

# Computer Vision with Deep Learning

## Computation Graphs

# PLAN

Logistic Regression

Computation Graphs

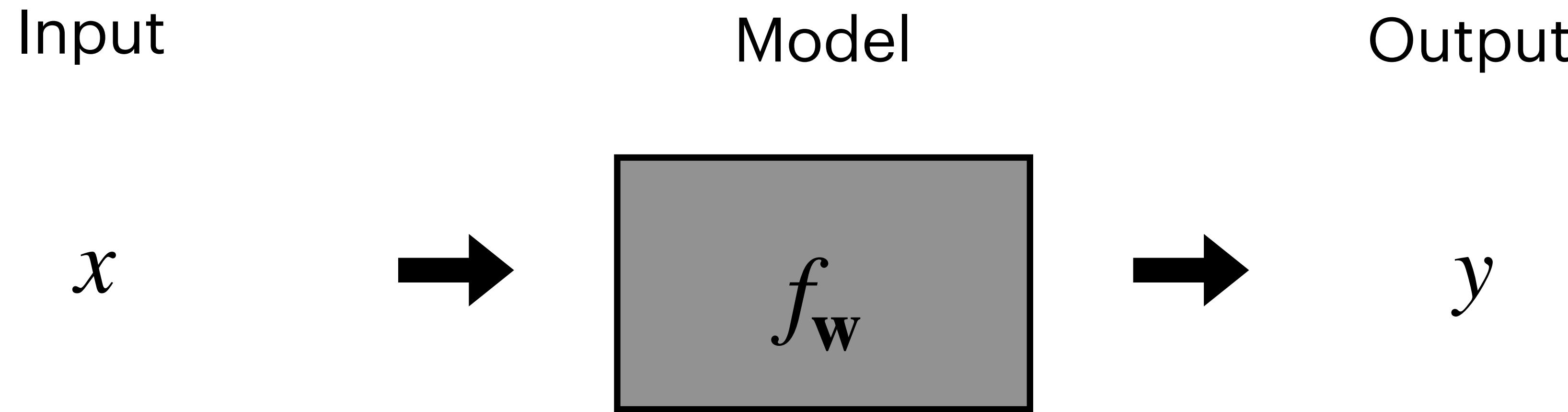
Backpropagation

Hands on



# Quick Recap

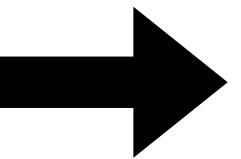
# Supervised Learning



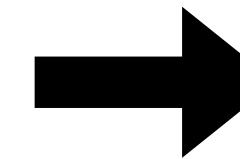
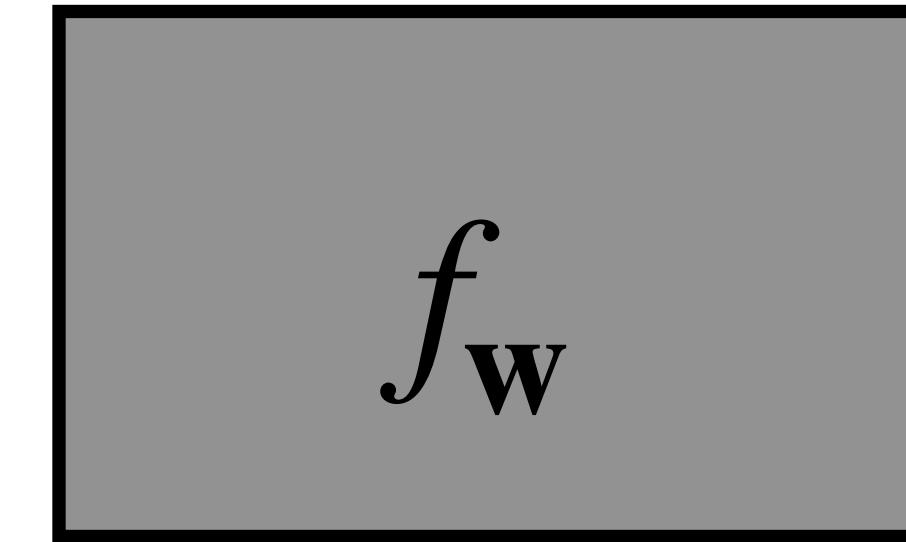
- ◆ **Learning:** Estimate parameters  $\mathbf{w}$  from training data  $\{(x_i, y_i)\}_{i=1}^N$
- ◆ **Inference:** Make novel predictions  $y = f_{\mathbf{w}}(x)$

# Regression

Input



Model



Output

143,52 €

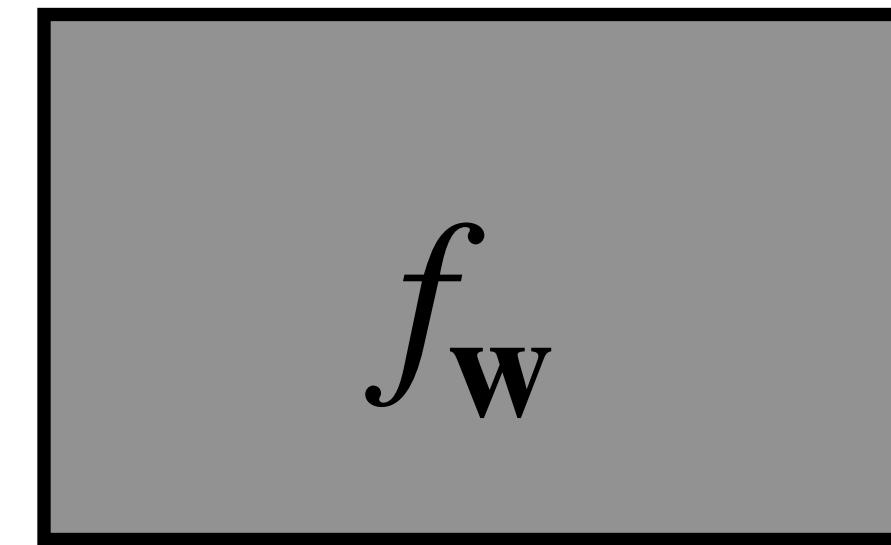
♦ **Mapping:**  $f_w : \mathcal{R}^N \mapsto \mathcal{R}$

