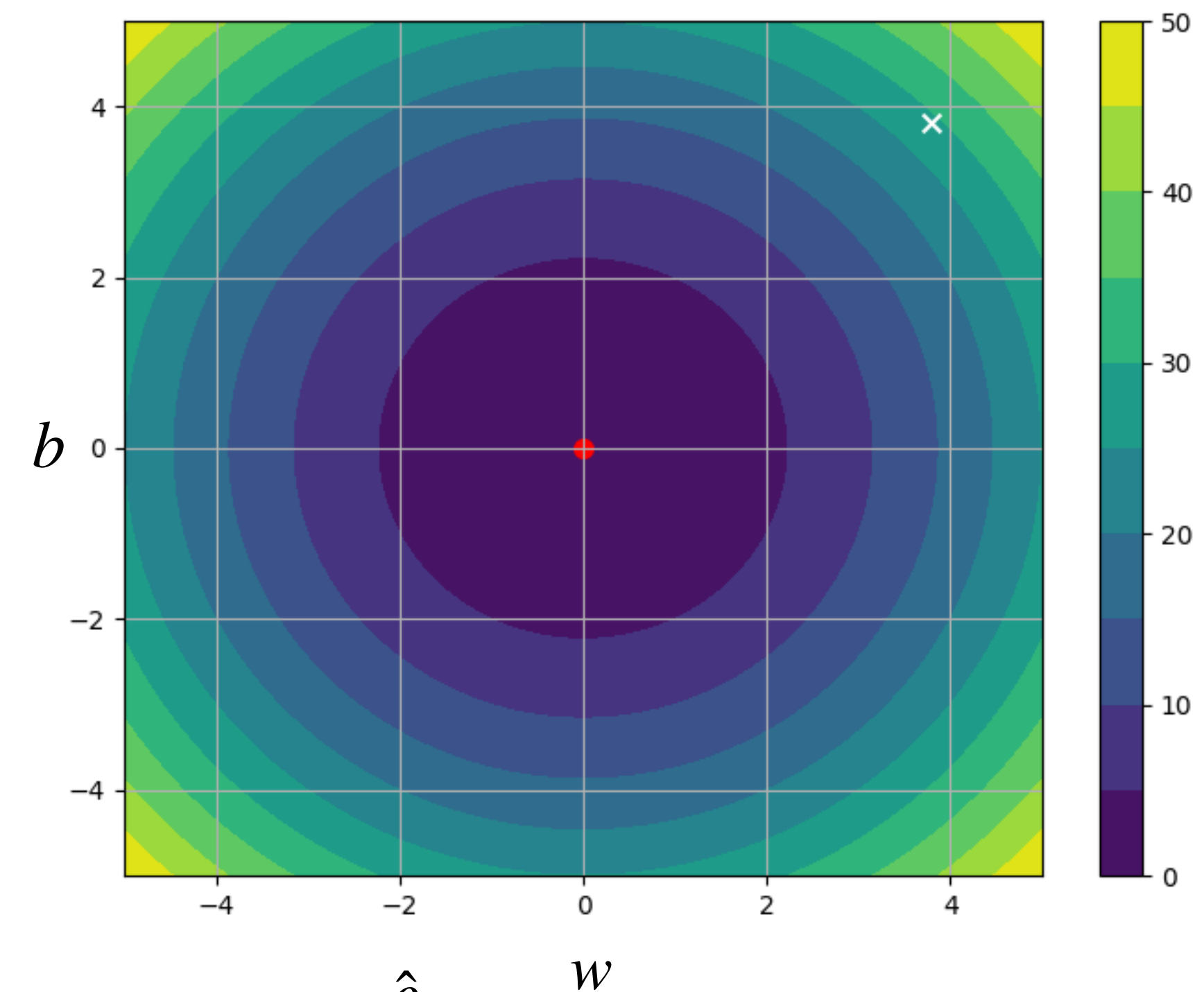
●   $\hat{\theta}$

✖   $\theta_0 = (w_0, b_0)$

# Loss Function Minimization

Solve the minimization problem $\hat{\theta} = argmin_{\theta} L(\theta)$ through gradient descent.

$$\theta = \{w, b\}$$



$$-\nabla L(\theta_0) = -\left( \frac{\partial L}{\partial w}(\theta_0), \frac{\partial L}{\partial b}(\theta_0) \right)$$

$\bullet$    $\hat{\theta}$

$\times$    $\theta_0 = (w_0, b_0)$
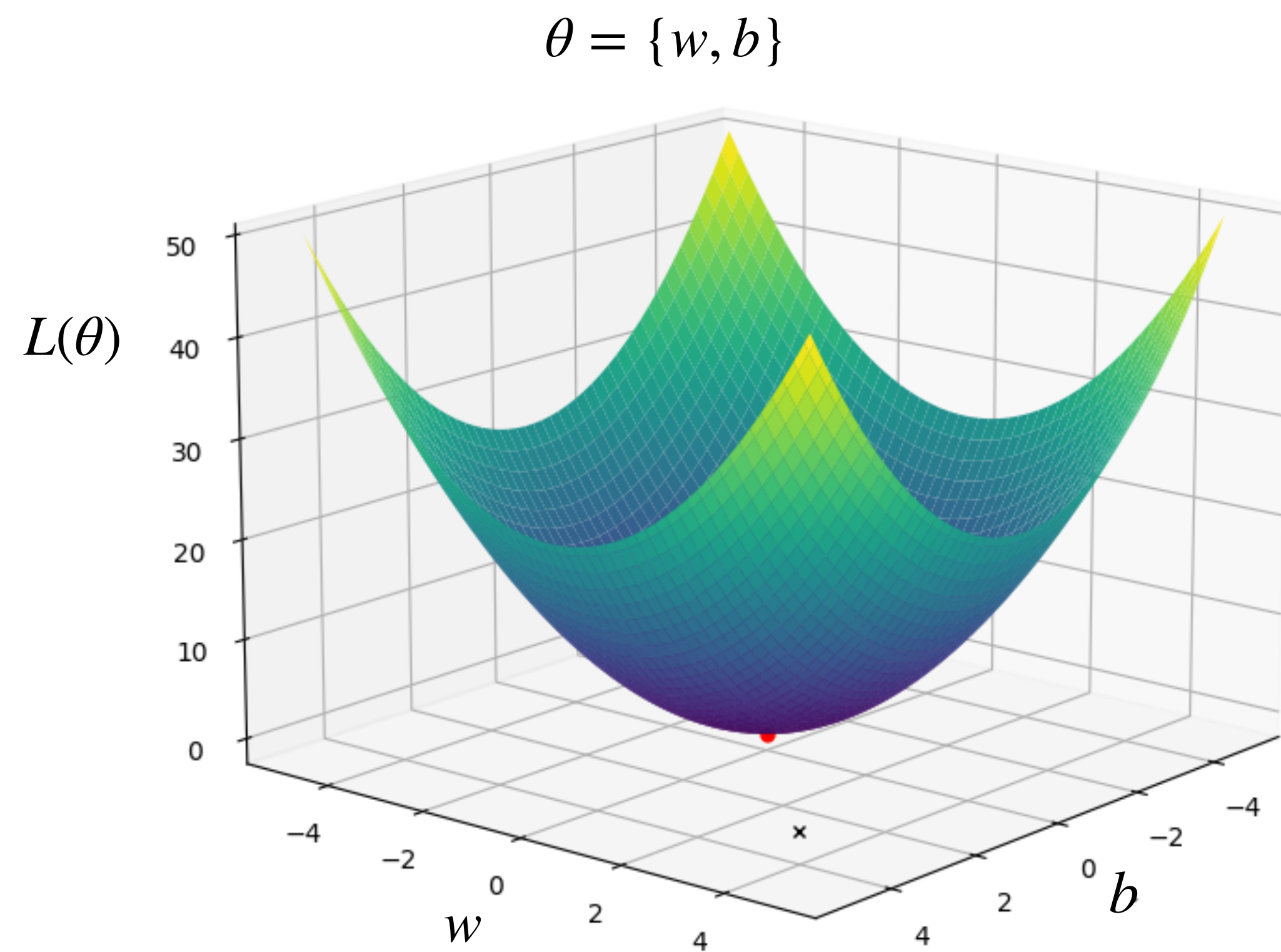
$\uparrow$    $-\alpha \nabla L(\theta_0)$

50

# Loss Function Minimization

Solve the minimization problem $\hat{\theta} = argmin_{\theta} L(\theta)$ through gradient descent.

$$\theta = \{w, b\}$$



$$\theta_1 = \theta_0 - \alpha \nabla L(\theta_0)$$



- $\bullet$    $\hat{\theta}$
- $\times$    $\theta_1 = (w_1, b_1)$
- $\uparrow$    $-\alpha \nabla L(\theta_0)$

# Loss Function Minimization

Solve the minimization problem $\hat{\theta} = argmin_\theta L(\theta)$ through gradient descent.

$$\theta = \{w, b\}$$
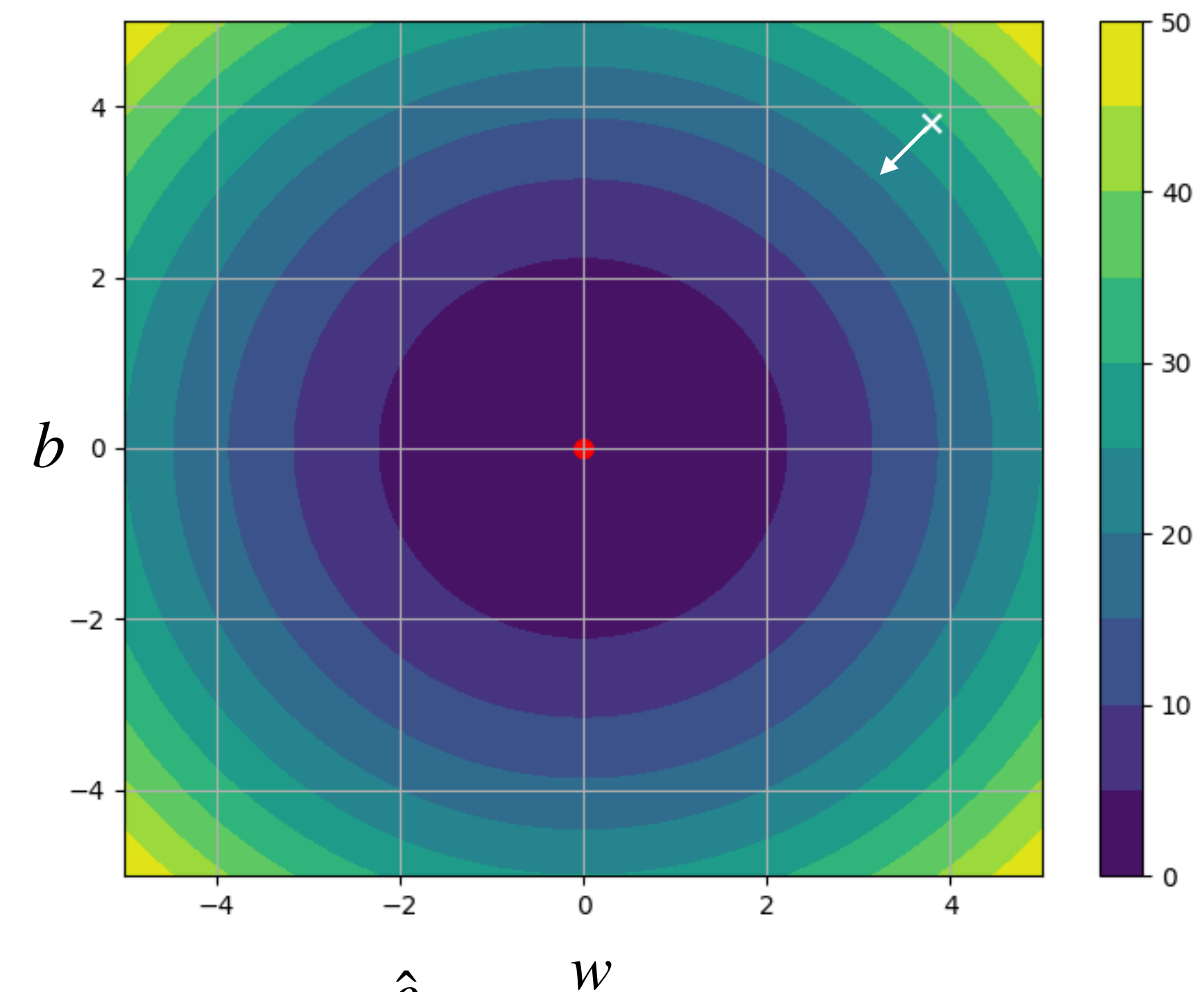


$$-\nabla L(\theta_1) = -\left( \frac{\partial L}{\partial w}(\theta_1), \frac{\partial L}{\partial b}(\theta_1) \right)$$

$\bullet$   $\hat{\theta}$

$\times$   $\theta_1 = (w_1, b_1)$

$\uparrow$   $-\alpha \nabla L(\theta_1)$

# Loss Function Minimization

Solve the minimization problem $\hat{\theta} = argmin_{\theta}L(\theta)$ through gradient descent.

$$\theta = \{w, b\}$$
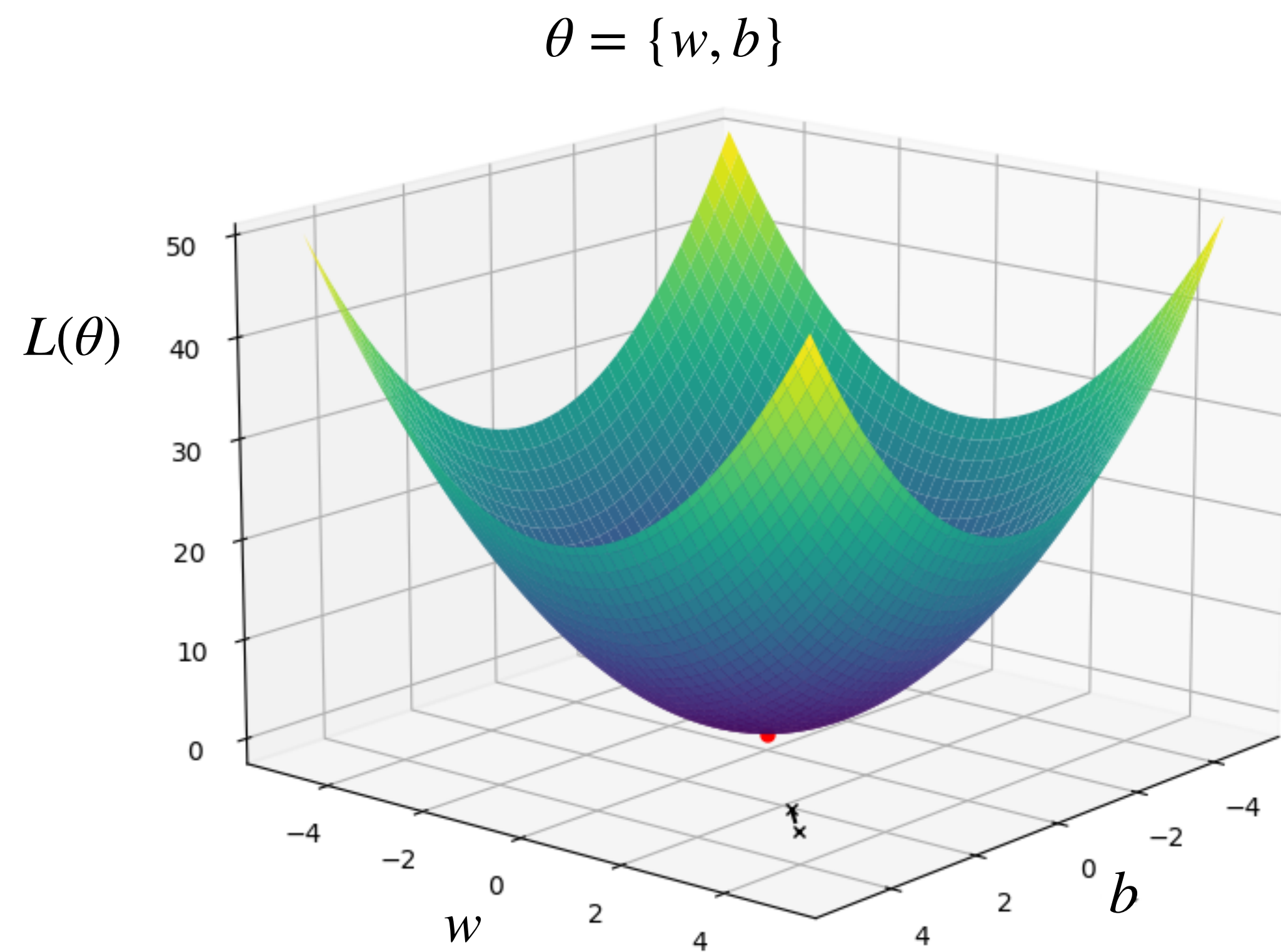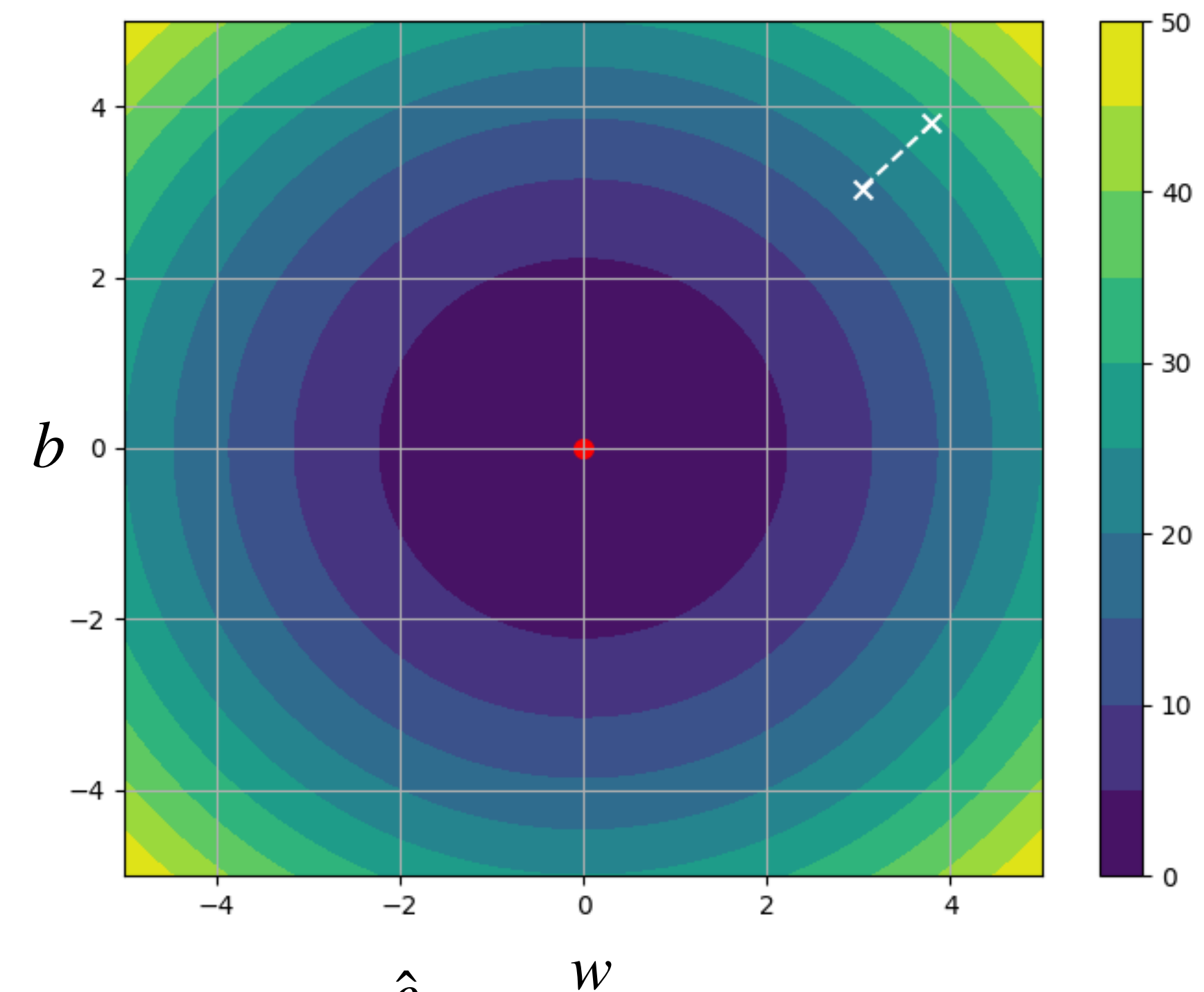