# Classification

Input



Model



Output

“Beach”

◆ **Mapping:**  $f_w : \mathcal{R}^{W \times H} \mapsto \{"Beach", "No beach"\}$

# Logistic Regression

# Logistic Regression

Conditional **Maximum Likelihood Estimator** for  $\mathbf{w}$ :

$$\hat{\mathbf{w}}_{ML} = \arg \max_{\mathbf{w}} \sum_{i=1}^N \log p_{model}(y_i | \mathbf{x}_i, \mathbf{w})$$

- ◆ Let's perform binary classification:  $y_i \in \{0,1\}$
- ◆ How to choose  $p_{model}(y | \mathbf{x}, \mathbf{w})$  for this problem?
- ◆ Answer: Bernoulli distribution

$$p_{model}(y | \mathbf{x}, \mathbf{w}) = \hat{y}^y (1 - \hat{y})^{(1-y)}$$

where  $\hat{y}$  is predicted by a model:  $\hat{y} = f_{\mathbf{w}}(\mathbf{x})$

# Logistic Regression

We assumed a Bernoulli distribution

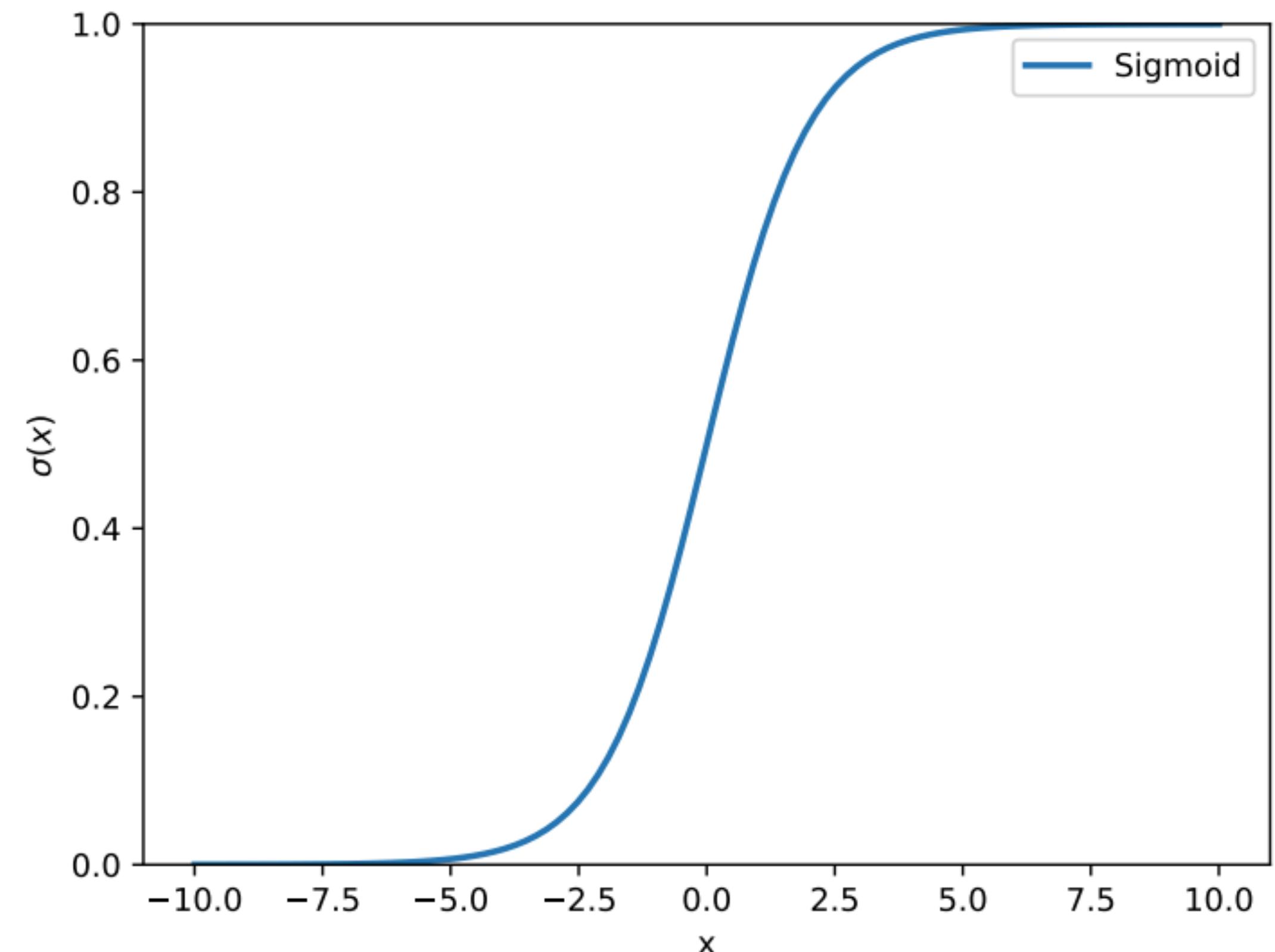
$$p_{model}(y | \mathbf{x}, \mathbf{w}) = \hat{y}^y (1 - \hat{y})^{(1-y)}$$

where  $\hat{y}$  is predicted by a model:  $\hat{y} = f_{\mathbf{w}}(\mathbf{x})$ .