$$\theta_2 = \theta_1 - \alpha \nabla L(\theta_1)$$

•   $\hat{\theta}$

✖   $\theta_2 = (w_2, b_2)$

↖   $-\alpha \nabla L(\theta_1)$

# Loss Function Minima







Convex function ⟶ global minimum

Realistic loss functions are not-convex ⟶ search a local minimum

# Gradient Descent

Iterative optimization algorithm to find $\hat{\theta} = argmin_\theta L(\theta)$

Step 1. Compute the derivatives $\nabla L(\theta) = \left( \dfrac{\partial L}{\partial w_1}(\theta), \dfrac{\partial L}{\partial w_2}(\theta)\ldots \right)$

Step 2. Update the parameters $\theta \leftarrow \theta - \alpha \nabla L(\theta)$ $\longleftrightarrow$ $w_j \leftarrow w_j - \alpha \dfrac{\partial L}{\partial w_j}(\theta)$

Terminate if $\|\nabla L(\theta)\|$ is small enough

The *learning rate* $\alpha$ controls how fast we are changing the weights.

# Deterministic Gradient Descent

Iterative optimization algorithm to find $\hat{\theta} = argmin_\theta L(\theta)$

Step 1. Compute the derivatives $\nabla L(\theta) = \left( \dfrac{\partial L}{\partial w_1}(\theta), \dfrac{\partial L}{\partial w_2}(\theta) \ldots \right)$

**(Deterministic/full-batch)**

$$\frac{\partial L}{\partial w_j}(\theta) = \frac{\partial}{\partial w_j}\left( \frac{1}{n}\sum_{i=1}^{n} l^{(i)} \right)(\theta) = \frac{1}{n}\sum_{i=1}^{n}\frac{\partial l^{(i)}}{\partial w_j}(\theta)$$

# Deterministic Gradient Descent

Iterative optimization algorithm to find $\hat{\theta} = argmin_{\theta} L(\theta)$

Step 1. Compute the derivatives $\nabla L(\theta) = \left( \dfrac{\partial L}{\partial w_1}(\theta), \dfrac{\partial L}{\partial w_2}(\theta) \dots \right)$

**(Deterministic/full-batch)**

$$\frac{\partial L}{\partial w_j}(\theta) = \frac{\partial}{\partial w_j} \left( \frac{1}{n} \sum_{i=1}^{n} l^{(i)} \right)(\theta) = \frac{1}{n} \sum_{i=1}^{n} \frac{\partial l^{(i)}}{\partial w_j}(\theta)$$

Problems:

- One update might be time consuming and require a lot of memory for large datasets

- The local minimum found highly depends on the initial $\theta_0$

# Stochastic Gradient Descent

Iterative optimization algorithm to find $\hat{\theta} = argmin_\theta L(\theta)$

Step 1. Compute the derivatives $\nabla L(\theta) = \left( \dfrac{\partial L}{\partial w_1}(\theta), \dfrac{\partial L}{\partial w_2}(\theta) \ldots \right)$

**(Stochastic) Gradient Descent**

$$\frac{\partial L}{\partial w_j}(\theta) \approx \frac{1}{m} \sum_{k=1}^{m} \frac{\partial l^{(i_k)}}{\partial w_j}(\theta)$$

Randomly sample mini-batches with m data points $\{x^{(i_k)}, y^{(i_k)}\}_{k=1,\ldots,m}$

Problems:
- One update might be time consuming and require a lot of memory for large datasets
- The local minimum found highly depends on the initial $\theta_0$

58

# Stochastic Gradient Descent

Iterative optimization algorithm to find $\hat{\theta} = argmin_\theta L(\theta)$

Step 1. Compute the derivatives $\nabla L(\theta) = \left( \dfrac{\partial L}{\partial w_1}(\theta), \dfrac{\partial L}{\partial w_2}(\theta) \dots \right)$

**(Stochastic) Gradient Descent**

$$\frac{\partial L}{\partial w_j}(\theta) \approx \frac{1}{m} \sum_{k=1}^{m} \frac{\partial l^{(i_k)}}{\partial w_j}(\theta)$$

Randomly sample mini-batches with m data points $\{x^{(i_k)}, y^{(i_k)}\}_{k=1,\dots,m}$

How to compute the derivatives?

# Computing Derivatives

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

# Computing Derivatives

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$



$$\frac{\partial l}{\partial w} =$$

$$\partial l \left\{ \right.$$

$$\frac{\partial l}{\partial b} =$$

$$\left. \right\} \partial l$$

$$\partial w$$

$$\partial b$$

How sensitive is the loss to small changes of a specific parameter

# Computing Derivatives

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$



$$\frac{\partial l}{\partial w} = \frac{\partial}{\partial w}\left(\frac{1}{2}(y - \sigma(xw + b))^2\right)$$

$$= (y - \sigma(xw + b))\frac{\partial}{\partial w}\left(y - \sigma(xw + b)\right)$$

$$= -(y - \sigma(xw + b))\sigma'(xw + b)x$$

$$\frac{\partial l}{\partial b} = \frac{\partial}{\partial b}\left(\frac{1}{2}(y - \sigma(xw + b))^2\right)$$

$$= (y - \sigma(xw + b))\frac{\partial}{\partial b}\left(y - \sigma(xw + b)\right)$$

$$= -(y - \sigma(xw + b))\sigma'(xw + b)$$

Why?

# Computing Derivatives

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$



$$\frac{\partial l}{\partial w} = \frac{\partial}{\partial w}\left(\frac{1}{2}(y - \sigma(xw + b))^2\right)$$

$$= (y - \sigma(xw + b))\frac{\partial}{\partial w}\left(y - \sigma(xw + b)\right)$$

$$= -(y - \sigma(xw + b))\sigma'(xw + b)x$$

$$\frac{\partial l}{\partial b} = \frac{\partial}{\partial b}\left(\frac{1}{2}(y - \sigma(xw + b))^2\right)$$

$$= (y - \sigma(xw + b))\frac{\partial}{\partial b}\left(y - \sigma(xw + b)\right)$$

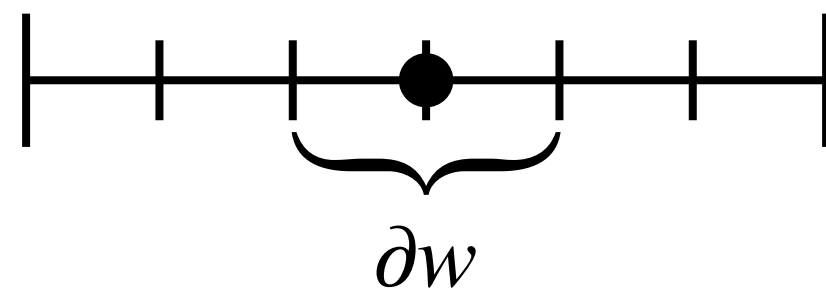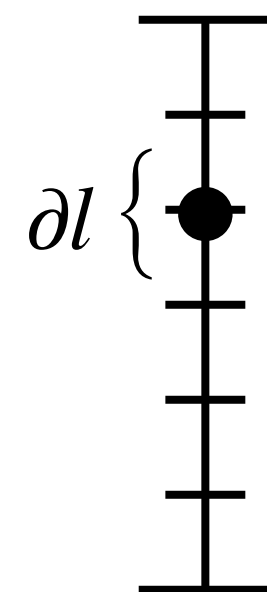$$= -(y - \sigma(xw + b))\sigma'(xw + b)$$

Chain Rule*

63

# Computing Derivatives

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$



$$\frac{\partial l}{\partial w} = \frac{\partial}{\partial w}\left(\frac{1}{2}(y - \sigma(xw + b))^2\right)$$

$$= (y - \sigma(xw + b))\frac{\partial}{\partial w}\left(y - \sigma(xw + b)\right)$$

$$= -(y - \sigma(xw + b))\sigma'(xw + b)x$$

$$\frac{\partial l}{\partial b} = \frac{\partial}{\partial b}\left(\frac{1}{2}(y - \sigma(xw + b))^2\right)$$

$$= (y - \sigma(xw + b))\frac{\partial}{\partial b}\left(y - \sigma(xw + b)\right)$$

$$= -(y - \sigma(xw + b))\sigma'(xw + b)$$

- Repeated computations

# Computing Derivatives

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$



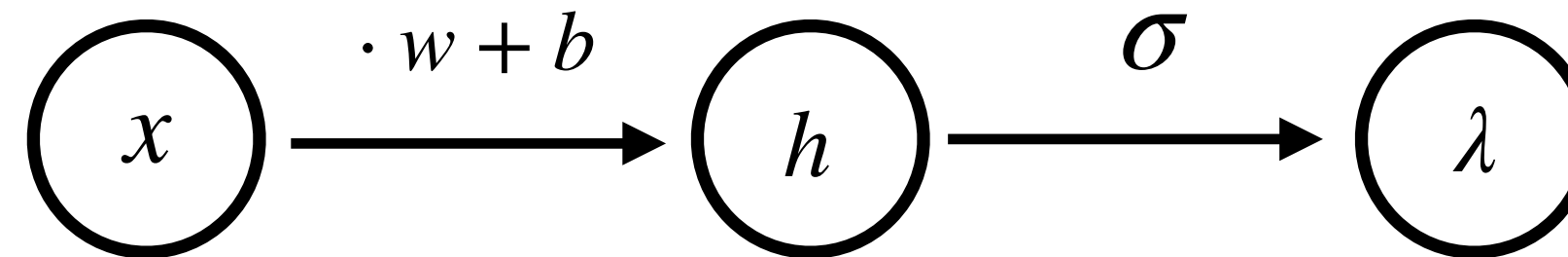$$\frac{\partial l}{\partial w} = \frac{\partial}{\partial w}\left(\frac{1}{2}(y - \sigma(xw + b))^2\right)$$

$$= (y - \sigma(xw + b))\frac{\partial}{\partial w}\left(y - \sigma(xw + b)\right)$$

$$= \boxed{-(y - \sigma(xw + b))\sigma'(xw + b)x}$$
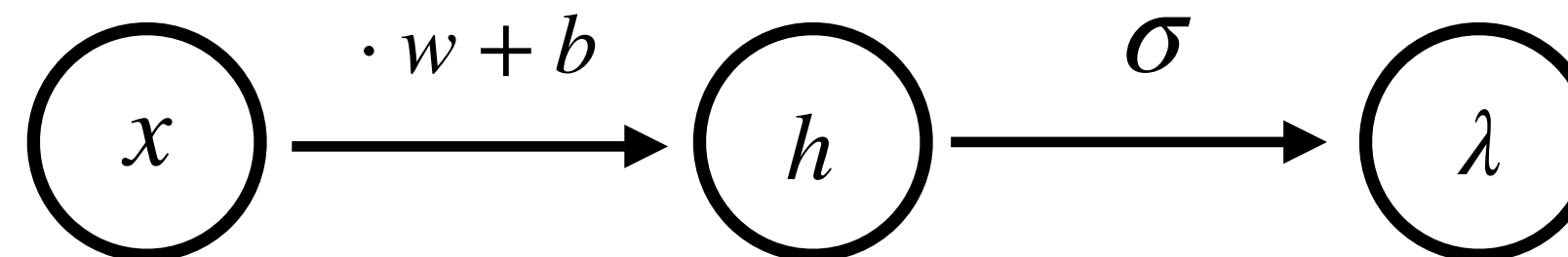
$$\frac{\partial l}{\partial b} = \frac{\partial}{\partial b}\left(\frac{1}{2}(y - \sigma(xw + b))^2\right)$$

$$= (y - \sigma(xw + b))\frac{\partial}{\partial b}\left(y - \sigma(xw + b)\right)$$

$$= \boxed{-(y - \sigma(xw + b))\sigma'(xw + b)}$$

- Repeated computations
- Identical terms

# Backpropagation

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$



Efficient way to use chain rule to compute derivatives

# Backpropagation

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

# Backpropagation

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

Forward



$$x^{(i)}$$

$$\lambda^{(i)} = 0.8$$

# Backpropagation

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

$x^{(i)}$

$\cdot w + b$

$\sigma$

0.3 → 0.67 → 0.8

$l^{(i)} = \frac{1}{2}(1 - 0.8)^2 = 0.02$

$\lambda^{(i)} = 0.8$

$y^{(i)} = 1$

# Backpropagation

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$



Forward

$$l^{(i)} = \frac{1}{2}(1 - 0.8)^2 = 0.02$$

$x^{(i)}$

$\lambda^{(i)} = 0.8$    $y^{(i)} = 1$



constant

$x$

$w$

$b$

$h$

constant

$y$

$\lambda$

$l$

*Computation graph*

# Backpropagation

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

**Forward pass equations**

$$l = \frac{1}{2}(y - \lambda)^2$$

$$\lambda = \sigma(h)$$

$$h = xw + b$$

Forward



$$l^{(i)} = \frac{1}{2}(1 - 0.8)^2 = 0.02$$

$$\lambda^{(i)} = 0.8 \qquad y^{(i)} = 1$$

constant $\quad$ constant

$x \qquad y$

$w \longrightarrow h \longrightarrow \lambda \longrightarrow l$

$b$

*Computation graph*

# Chain Rule in Neural Network

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

**Forward pass equations**

$$l = \frac{1}{2}(y - \lambda)^2$$

$$\lambda = \sigma(h)$$

$$h = xw + b$$

$$\frac{\partial l}{\partial w} = \cdots$$

$$\frac{\partial l}{\partial b} = \cdots$$

Forward

$\cdot w + b$     $\sigma$

$0.3$    $0.67$    $0.8$    $l^{(i)} = \frac{1}{2}(1 - 0.8)^2 = 0.02$

$x^{(i)}$

$\lambda^{(i)} = 0.8$     $y^{(i)} = 1$

$x$     $y$

$w$    $h$    $\lambda$    $l$    $\partial l$

$\partial w$

$b$

*Computation graph*

# Chain Rule in Neural Network

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

**Forward pass equations**

$$l = \frac{1}{2}(y - \lambda)^2$$

$$\lambda = \sigma(h)$$

$$h = xw + b$$

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial \lambda} \cdots$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial \lambda} \cdots$$

Forward

$$x^{(i)}$$

$$\cdot w + b \qquad \sigma$$

0.3 → 0.67 → 0.8

$$l^{(i)} = \frac{1}{2}(1 - 0.8)^2 = 0.02$$

$$\lambda^{(i)} = 0.8 \qquad y^{(i)} = 1$$

$x$

$w$

$b$

$h$

$y$

$\lambda$

$l$

$\} \partial \lambda$

$\} \partial l$

*Computation graph*

# Chain Rule in Neural Network

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

**Forward pass equations**

$$l = \frac{1}{2}(y - \lambda)^2$$

$$\lambda = \sigma(h)$$

$$h = xw + b$$

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\cdots$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\cdots$$