- ◆ How to choose  $f_{\mathbf{w}}(\mathbf{x})$ ?
- ◆ Requirement:  $f_{\mathbf{w}}(\mathbf{x}) \in [0,1]$
- ◆ Choose  $f_{\mathbf{w}}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$

where  $\sigma$  is the sigmoid function:

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



# Logistic Regression

Putting it together:

$$\begin{aligned}\hat{\mathbf{w}}_{ML} &= \arg \max_{\mathbf{w}} \sum_{i=1}^N \log p_{model}(y_i | \mathbf{x}_i, \mathbf{w}) \\ &= \arg \max_{\mathbf{w}} \sum_{i=1}^N \log [\hat{y}_i^{y_i} (1 - \hat{y}_i)^{(1-y_i)}] \\ &= \arg \min_{\mathbf{w}} \sum_{i=1}^N -y_i \log \hat{y}_i - (1 - y_i) \log (1 - \hat{y}_i)\end{aligned}$$

Binary Cross  
Entropy Loss

$\mathcal{L}(\hat{y}_i, y_i)$

- ♦ In ML, we use the more general term “loss function” rather than “error function”
- ♦ Interpretation: We minimize the dissimilarity between the empirical data distribution  $p_{data}$  (defined by the training set) and the model distribution  $p_{model}$

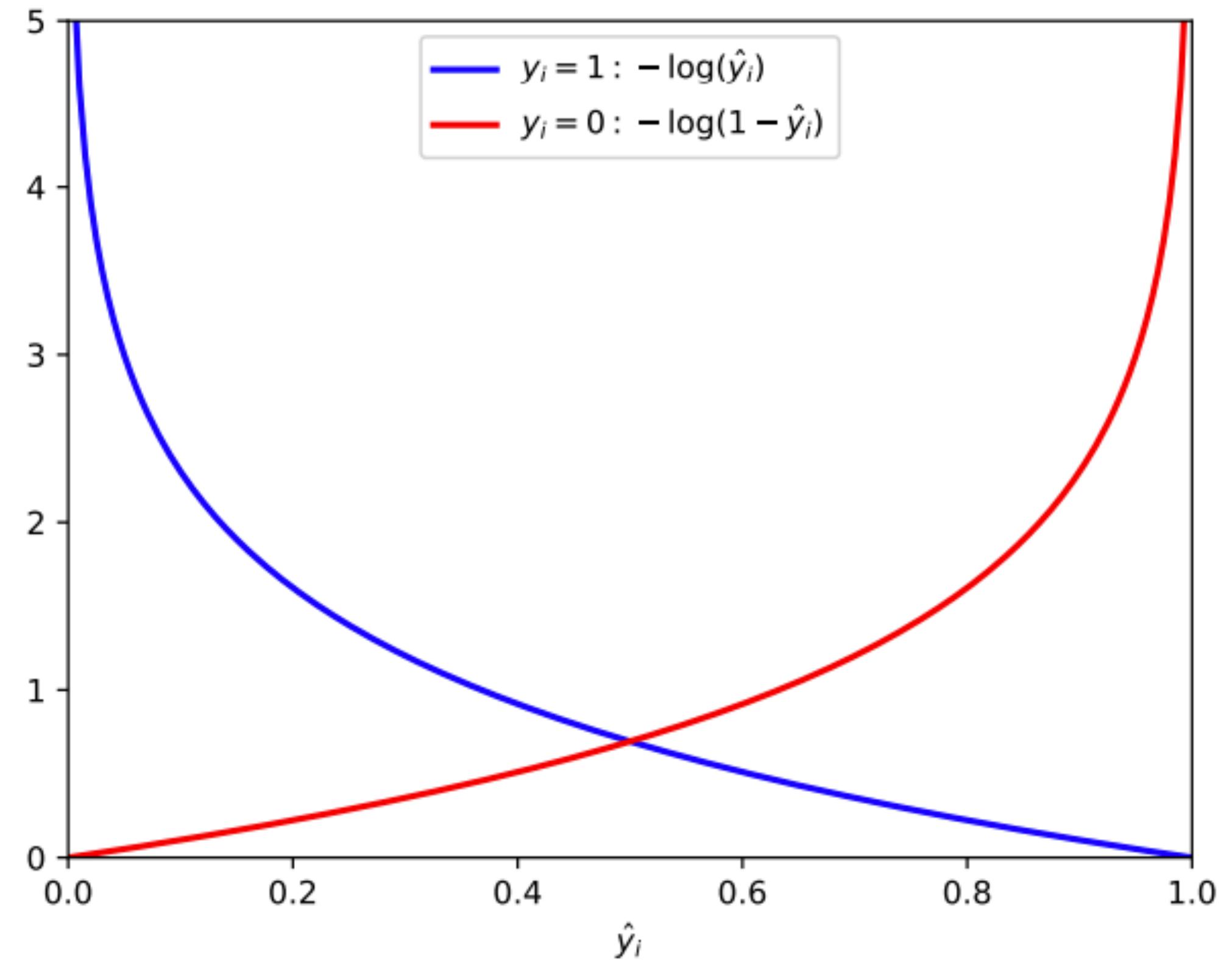
# Logistic Regression

## Binary Cross Entropy Loss:

$$\mathcal{L}(\hat{y}_i, y_i) = -y_i \log \hat{y}_i - (1 - y_i) \log (1 - \hat{y}_i)$$