Forward



$$l^{(i)} = \frac{1}{2}(1 - 0.8)^2 = 0.02$$

$$\lambda^{(i)} = 0.8 \qquad y^{(i)} = 1$$

*Computation graph*

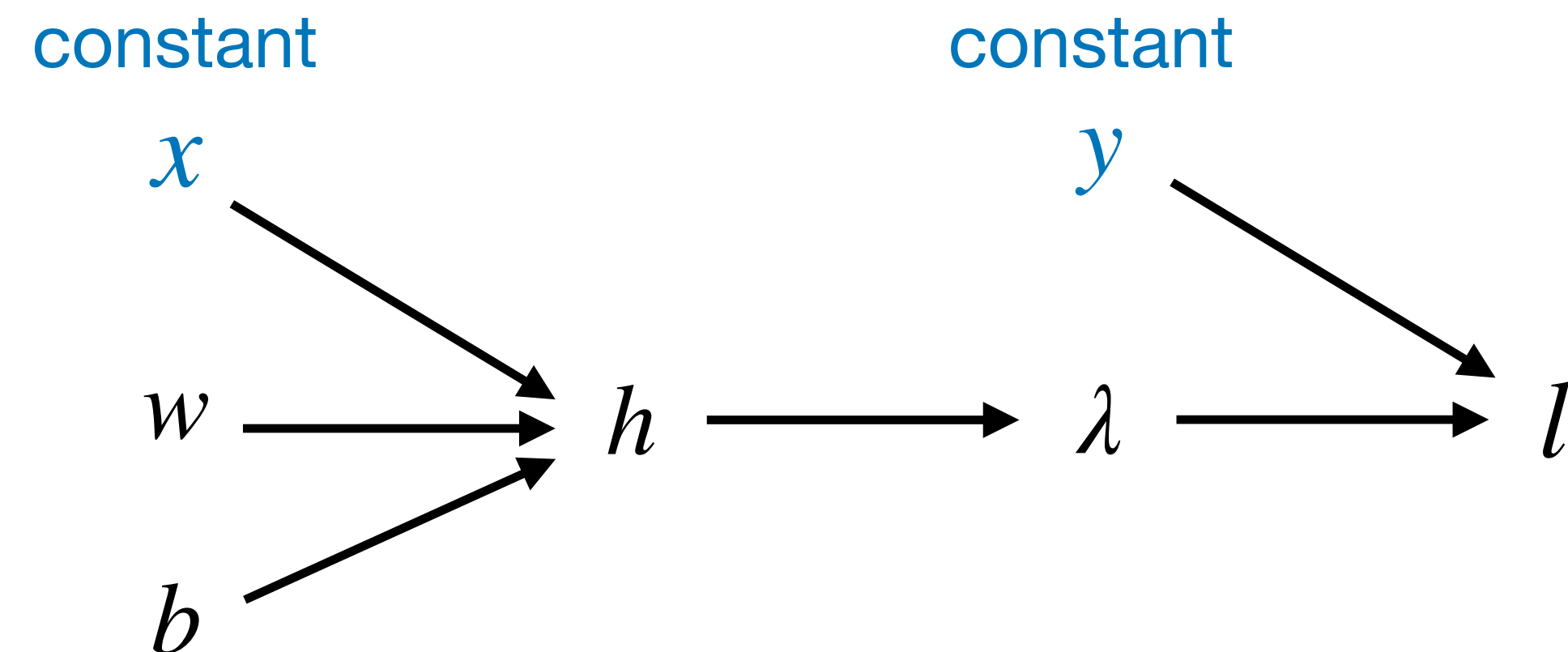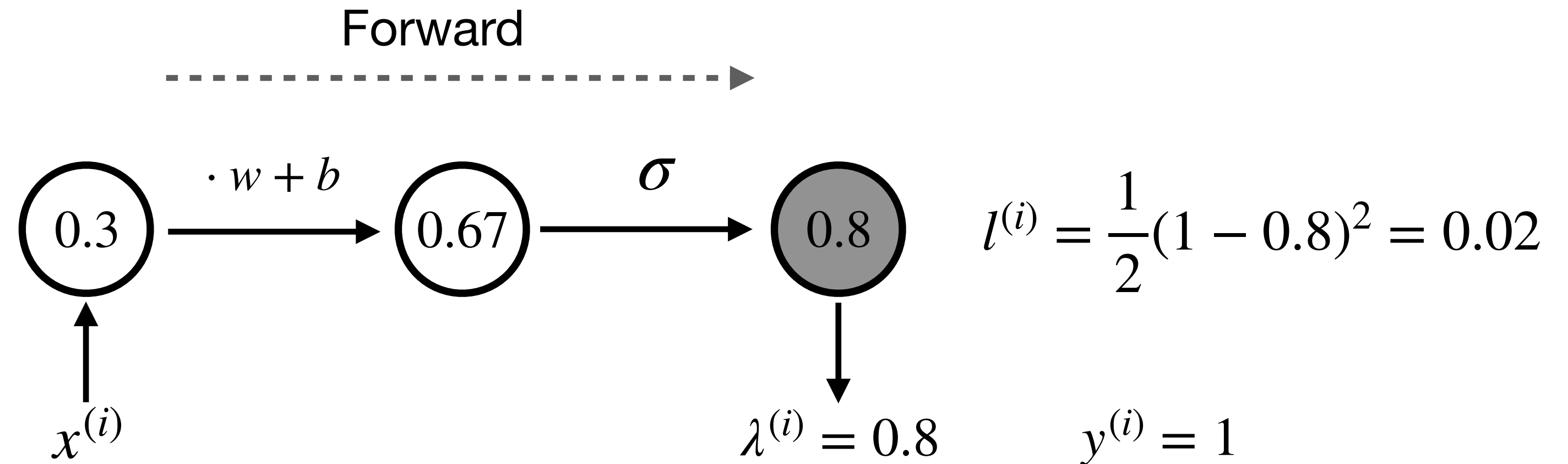# Chain Rule in Neural Network

**Example: univariate regression**

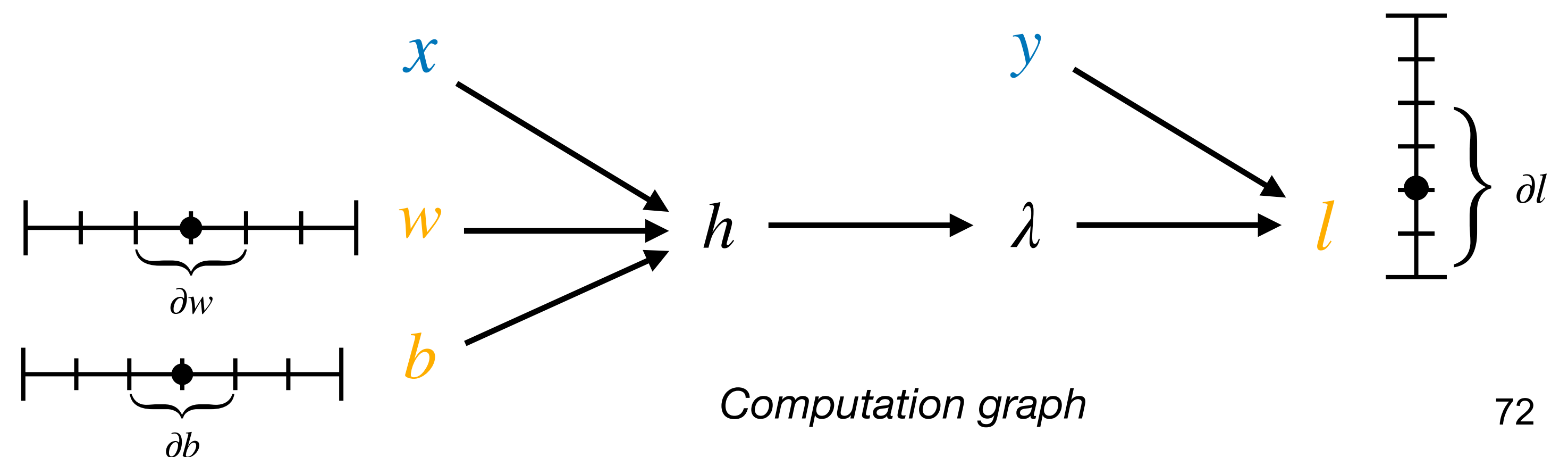$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

**Forward pass equations**

$$l = \frac{1}{2}(y - \lambda)^2$$

$$\lambda = \sigma(h)$$

$$h = xw + b$$

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\frac{\partial h}{\partial w}$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\frac{\partial h}{\partial b}$$

Forward

$$\cdot w + b \qquad \sigma$$

$0.3 \qquad 0.67 \qquad 0.8 \qquad l^{(i)} = \frac{1}{2}(1 - 0.8)^2 = 0.02$

$x^{(i)}$

$\lambda^{(i)} = 0.8 \qquad y^{(i)} = 1$

$\partial h$

$x$

$y$

$w \qquad h \qquad \lambda \qquad l$

$\partial w$

$b$

*Computation graph*

# Chain Rule in Neural Network

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

$l^{(i)} = \frac{1}{2}(1 - 0.8)^2 = 0.02$

0.3 $\xrightarrow{\cdot w + b}$ 0.67 $\xrightarrow{\sigma}$ 0.8

$x^{(i)}$

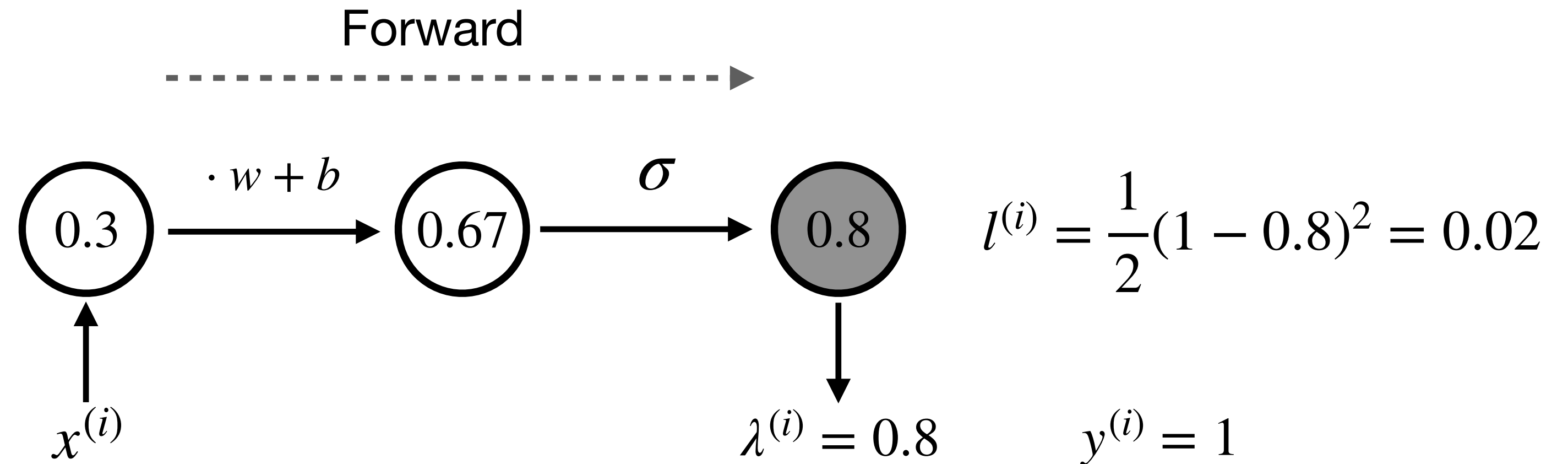$\lambda^{(i)} = 0.8$      $y^{(i)} = 1$
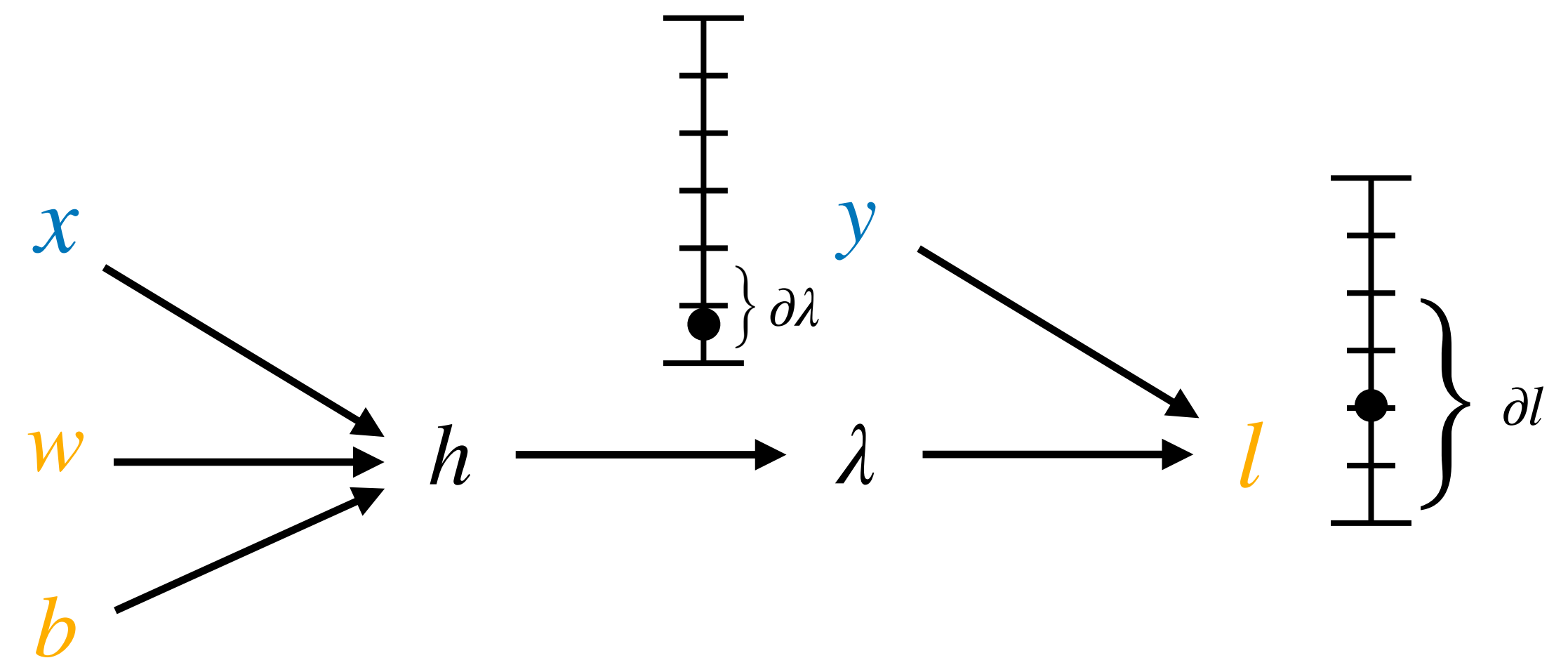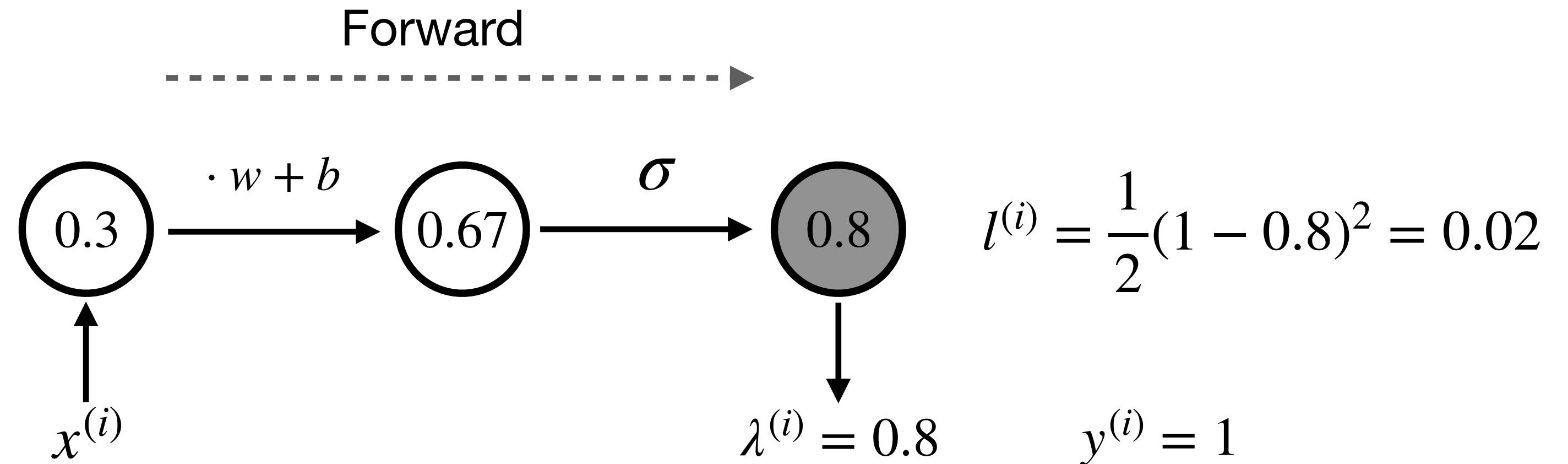
**Forward pass equations**

$$l = \frac{1}{2}(y - \lambda)^2$$

$$\lambda = \sigma(h)$$

$$h = xw + b$$

$$\frac{\partial l}{\partial w} = \boxed{\frac{\partial l}{\partial \lambda} \frac{\partial \lambda}{\partial h} \frac{\partial h}{\partial w}}$$

3 derivatives

+

?