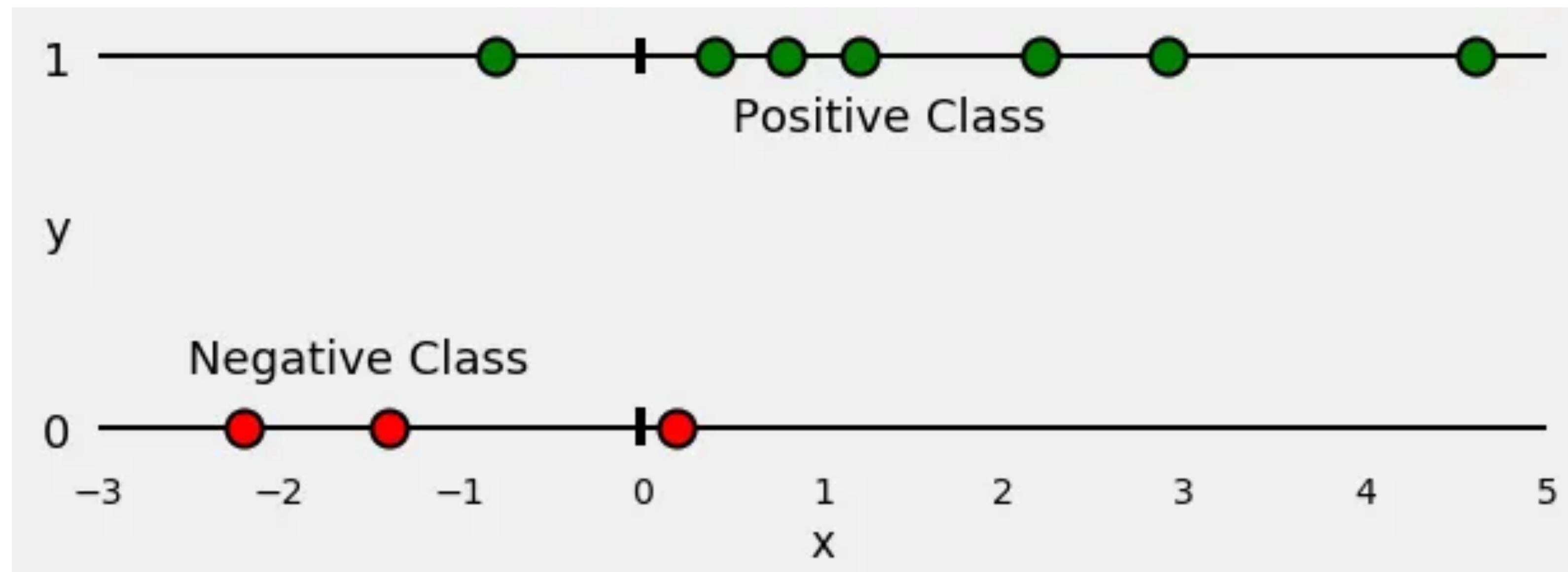
$$= \begin{cases} -\log \hat{y}_i & \text{if } y_i = 1 \\ -\log (1 - \hat{y}_i) & \text{if } y_i = 0 \end{cases}$$

- ◆ For  $y_i = 1$ , the loss  $\mathcal{L}$  is minimized if  $\hat{y}_i = 1$
- ◆ For  $y_i = 0$ , the loss  $\mathcal{L}$  is minimized if  $\hat{y}_i = 0$
- ◆ Thus  $\mathcal{L}$  is minimal if  $\hat{y}_i = y_i$
- ◆ Can be extended to  $> 2$  classes



# Logistic Regression

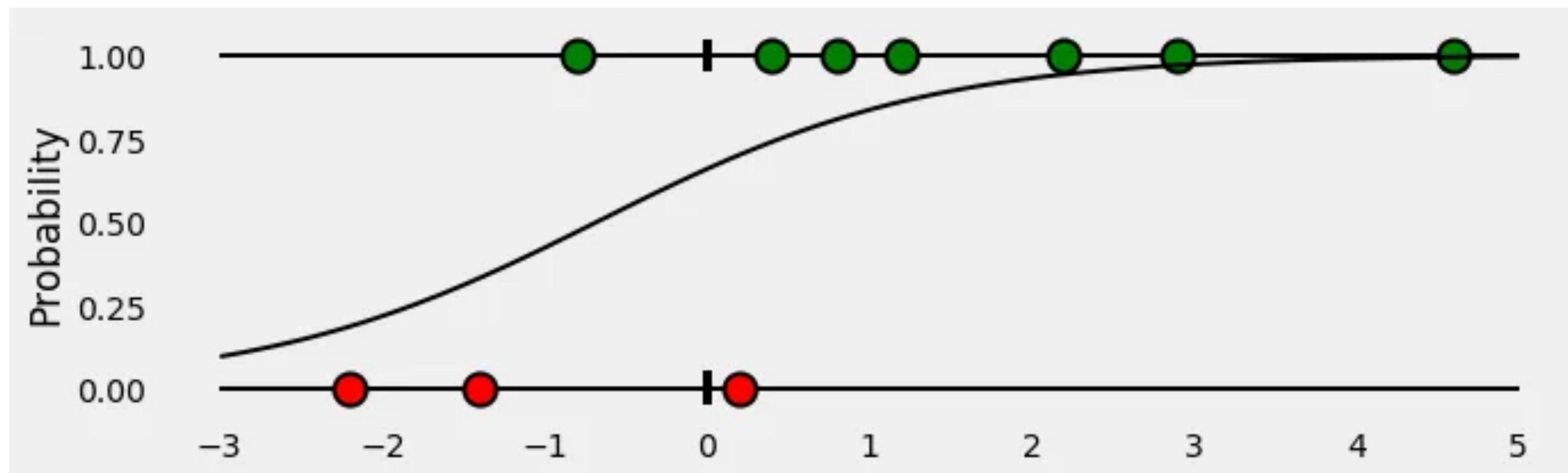
**A simple 1D example:**



- ◆ Dataset  $\mathcal{X}$  with positive ( $y_i = 1$ ) and negative ( $y_i = 0$ ) samples