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial \lambda} \frac{\partial \lambda}{\partial h} \frac{\partial h}{\partial b}$$

$x$          $y$

$w \longrightarrow h \longrightarrow \lambda \longrightarrow l$

$b$

*Computation graph*

76

# Chain Rule in Neural Network

**Example: univariate regression**

Forward

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$



$$l^{(i)} = \frac{1}{2}(1 - 0.8)^2 = 0.02$$

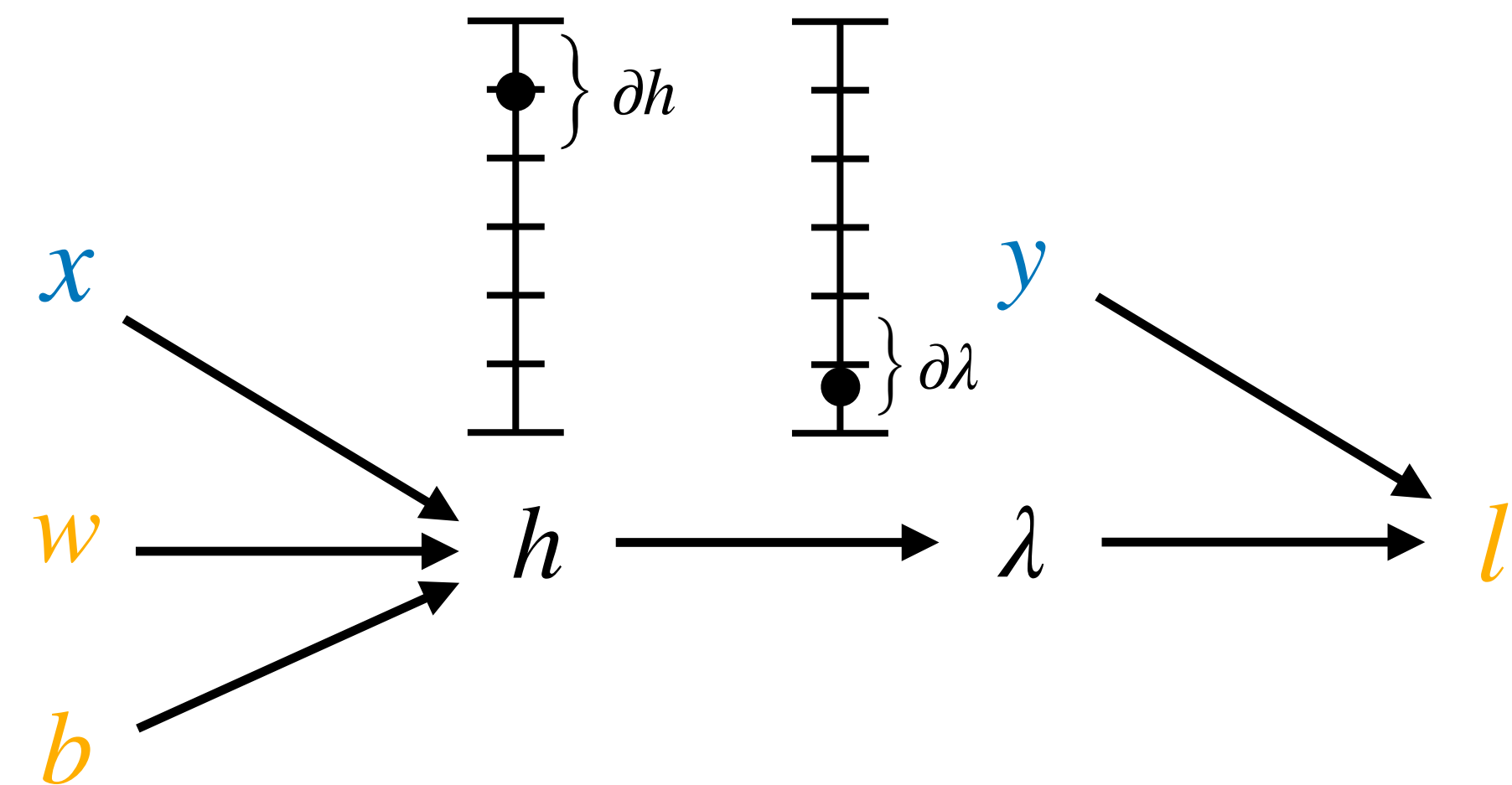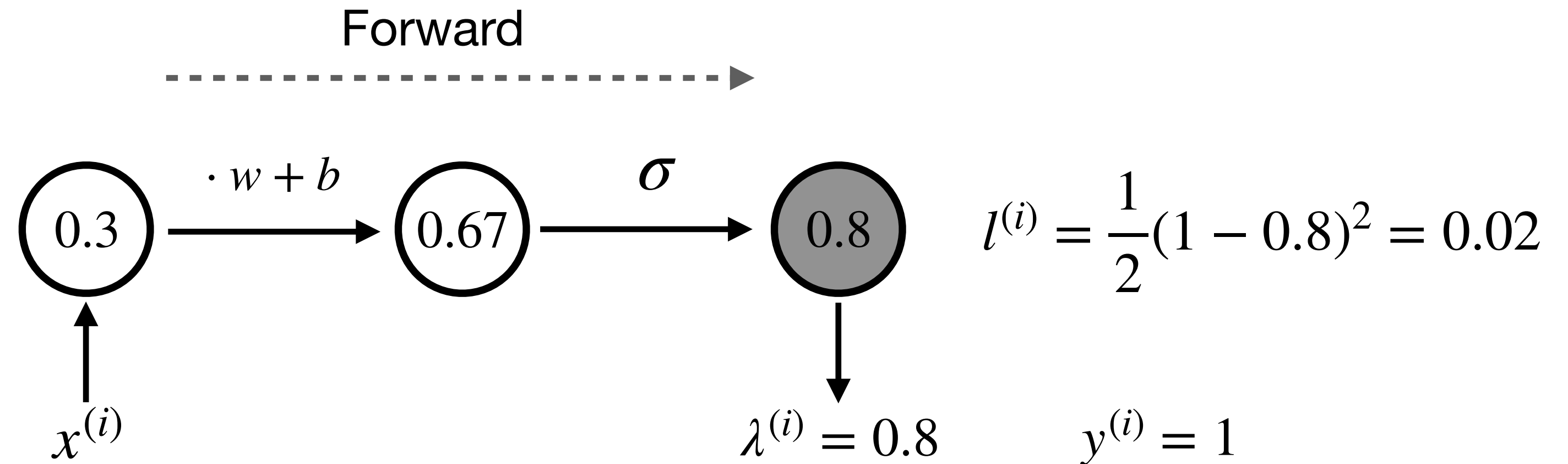$\lambda^{(i)} = 0.8$      $y^{(i)} = 1$

**Forward pass equations**

$$l = \frac{1}{2}(y - \lambda)^2$$

$$\lambda = \sigma(h)$$

$$h = xw + b$$

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\frac{\partial h}{\partial w}$$

$$=$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\boxed{\frac{\partial h}{\partial b}}$$

3 derivatives

+

1 derivative

*Computation graph*

# Forward Pass

**Example: univariate regression**

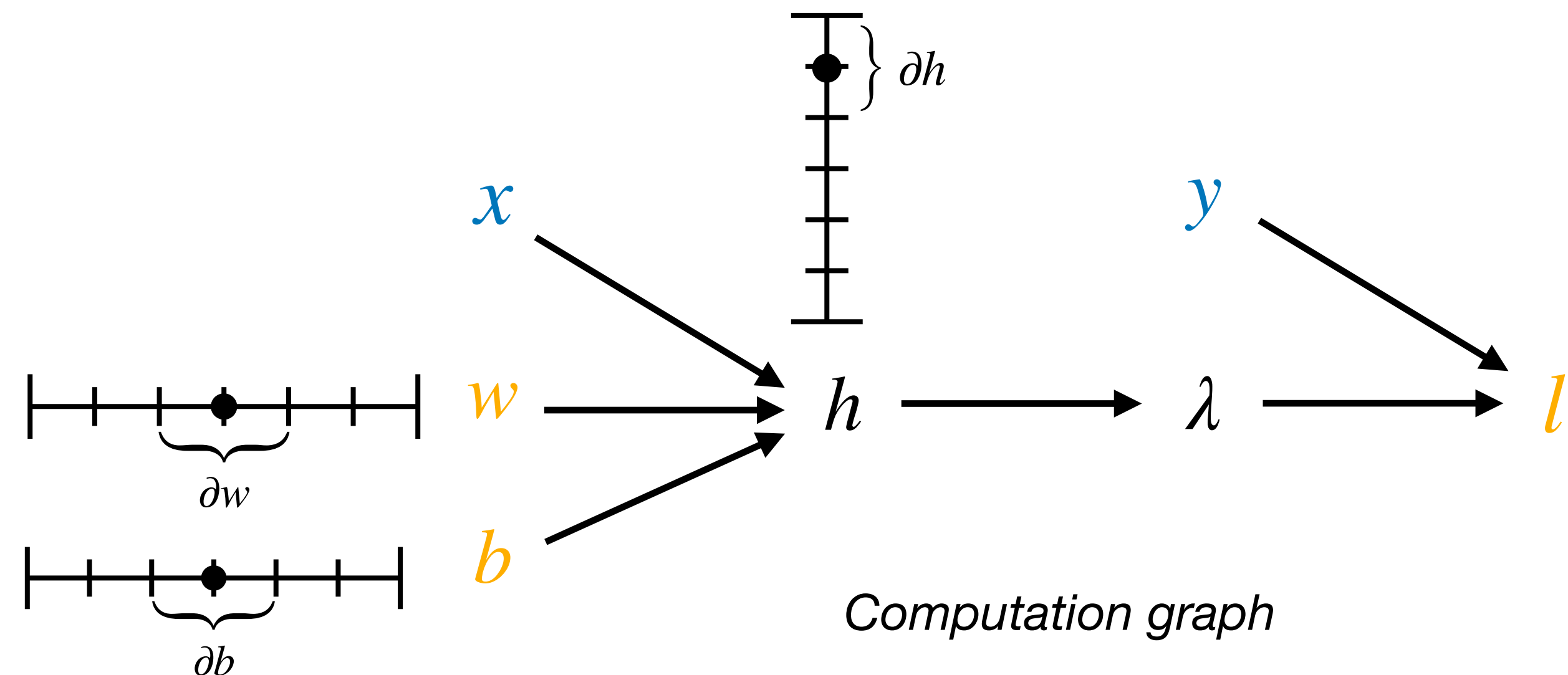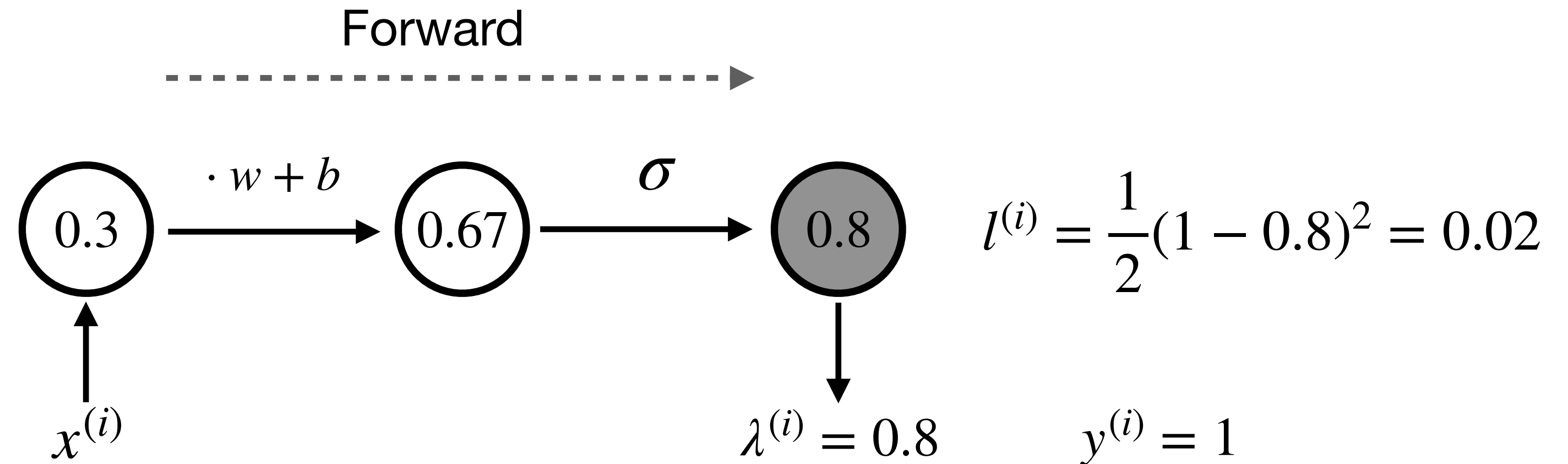$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

**Forward pass equations**

$$l = \frac{1}{2}(y - \lambda)^2$$

$$\lambda = \sigma(h)$$

$$h = xw + b$$

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\frac{\partial h}{\partial w}$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\frac{\partial h}{\partial b}$$

Forward

$$\cdot\, w + b \qquad \sigma$$

$0.3 \longrightarrow 0.67 \longrightarrow 0.8 \qquad l^{(i)} = \frac{1}{2}(1 - 0.8)^2 = 0.02$

$x^{(i)}$

$\lambda^{(i)} = 0.8 \qquad y^{(i)} = 1$

Stored
$$(x^{(i)}, h^{(i)}, \lambda^{(i)}, y^{(i)}, l^{(i)}) \leftarrow (0.3, 0.67, 0.8, 1, 0.02)$$

$$x \quad y$$
$$w \longrightarrow h \longrightarrow \lambda \longrightarrow l$$
$$b$$

*Computation graph*

# Backward Pass

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

**Forward pass equations**   **Backward pass equations**

$$l = \frac{1}{2}(y - \lambda)^2$$

$$\lambda = \sigma(h)$$

$$h = xw + b$$

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial \lambda} \frac{\partial \lambda}{\partial h} \frac{\partial h}{\partial w}$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial \lambda} \frac{\partial \lambda}{\partial h} \frac{\partial h}{\partial b}$$

Backward



Stored
$$(x^{(i)}, h^{(i)}, \lambda^{(i)}, y^{(i)}, l^{(i)}) \leftarrow (0.3, 0.67, 0.8, 1, 0.02)$$

*Computation graph*

# Backward Pass

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

**Forward pass equations**

$$l = \frac{1}{2}(y - \lambda)^2$$
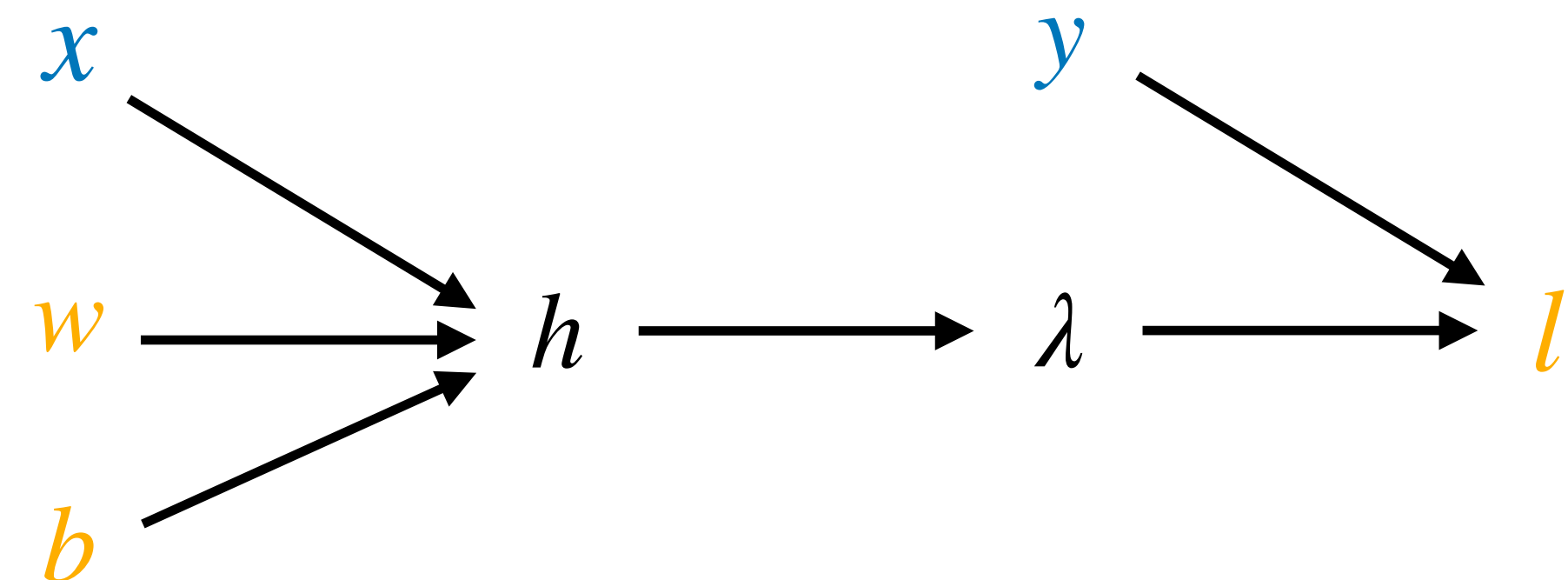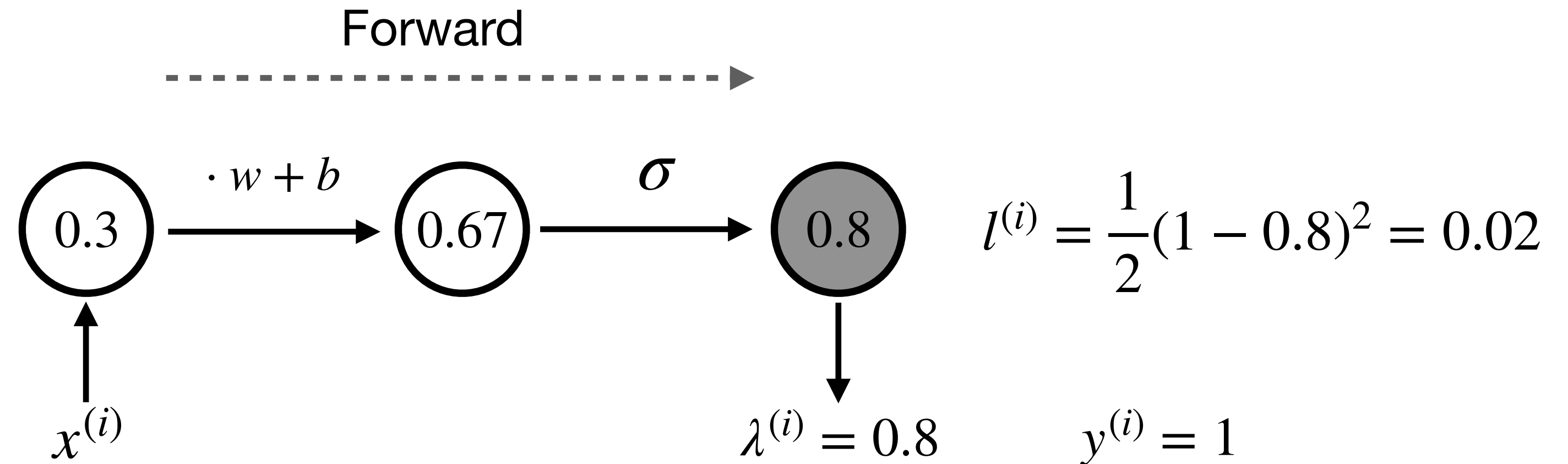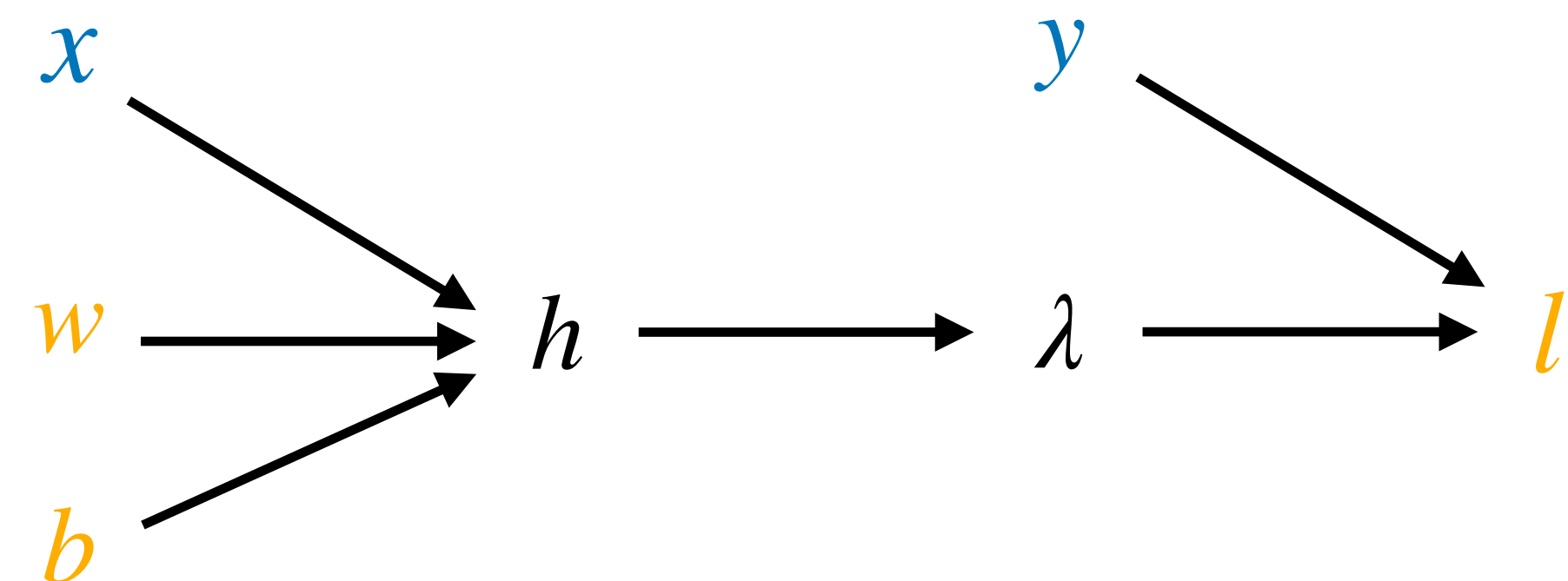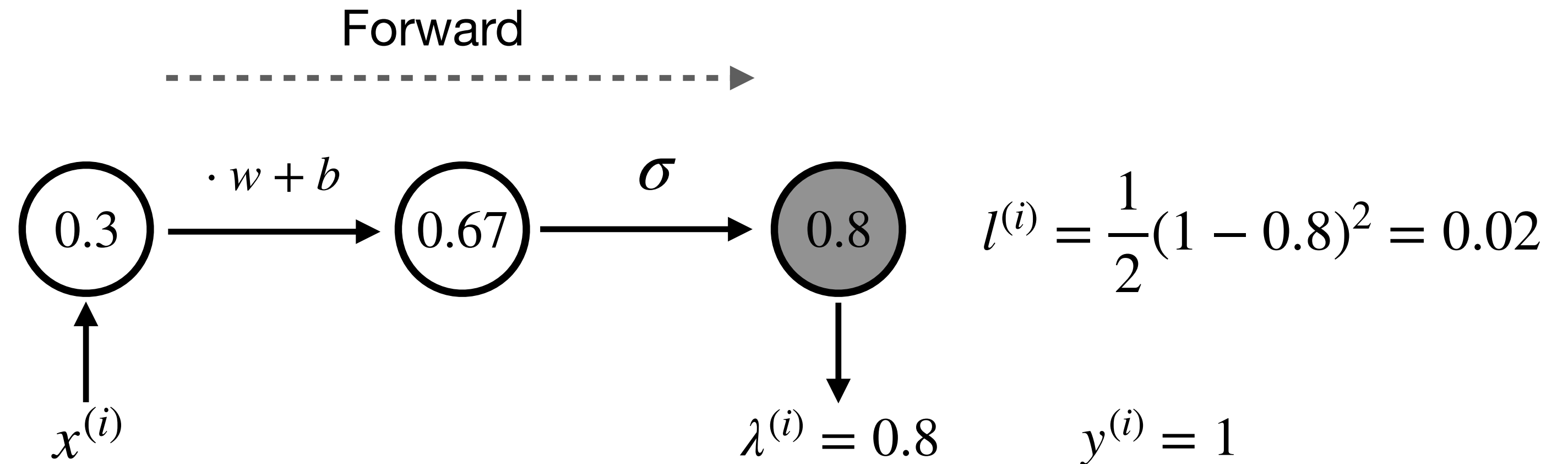
$$\lambda = \sigma(h)$$

$$h = xw + b$$

**Backward pass equations**

$$\frac{\partial l}{\partial \lambda} = (\lambda - y)$$

Stored
$$(x^{(i)}, h^{(i)}, \lambda^{(i)}, y^{(i)}, l^{(i)}) \leftarrow (0.3, 0.67, 0.8, 1, 0.02)$$

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial \lambda} \frac{\partial \lambda}{\partial h} \frac{\partial h}{\partial w}$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial \lambda} \frac{\partial \lambda}{\partial h} \frac{\partial h}{\partial b}$$

*Computation graph*

80

# Backward Pass

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

Backward



**Forward pass equations**

$$l = \frac{1}{2}(y - \lambda)^2$$

$$\lambda = \sigma(h)$$

$$h = xw + b$$

**Backward pass equations**

$$\frac{\partial l}{\partial \lambda} = (\lambda - y)$$

Stored
$$(x^{(i)}, h^{(i)}, \lambda^{(i)}, y^{(i)}, l^{(i)}) \leftarrow (0.3, 0.67, 0.8, 1, 0.02)$$

$$\frac{\partial l^{(i)}}{\partial \lambda} \leftarrow (\lambda^{(i)} - y^{(i)}) = -0.2$$

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial \lambda} \frac{\partial \lambda}{\partial h} \frac{\partial h}{\partial w}$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial \lambda} \frac{\partial \lambda}{\partial h} \frac{\partial h}{\partial b}$$

*Computation graph*

81

# Backward Pass

**Example: univariate regression**

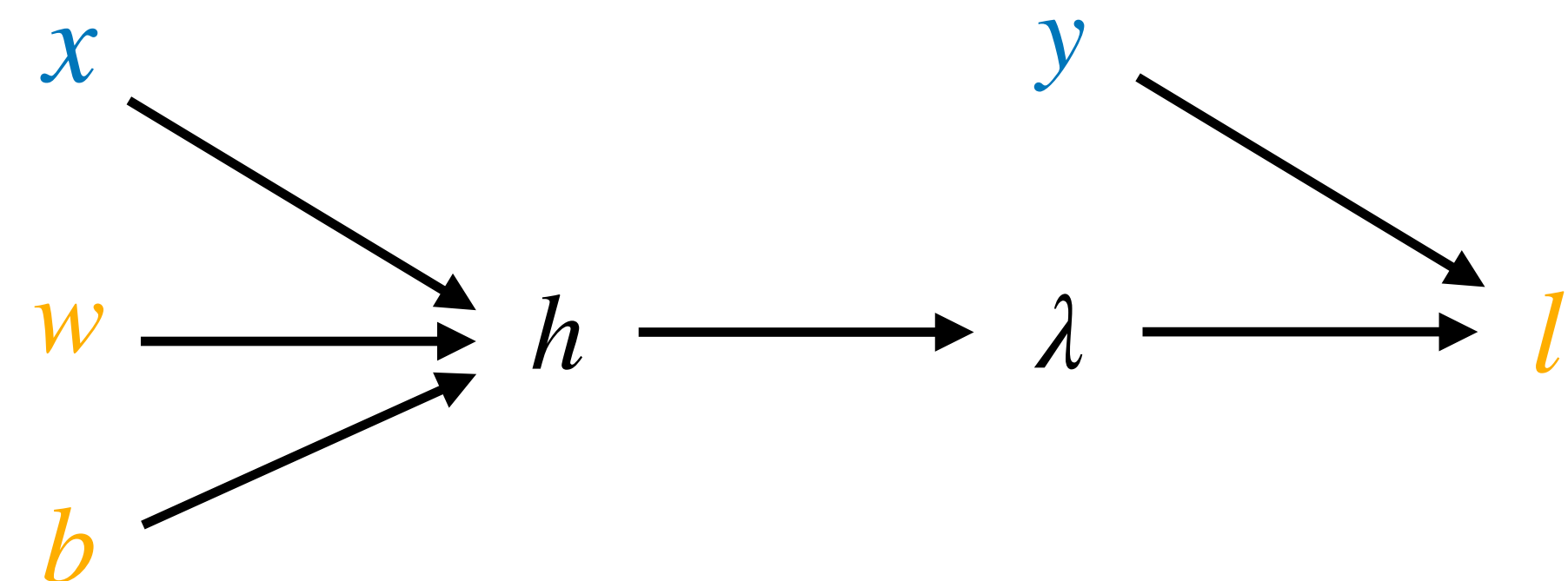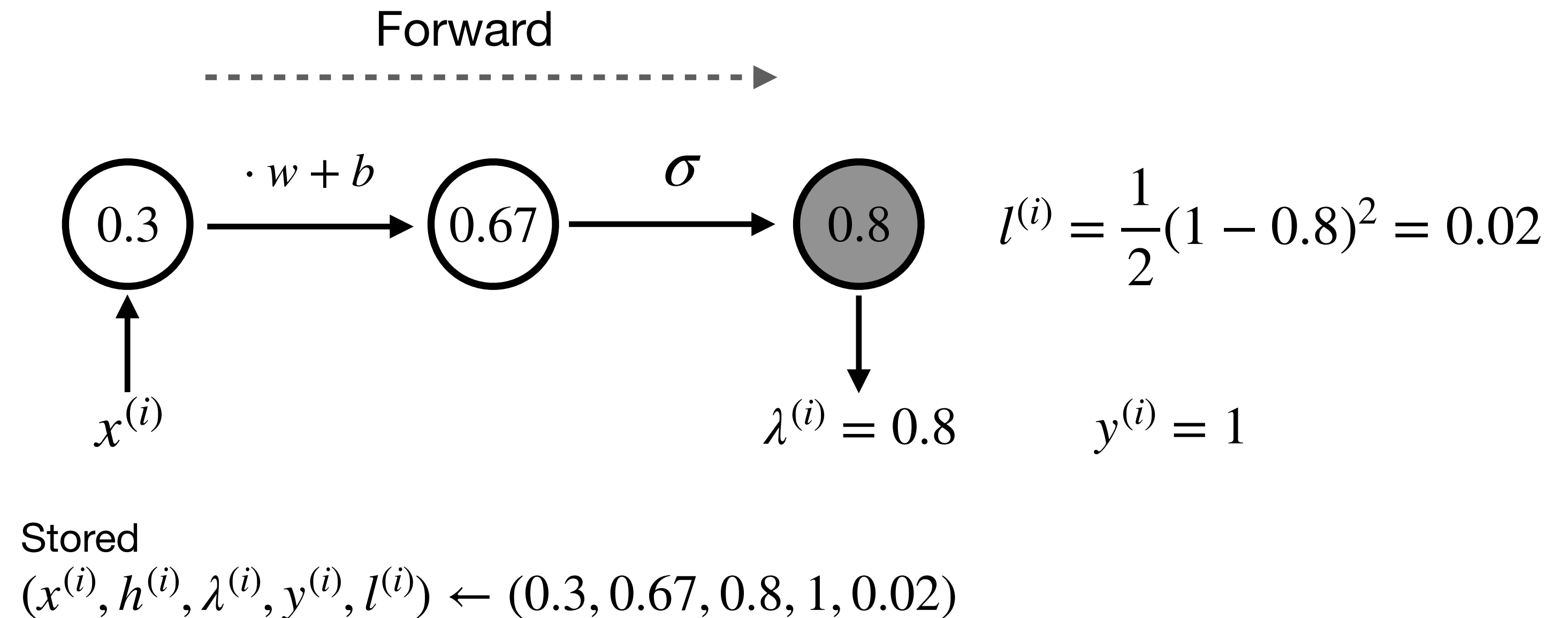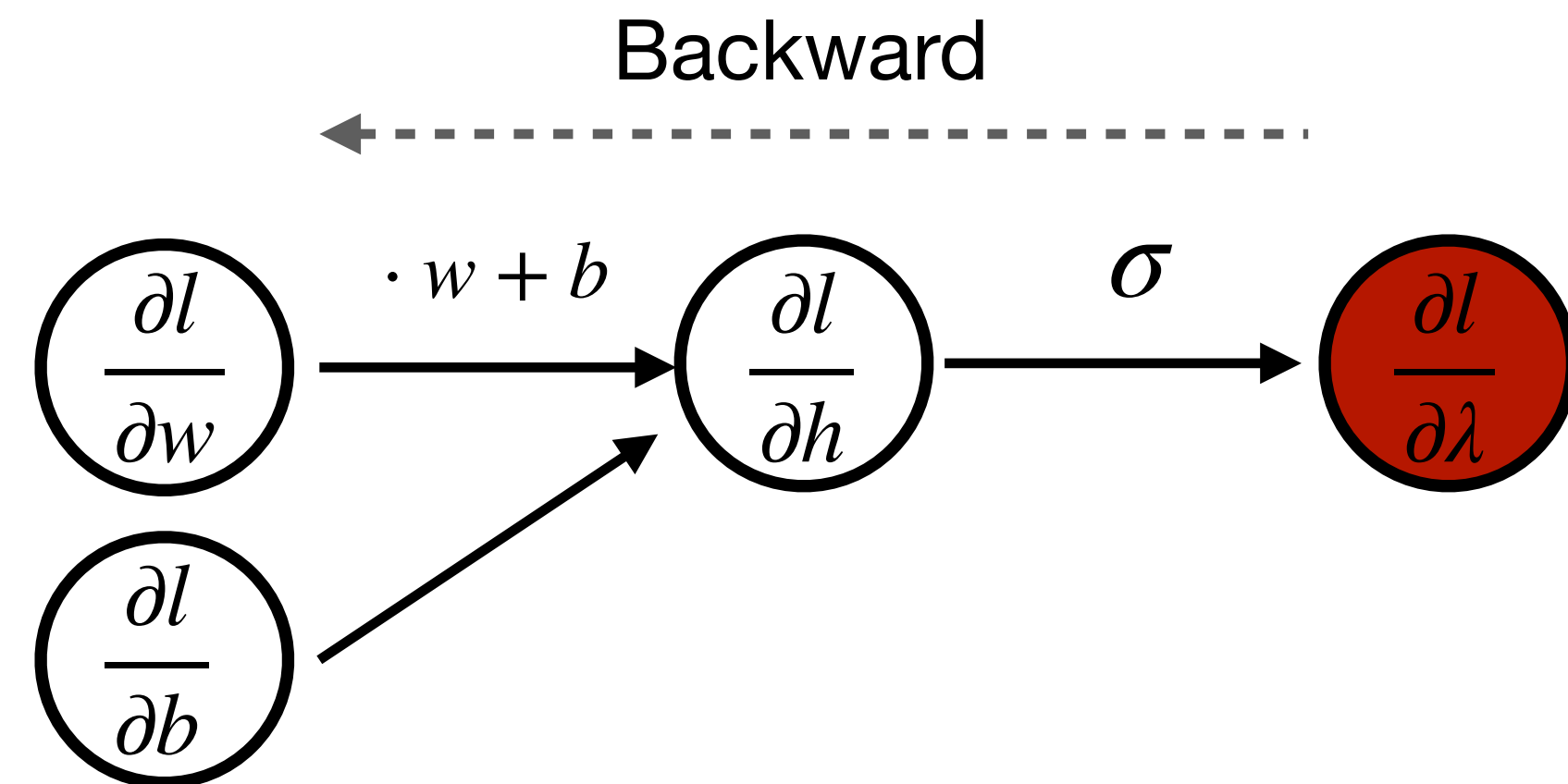$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

**Forward pass equations**

$$l = \frac{1}{2}(y - \lambda)^2$$

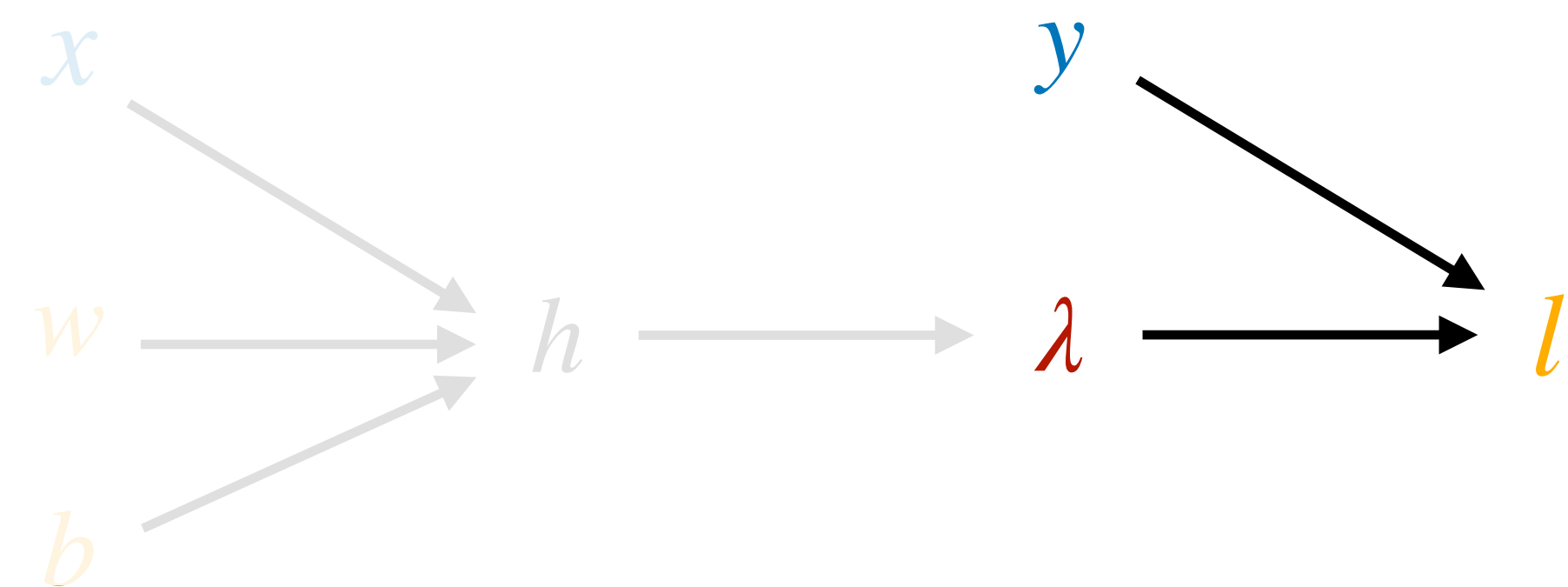$$\lambda = \sigma(h)$$

$$h = xw + b$$

**Backward pass equations**

$$\frac{\partial l}{\partial \lambda} = (\lambda - y)$$

$$\frac{\partial l}{\partial h} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}$$