# Logistic Regression

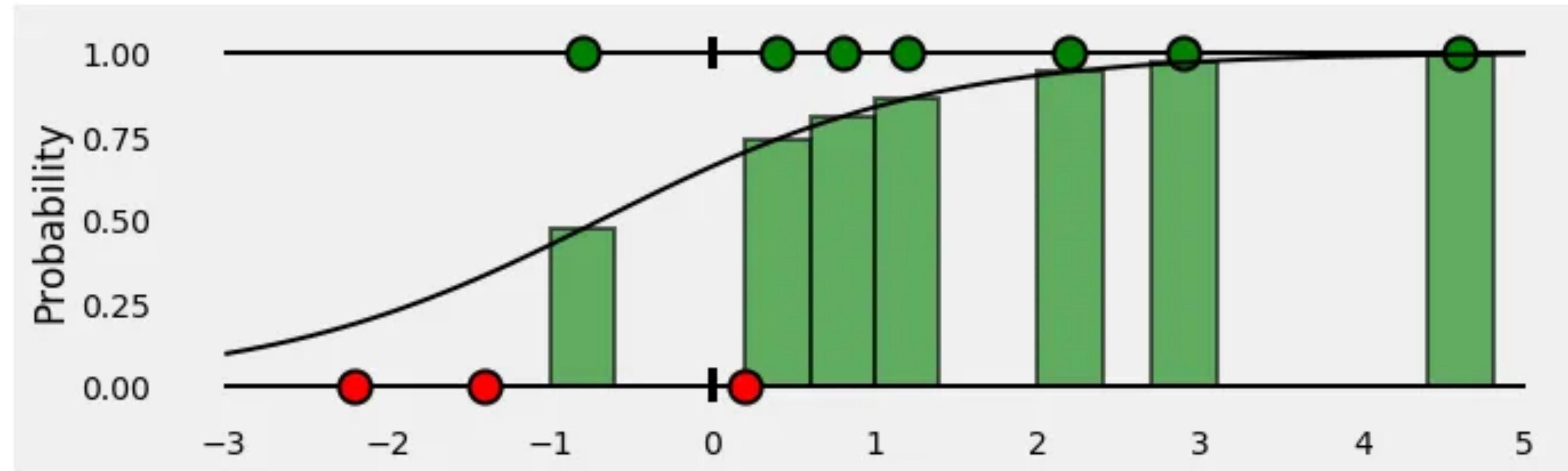
**A simple 1D example:**



- ◆ Logistic regressor  $f_{\mathbf{w}}(x) = \sigma(w_0 + w_1 x)$  fit to dataset  $\mathcal{X}$

# Logistic Regression

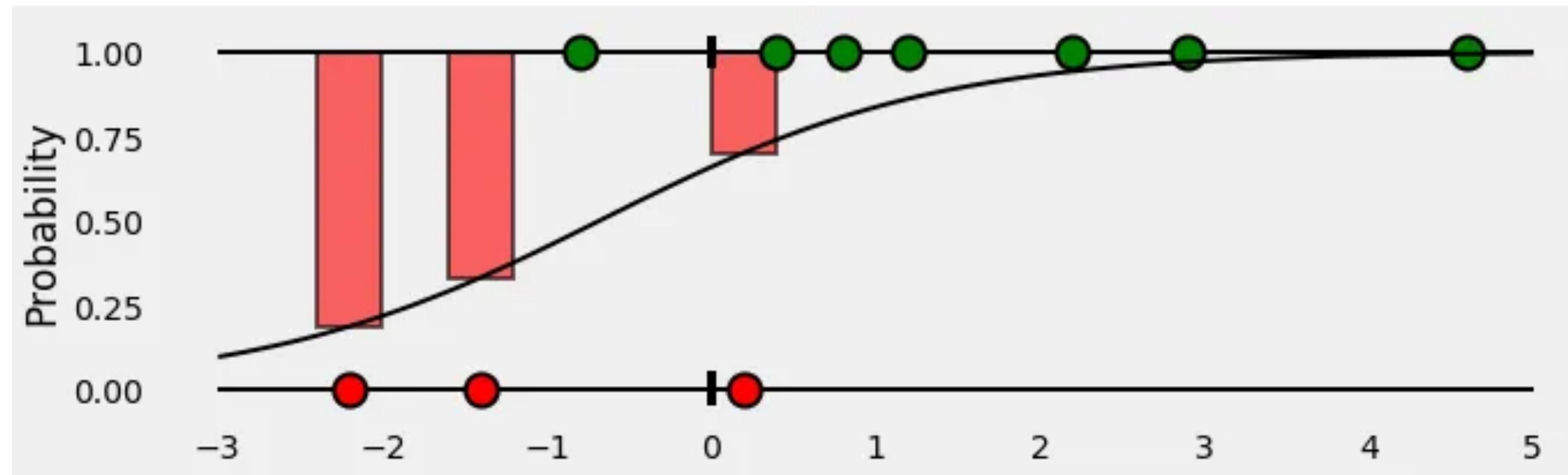
**A simple 1D example:**



- ◆ Probabilities of classifier  $f_{\mathbf{w}}(x_i)$  for positive samples ( $y_i = 1$ )

# Logistic Regression

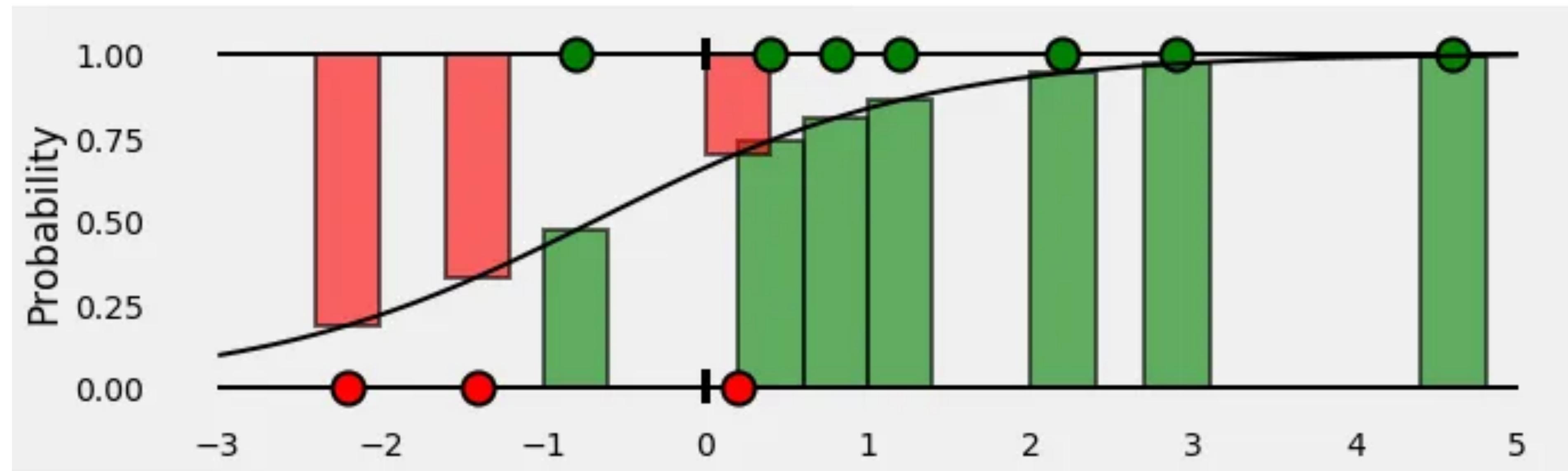
**A simple 1D example:**



- ◆ Probabilities of classifier  $f_w(x_i)$  for negative samples ( $y_i = 0$ )