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\frac{\partial h}{\partial w}$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\frac{\partial h}{\partial b}$$

Backward



Stored
$$(x^{(i)}, h^{(i)}, \lambda^{(i)}, y^{(i)}, l^{(i)}) \leftarrow (0.3, 0.67, 0.8, 1, 0.02)$$

$$\frac{\partial l^{(i)}}{\partial \lambda} \leftarrow -0.2$$

*Computation graph*

# Backward Pass

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$
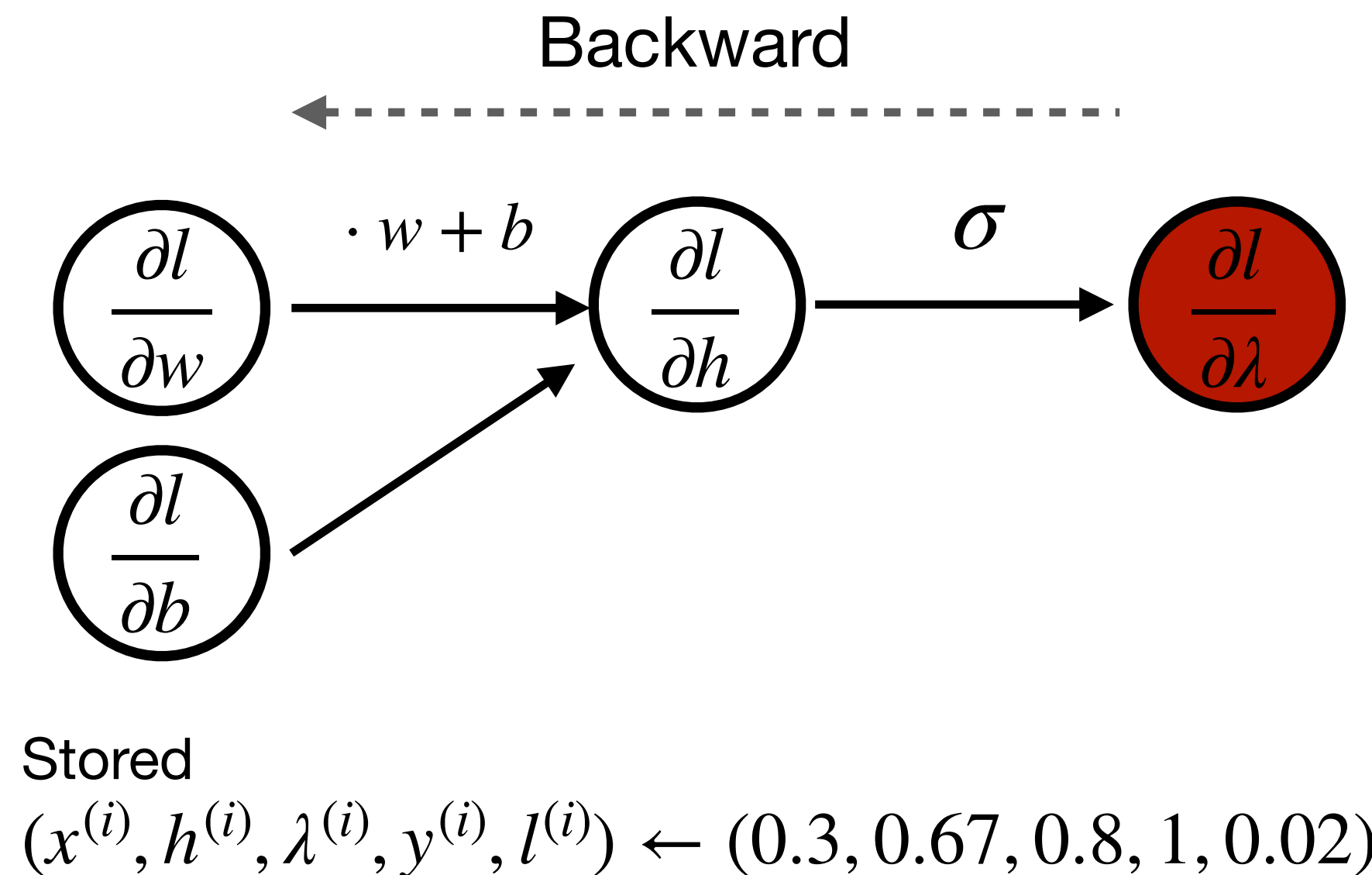
**Forward pass equations**

$$l = \frac{1}{2}(y - \lambda)^2$$
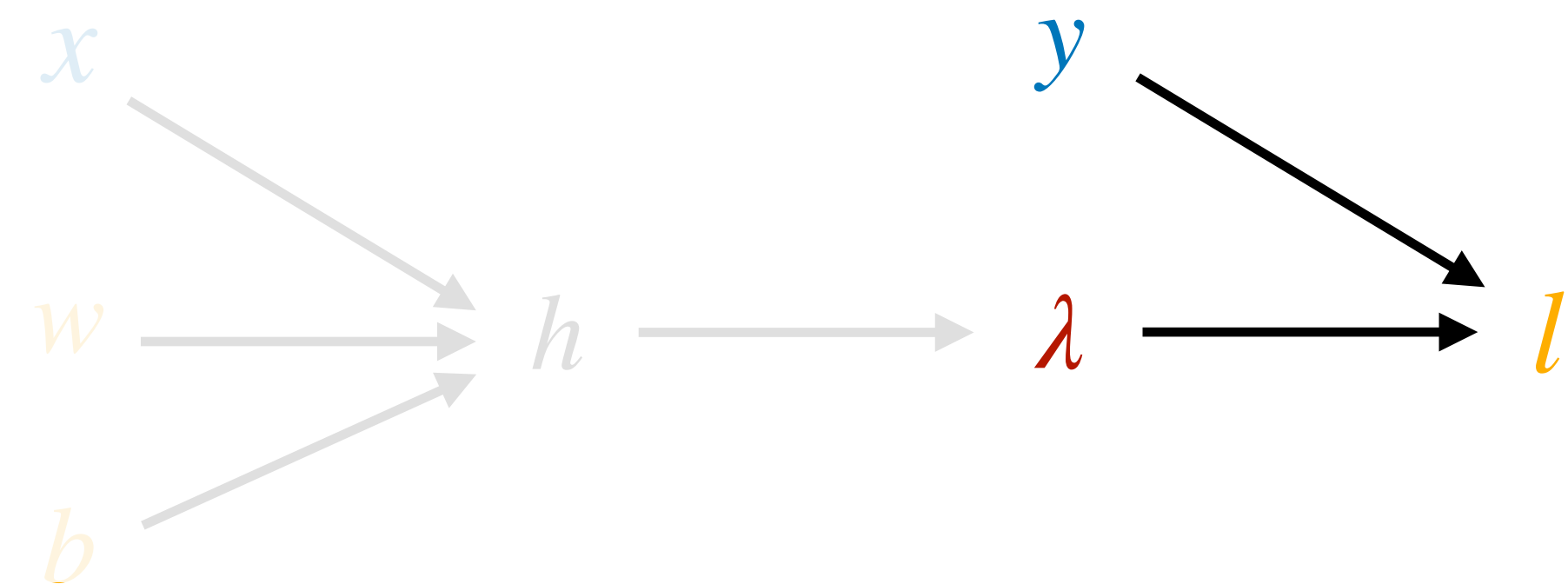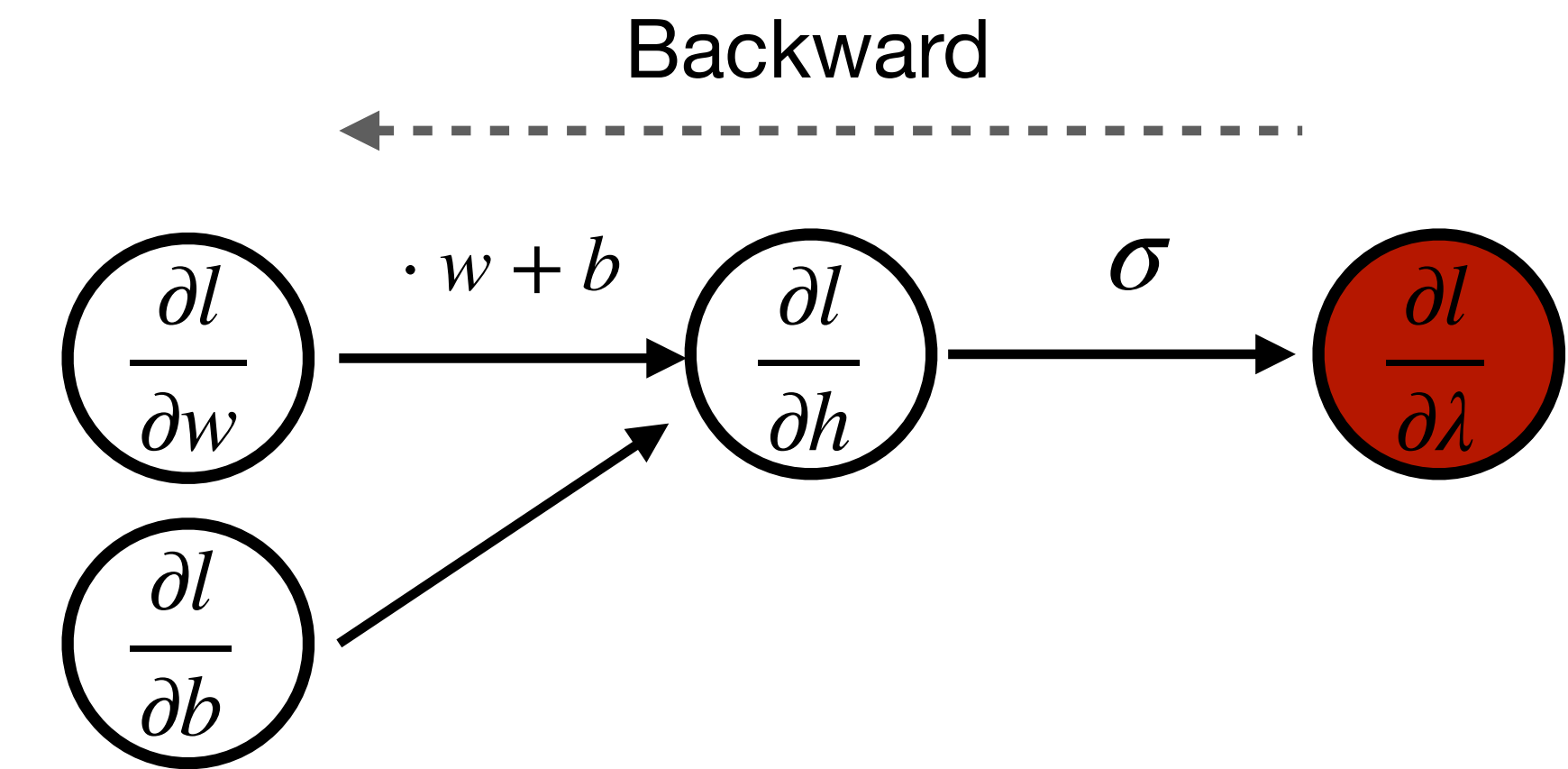
$$\lambda = \sigma(h)$$

$$h = xw + b$$

**Backward pass equations**

$$\frac{\partial l}{\partial \lambda} = (\lambda - y)$$

$$\frac{\partial l}{\partial h} = \frac{\partial l}{\partial \lambda}\sigma'(h)$$

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\frac{\partial h}{\partial w}$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\frac{\partial h}{\partial b}$$

Backward

$$\frac{\partial l}{\partial w} \quad \cdot w + b \quad \frac{\partial l}{\partial h} \quad \sigma \quad -0.2$$

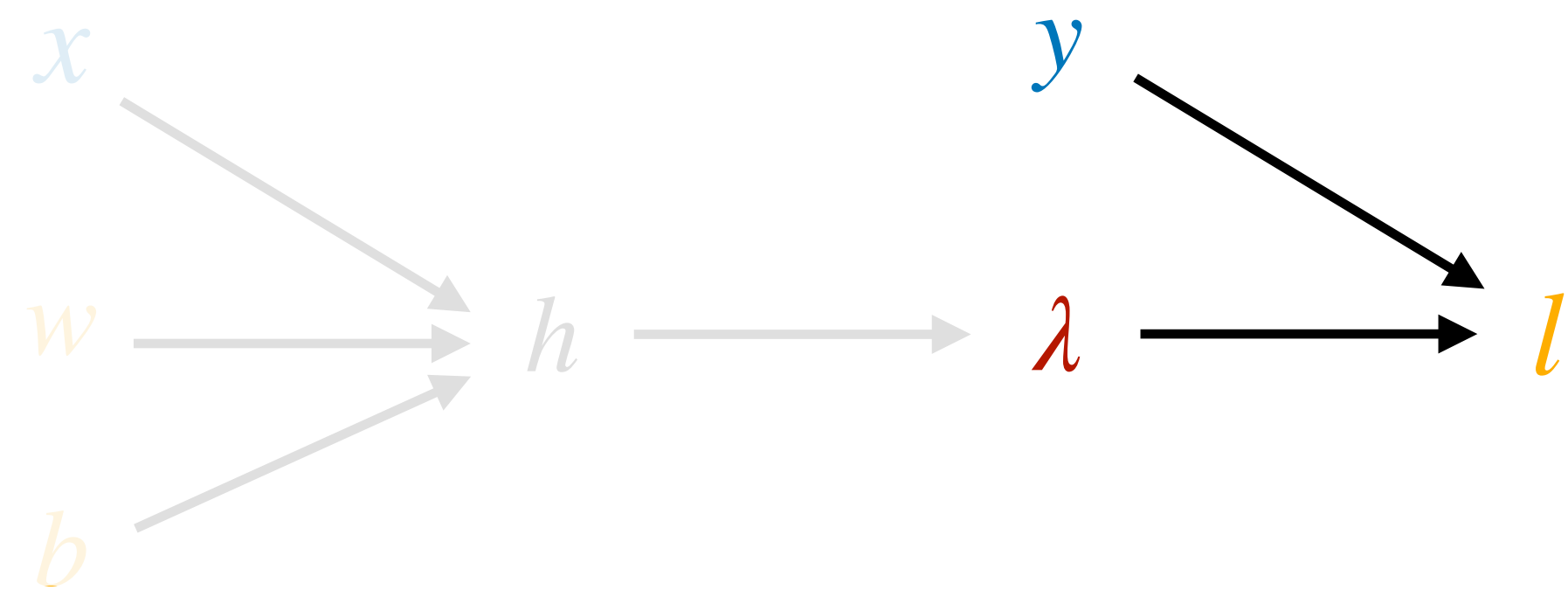$$\frac{\partial l}{\partial b}$$

Stored
$$(x^{(i)}, h^{(i)}, \lambda^{(i)}, y^{(i)}, l^{(i)}) \leftarrow (0.3, 0.67, 0.8, 1, 0.02)$$

$$\frac{\partial l^{(i)}}{\partial \lambda} \leftarrow -0.2$$

$$x \quad\quad y$$

$$w \quad h \quad \lambda \quad l$$

$$b$$

*Computation graph*

# Backward Pass

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

**Forward pass equations**

$$l = \frac{1}{2}(y - \lambda)^2$$

$$\lambda = \sigma(h)$$

$$h = xw + b$$

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\frac{\partial h}{\partial w}$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\frac{\partial h}{\partial b}$$
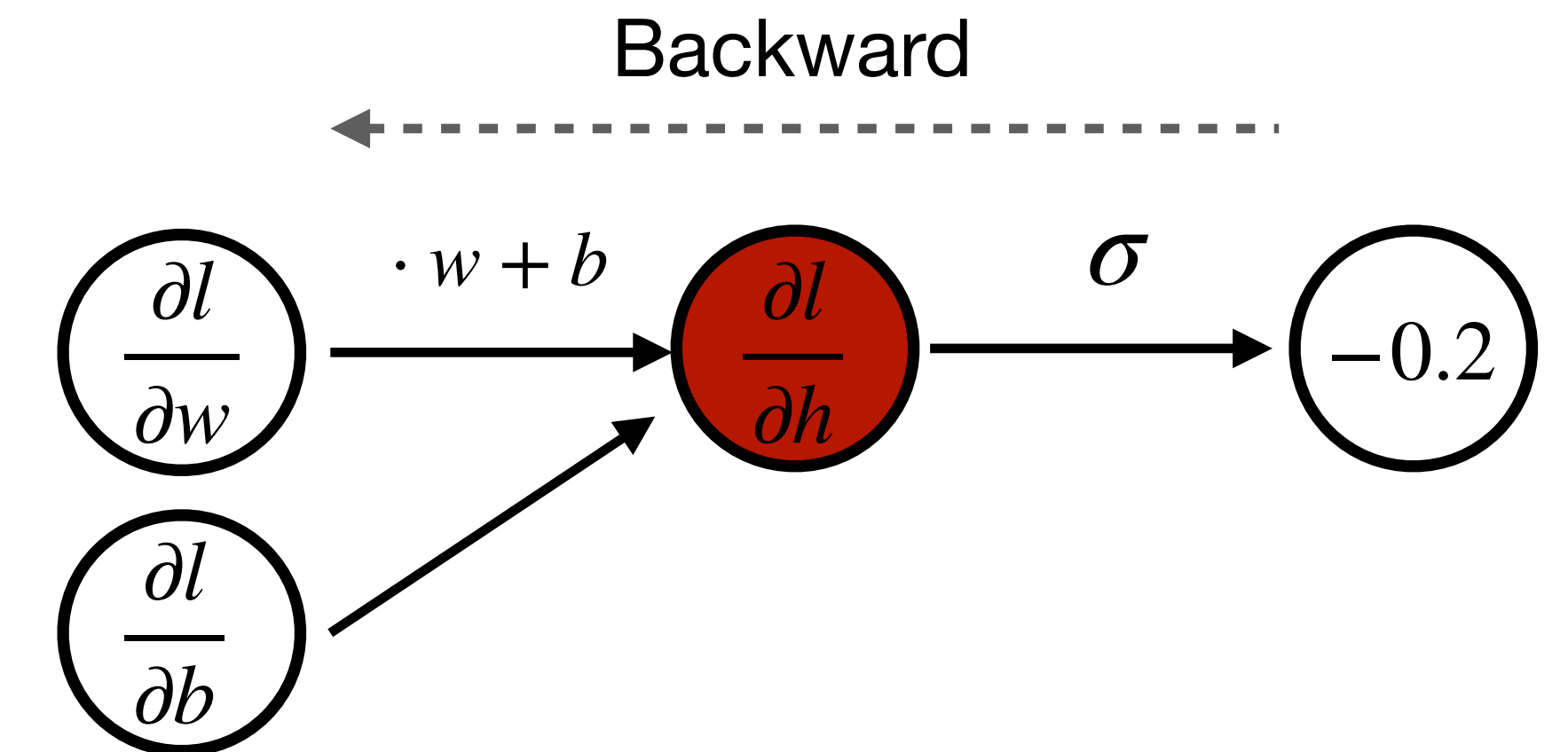
**Backward pass equations**

$$\frac{\partial l}{\partial \lambda} = (\lambda - y)$$

$$\frac{\partial l}{\partial h} = \frac{\partial l}{\partial \lambda}\lambda(1 - \lambda)$$