# Logistic Regression

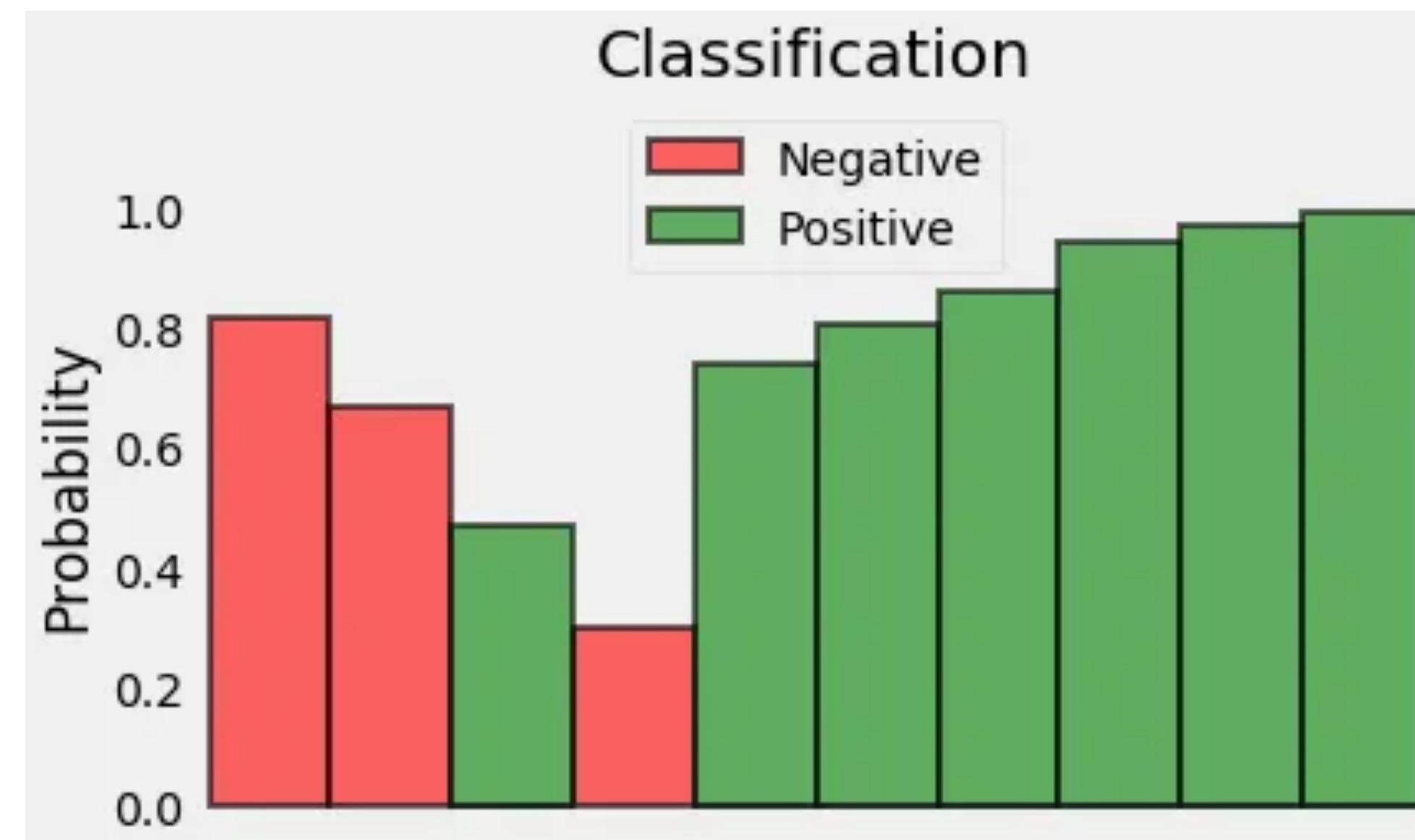
**A simple 1D example:**



◆ Putting it together

# Logistic Regression

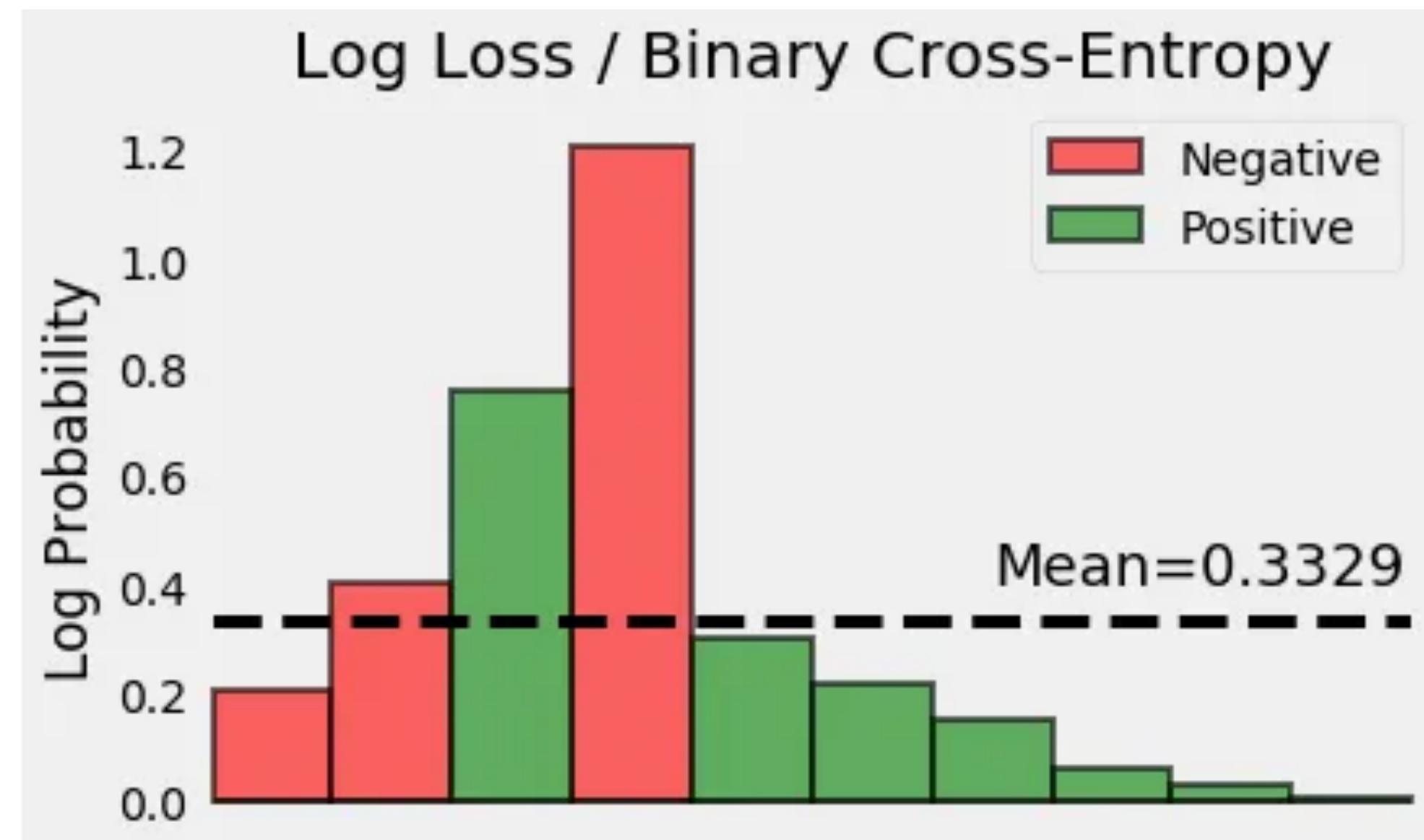
**A simple 1D example:**



- ◆ Let's get rid of the x axis

# Logistic Regression

**A simple 1D example:**



- ◆ And finally compute the negative logarithm:  $-\log(f_w(x_i))$

# Logistic Regression

**Maximum Likelihood for Logistic Regression:**

$$\hat{\mathbf{w}}_{ML} = \arg \min_{\mathbf{w}} \sum_{i=1}^N -y_i \log \hat{y}_i - (1 - y_i) \log (1 - \hat{y}_i)$$

where  $\hat{y} = f_{\mathbf{w}}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$  and  $\sigma(x) = \frac{1}{1 + e^{-x}}$

How do we find the minimizer  $\hat{\mathbf{w}}$  ?

- ◆ In contrast to linear regression, the loss  $\mathcal{L}(\hat{y}_i, y_i)$  is not quadratic in  $\mathbf{w}$
- ◆ We must apply iterative gradient-based optimization. The gradient is given by:

$$\nabla_{\mathbf{w}} \mathcal{L}(\hat{y}_i, y_i) = (\hat{y}_i - y_i) \mathbf{x}_i$$

Binary Cross  
Entropy Loss  
 $\mathcal{L}(\hat{y}_i, y_i)$

# Logistic Regression

## Gradient Descent:

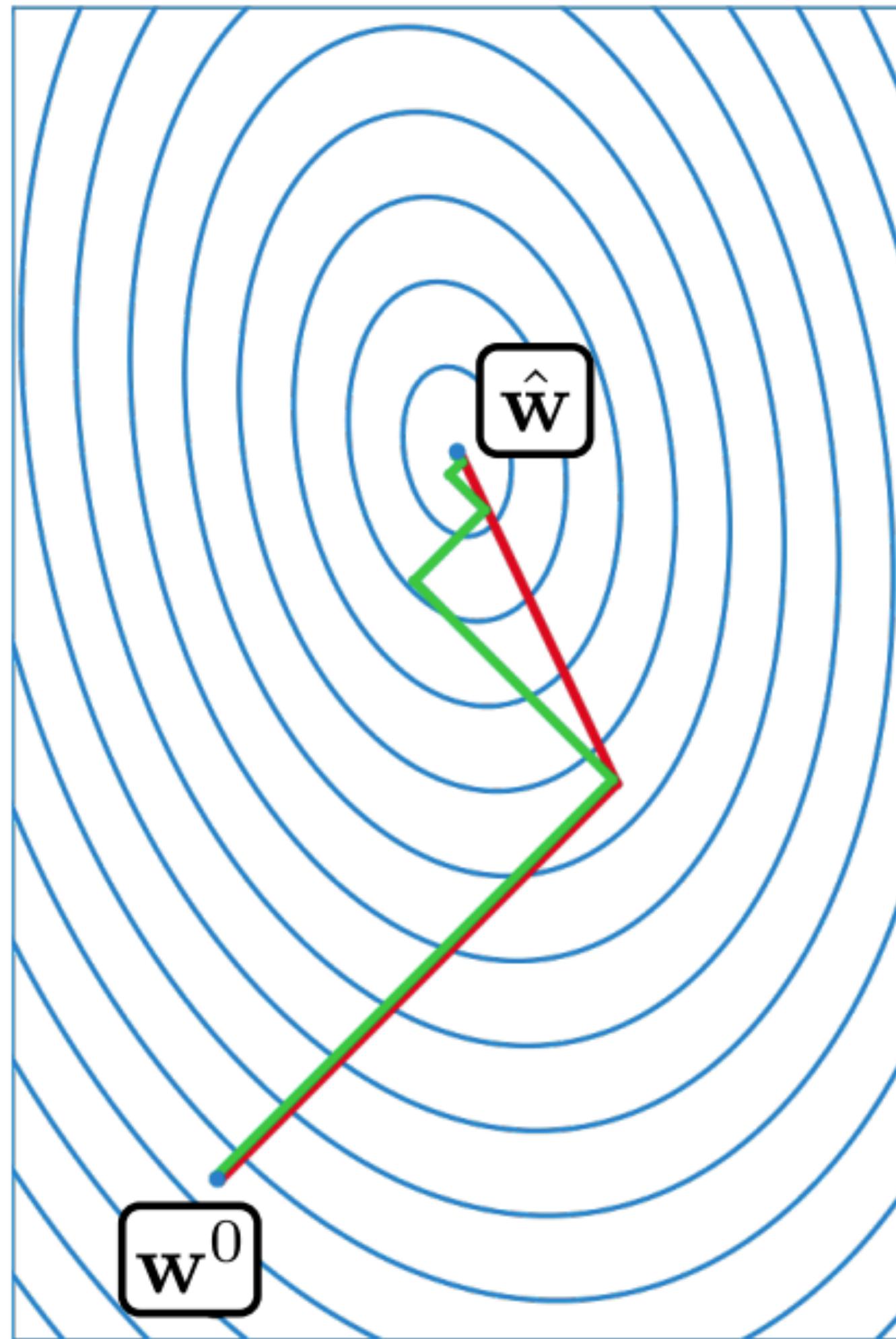
- ◆ Pick step size  $\eta$  and tolerance  $\epsilon$
- ◆ Initialize  $\mathbf{w}^0$
- ◆ Repeat until  $\|\mathbf{v}\| < \epsilon$