Backward



Stored

$$(x^{(i)}, h^{(i)}, \lambda^{(i)}, y^{(i)}, l^{(i)}) \leftarrow (0.3, 0.67, 0.8, 1, 0.02)$$

$$\frac{\partial l^{(i)}}{\partial \lambda} \leftarrow -0.2$$

$$\frac{\partial l^{(i)}}{\partial h} \leftarrow \frac{\partial l^{(i)}}{\partial \lambda}\lambda^{(i)}(1 - \lambda^{(i)}) = -0.04$$

*Computation graph*

# Backward Pass

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

**Forward pass equations**

$$l = \frac{1}{2}(y - \lambda)^2$$

$$\lambda = \sigma(h)$$

$$h = xw + b$$

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\frac{\partial h}{\partial w}$$

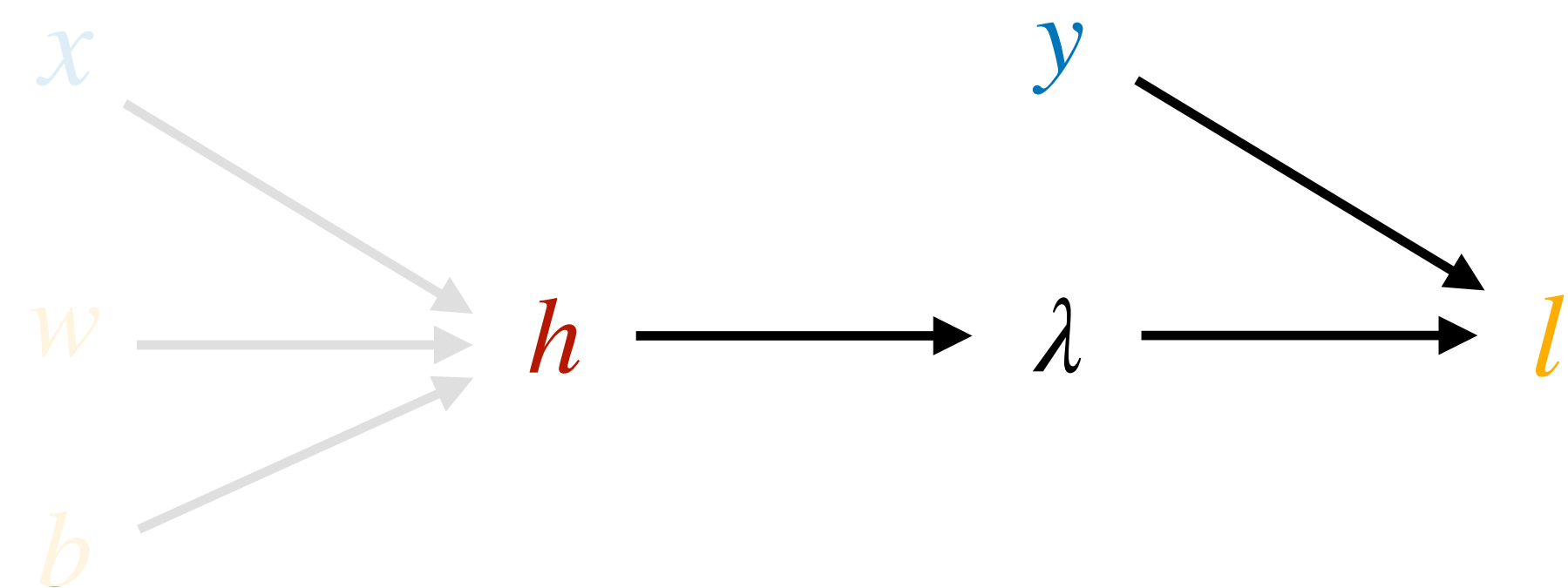$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\frac{\partial h}{\partial b}$$
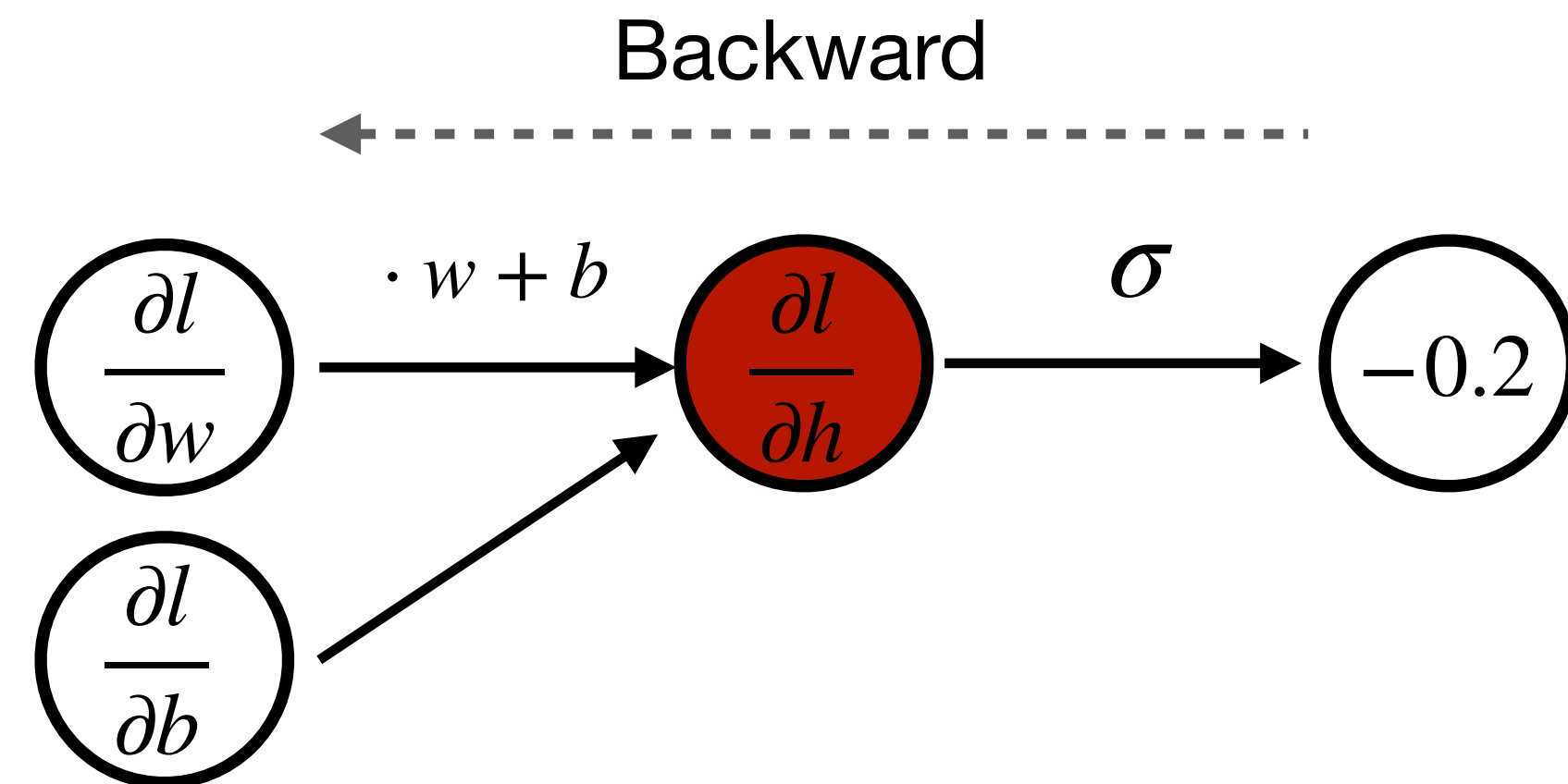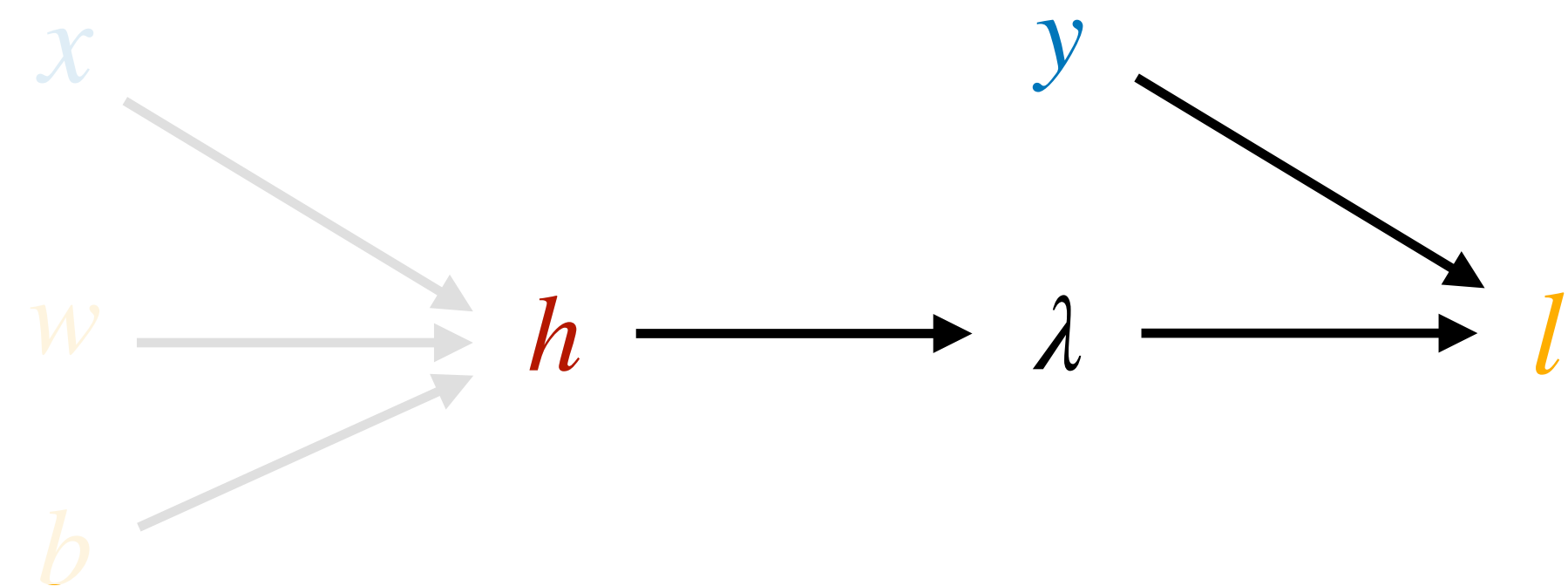
**Backward pass equations**

$$\frac{\partial l}{\partial \lambda} = (\lambda - y)$$

$$\frac{\partial l}{\partial h} = \frac{\partial l}{\partial \lambda}\lambda(1 - \lambda)$$

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial h}\frac{\partial h}{\partial w}$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial h}\frac{\partial h}{\partial b}$$

Backward

$$\frac{\partial l}{\partial w}$$  $\cdot w + b$  $-0.04$  $\sigma$  $-0.2$

$$\frac{\partial l}{\partial b}$$

Stored

$$(x^{(i)}, h^{(i)}, \lambda^{(i)}, y^{(i)}, l^{(i)}) \leftarrow (0.3, 0.67, 0.8, 1, 0.02)$$

$$\frac{\partial l^{(i)}}{\partial \lambda} \leftarrow -0.2$$

$$\frac{\partial l^{(i)}}{\partial h} \leftarrow -0.04$$

$x$  $y$

$w$  $h$  $\lambda$  $l$

$b$

*Computation graph*

# Backward Pass

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

**Forward pass equations**

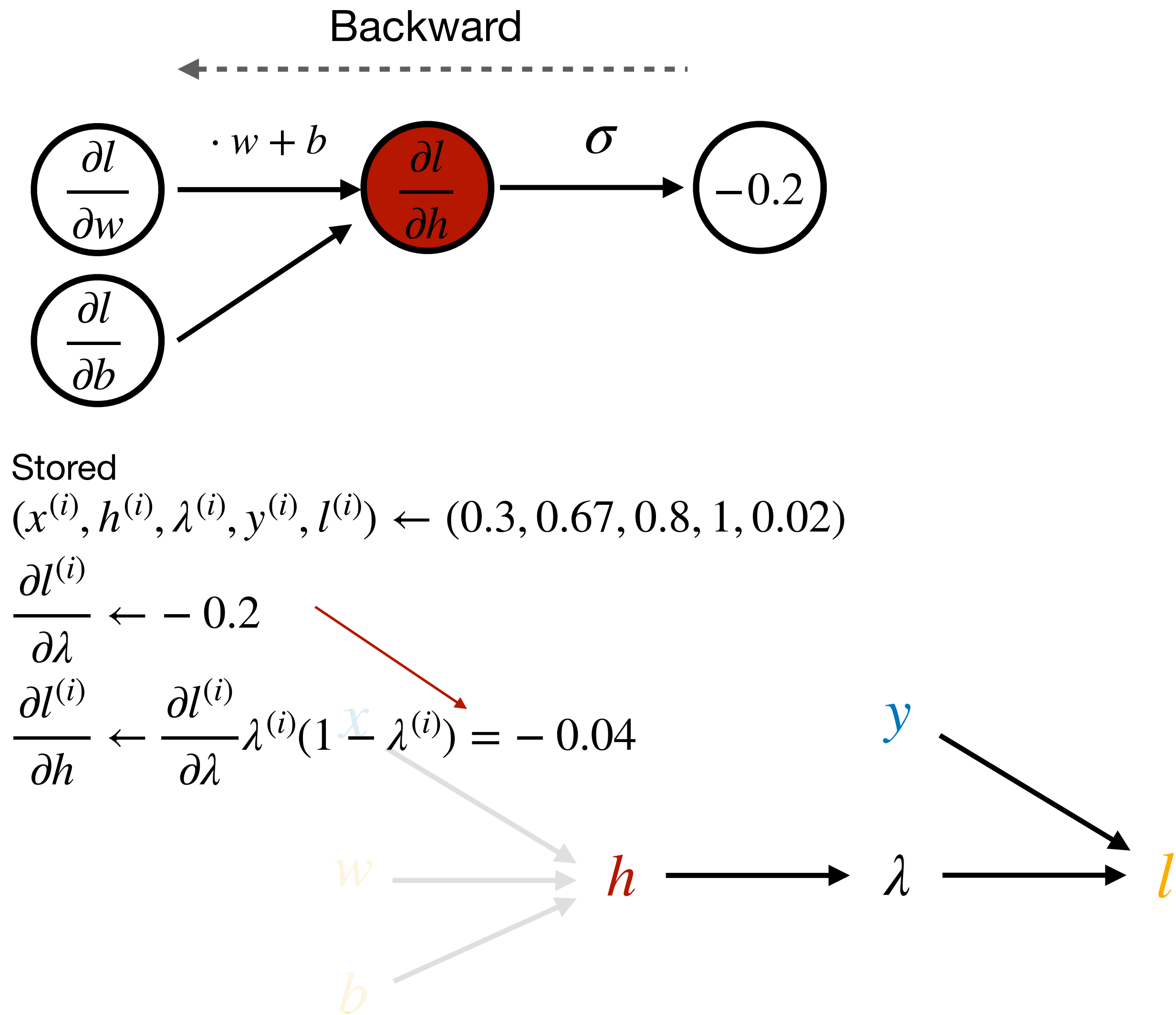$$l = \frac{1}{2}(y - \lambda)^2$$

$$\lambda = \sigma(h)$$

$$h = xw + b$$

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\frac{\partial h}{\partial w}$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial \lambda}\frac{\partial \lambda}{\partial h}\frac{\partial h}{\partial b}$$
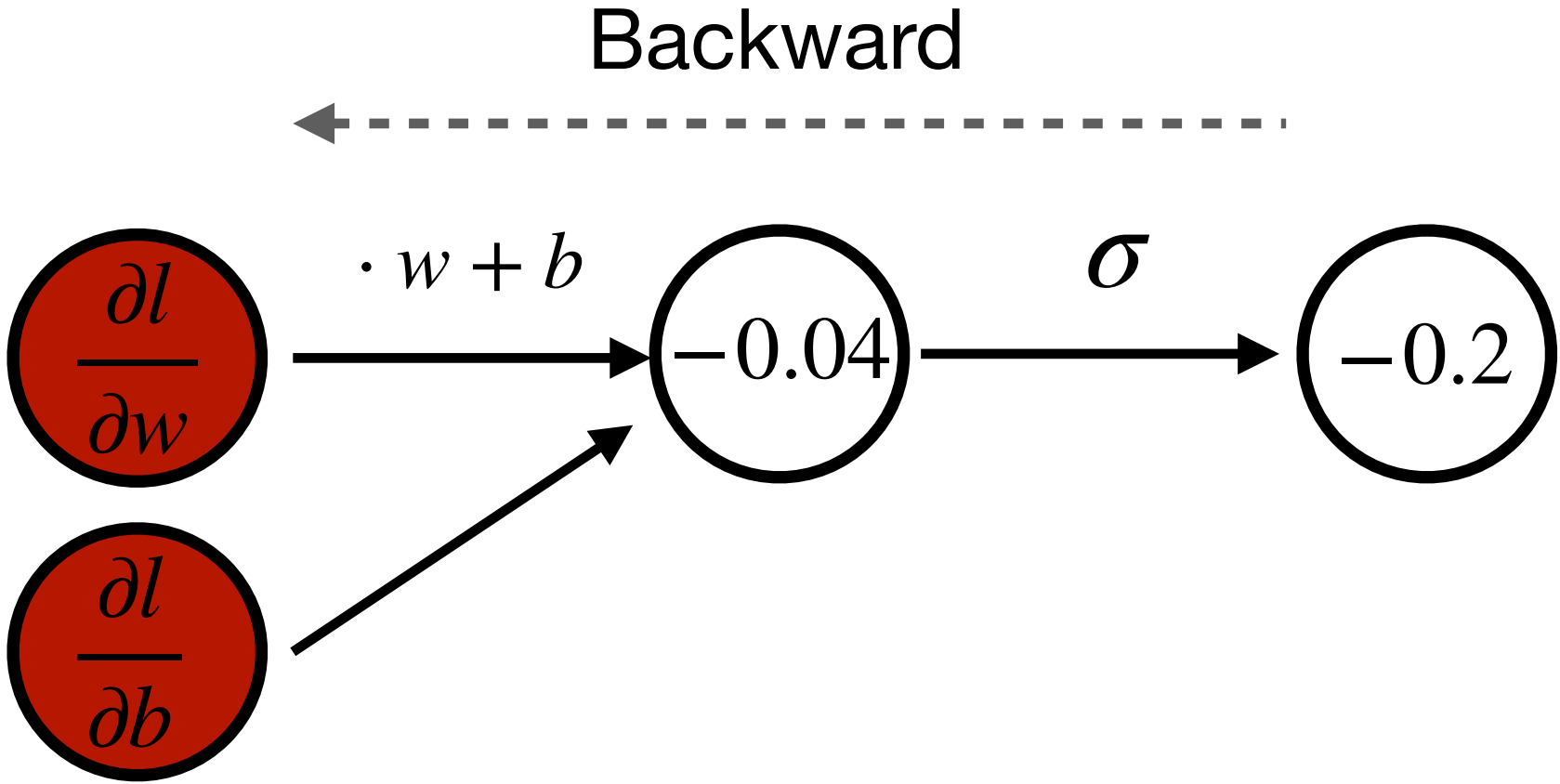
**Backward pass equations**

$$\frac{\partial l}{\partial \lambda} = (\lambda - y)$$

$$\frac{\partial l}{\partial h} = \frac{\partial l}{\partial \lambda}\lambda(1 - \lambda)$$

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial h}x$$
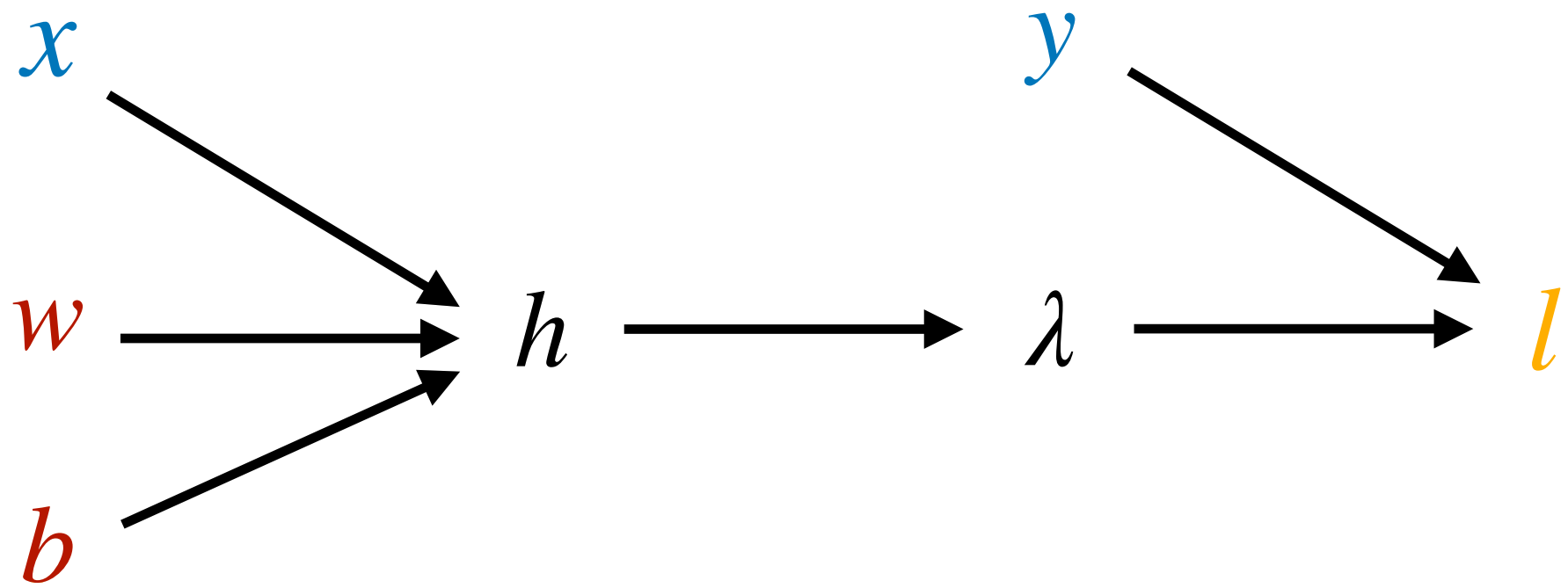
$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial h}$$

Backward

$$\frac{\partial l}{\partial w} \quad \cdot w + b \quad -0.04 \quad \sigma \quad -0.2$$

$$\frac{\partial l}{\partial b}$$

Stored

$$(x^{(i)}, h^{(i)}, \lambda^{(i)}, y^{(i)}, l^{(i)}) \leftarrow (0.3, 0.67, 0.8, 1, 0.02)$$

$$\frac{\partial l^{(i)}}{\partial \lambda} \leftarrow -0.2$$

$$\frac{\partial l^{(i)}}{\partial h} \leftarrow -0.04$$

$$\frac{\partial l^{(i)}}{\partial w} \leftarrow -0.26$$

$$\frac{\partial l^{(i)}}{\partial b} \leftarrow -0.04$$

$$x \quad y$$

$$w \quad h \quad \lambda \quad l$$

$$b$$

*Computation graph*

86

# Stochastic Gradient Descent

**Example: univariate regression**

$$l = \frac{1}{2}(y - \sigma(xw + b))^2$$

Step 1. Compute the derivatives $\quad \nabla L = \left( \dfrac{\partial L}{\partial w}, \dfrac{\partial L}{\partial b} \right)$

$$\frac{\partial l^{(i)}}{\partial w} = -0.26$$

$$\left\{ \frac{\partial l^{(i)}}{\partial w}, \frac{\partial l^{(i)}}{\partial b} \right\}_i \text{ for the mini-batch}$$

$$\frac{\partial L}{\partial w} \approx \frac{1}{m} \sum_{i=1}^{m} \frac{\partial l^{(i)}}{\partial w}$$

$$\frac{\partial l^{(i)}}{\partial b} = -0.04$$

$$\frac{\partial L}{\partial b} \approx \frac{1}{m} \sum_{i=1}^{m} \frac{\partial l^{(i)}}{\partial b}$$

Step 2. Update the parameters $\quad w \leftarrow w - \alpha \dfrac{\partial L}{\partial w}$

$$b \leftarrow b - \alpha \frac{\partial L}{\partial b}$$

# Backpropagation: exercise

**Exercise 1**

Consider the following forward pass equation for a two layer network

$$z_1 = x \cdot W_1 + b_1$$
$$a_1 = ReLU(z_1)$$
$$z_2 = a_1 \cdot W_2 + b_2$$
$$\lambda = \sigma(z_2)$$
$$l = -y \log(\lambda) - (1-y)\log(1-\lambda)$$

$$x^{(i)} \in \mathbb{R}^N, \ y^{(i)} \in \mathbb{R}^1, z_1 \in \mathbb{R}^{D_1}$$

1. What are the shapes of $W_1, b_1, W_2, b_2$? If we use the network across a mini-batch with $m$ samples, what would be the shape of $W_1, b_1, W_2, b_2$? and the shape of the input $x$ and output $y$?

2. Depict the computation graph

3. Write down the backward pass equations for $\dfrac{\partial l}{\partial W_2}, \dfrac{\partial l}{\partial b_2}$ [*hint: be careful with the dimensions!]**

# Backpropagation: solution

**Exercise 1**

Consider the following forward pass equation for a two layer network

$z_1 = x \cdot W_1 + b_1$

$a_1 = ReLU(z_1)$

$z_2 = a_1 \cdot W_2 + b_2$

$\lambda = \sigma(z_2)$

$l = -y \log(\lambda) - (1 - y)\log(1 - \lambda)$

$x^{(i)} \in \mathbb{R}^N, \, y^{(i)} \in \mathbb{R}^1, z_1 \in \mathbb{R}^{D_1}$

1. What are the shapes of $W_1, b_1, W_2, b_2$? If we use the network across a mini-batch with $m$ samples, what would be the shape of $W_1, b_1, W_2, b_2$? and the shape of the input $x$ and output $y$?

1. $W_1 \in \mathbb{R}^{N \times D_1}, \, b_1 \in \mathbb{R}^{D_1}, \, W_2 \in \mathbb{R}^{D_1}, \, b_2 \in \mathbb{R}^1$ . Weights and bias do not change shape with number of input samples.
$x \in \mathbb{R}^{m,N}, \, y \in \mathbb{R}^m$

# Backpropagation: solution

**Exercise 1**

Consider the following forward pass equation for a two layer network

$$z_1 = x \cdot W_1 + b_1$$

$$a_1 = ReLU(z_1)$$

$$z_2 = a_1 \cdot W_2 + b_2$$

$$\lambda = \sigma(z_2)$$

$$l = -y \log(\lambda) - (1-y)\log(1-\lambda)$$

$$x^{(i)} \in \mathbb{R}^N, \, y^{(i)} \in \mathbb{R}^1, \, z_1 \in \mathbb{R}^{D_1}$$

2. Depict the computation graph

# Backpropagation: solution

**Exercise 1**

Consider the following forward pass equation for a two layer network
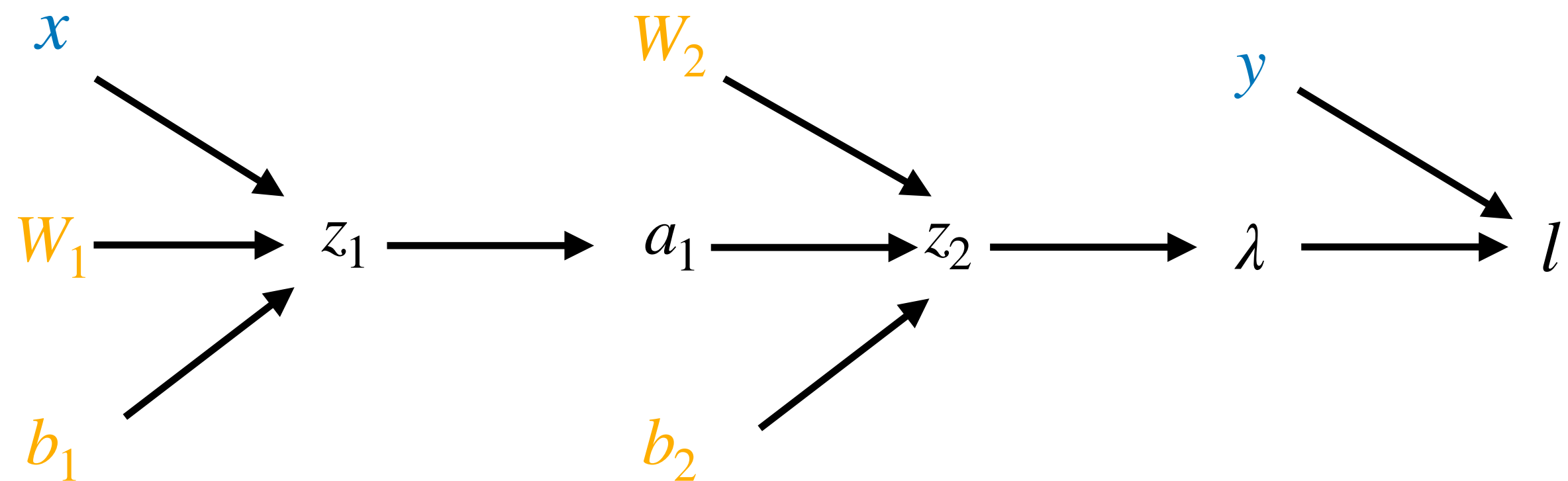
$$z_1 = x \cdot W_1 + b_1$$
$$a_1 = ReLU(z_1)$$
$$z_2 = a_1 \cdot W_2 + b_2$$
$$\lambda = \sigma(z_2)$$
$$l = -y \log(\lambda) - (1-y)\log(1-\lambda)$$

$$x^{(i)} \in \mathbb{R}^N, \, y^{(i)} \in \mathbb{R}^1, z_1 \in \mathbb{R}^{D_1}$$

$$\frac{\partial l}{\partial \lambda} = -\frac{y}{\lambda} + \frac{(1-y)}{1-\lambda} \in \mathbb{R}$$

$$\frac{\partial l}{\partial z_2} = \frac{\partial l}{\partial \lambda} \, \sigma(z_2)(1-\sigma(z_2)) \in \mathbb{R}$$

3. Write down the backward pass equations for $\dfrac{\partial l}{\partial W_2}, \dfrac{\partial l}{\partial b_2}$

$$\frac{\partial l}{\partial W_2} = \frac{\partial l}{\partial z_2} a_1 \in \mathbb{R}^{D_1} \qquad \frac{\partial l}{\partial b_2} = \frac{\partial l}{\partial z_2} \cdot 1 \in \mathbb{R}$$

# Maximum Likelihood for Regression (optional)

Assume you have a training set $\{x^{(i)}, y^{(i)}\}_{i=1,\dots,n}$ for $x, y \in \mathbb{R}$ and you want to train a neural network to predict the outcome $y$ given the value of $x$.

We follow the steps to build the loss function according to the maximum likelihood criterion:

1. We assume the data are generated by a Gaussian distribution $p_{data}(y \,|\, x) \sim \mathcal{N}(\mu_x, \sigma = 1)$

2. We want to learn $p_{model}(y \,|\, \lambda_x) \approx p_{data}(y \,|\, x)$

3. We set the network to output an approximation of the mean of $f(x; \theta) = \lambda_x \approx \mu_x$

What is the loss function build according to maximum likelihood criterion for this problem?

*where $p_{model}$ as $p_{data}$ are intended as probability density functions (it's a continuous space problem)

# Maximum Likelihood for Regression (optional)

Assume you have a training set $\{x^{(i)}, y^{(i)}\}_{i=1,\ldots,n}$ for $x, y \in \mathbb{R}$ and you want to train a neural network to predict the outcome $y$ given the value of $x$.

We follow the steps to build the loss function according to the maximum likelihood criterion:

1. We assume the data are generated by a Gaussian distribution $p_{data}(y \,|\, x) \sim \mathcal{N}(\mu_x, \sigma = 1)$

2. We want to learn $p_{model}(y \,|\, \lambda_x) \approx p_{data}(y \,|\, x)$

3. We set the network to output an approximation of the mean of $f(x; \theta) = \lambda_x \approx \mu_x$

What is the loss function build according to maximum likelihood criterion for this problem?

$$L(\theta) = \frac{1}{n}\sum_{i=1}^{n} -\log(p_{model}(y^{(i)} \,|\, f(x^{(i)}; \theta))) = \frac{1}{n}\sum_{i=1}^{n} -\log\left(\frac{1}{\sqrt{2\pi}} e^{-\frac{\left(y^{(i)} - \lambda^{(i)}\right)^2}{2}}\right) = \frac{1}{n}\sum_{i=1}^{n} \log(\sqrt{2\pi}) + \frac{\left(y^{(i)} - \lambda^{(i)}\right)^2}{2}$$

$$= \text{const} + \frac{1}{n}\sum_{i=1}^{n} \frac{\left(y^{(i)} - \lambda^{(i)}\right)^2}{2} = L_{SE}$$

This is the maximum likelihood formulation for a regression problem and indeed the loss is the mean squared error :)

*where $p_{model}$ as $p_{data}$ are intended as probability density functions (it's a continuous space problem)

# Summary

**Topics**

- Loss function

- Maximum-likelihood

- Stochastic Gradient Discent

- Backpropagation algorithm

**Reading material**

- *Understanding Deep Learning* - Chapter 5 (5.1, 5.2, 5.7), Chapter 6 (6.1)

- *Deep Learning book* - Chapter 5.5

- [Backpropagation1](#) [Backpropagation2](#)