$$\blacktriangleright \mathbf{v} = \nabla_{\mathbf{w}} \mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}) = \sum_{i=1}^N \nabla_{\mathbf{w}} \mathcal{L}(\hat{y}_i, y_i)$$

$$\blacktriangleright \mathbf{w}^{t+1} = \mathbf{w}^t - \eta \mathbf{v}$$

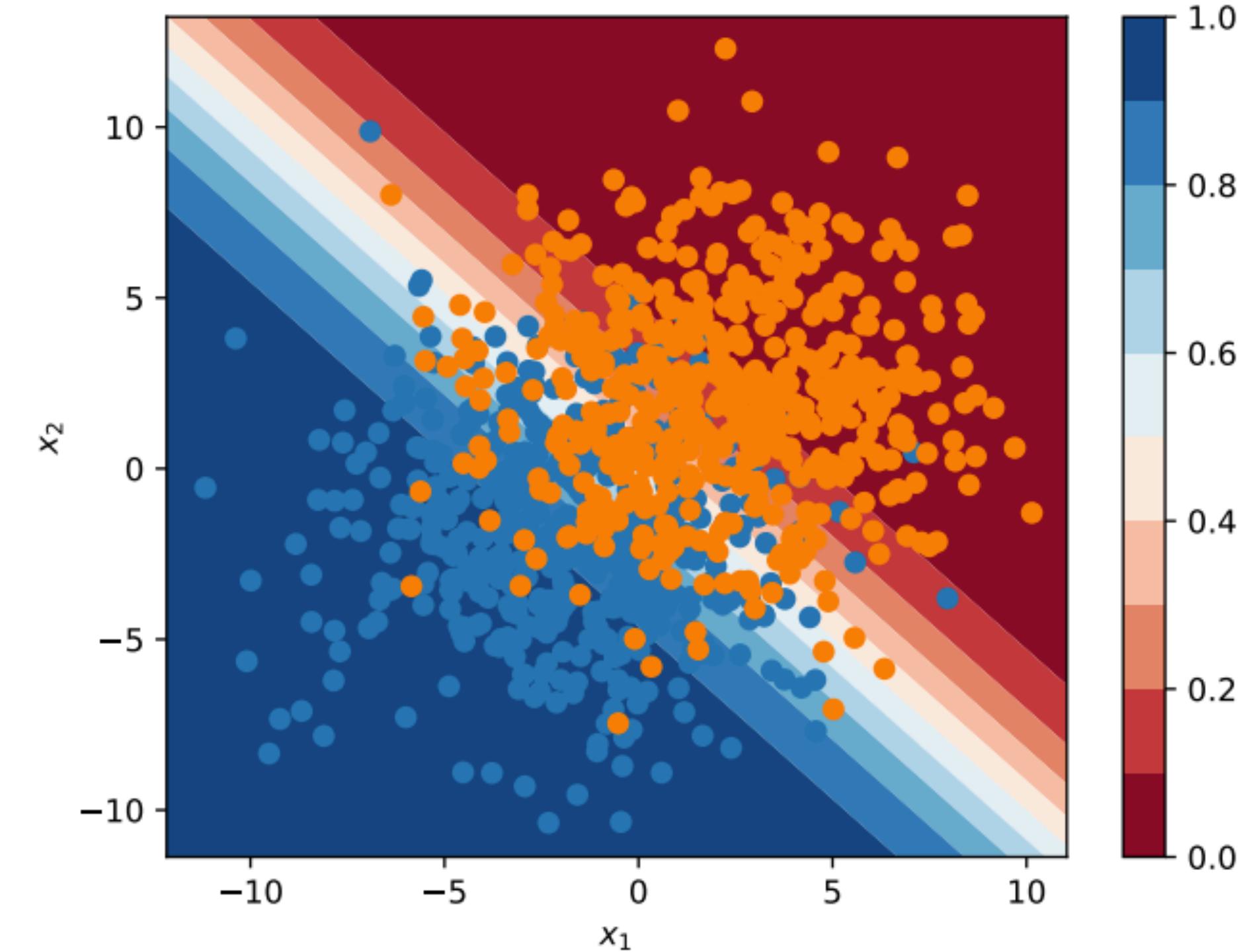
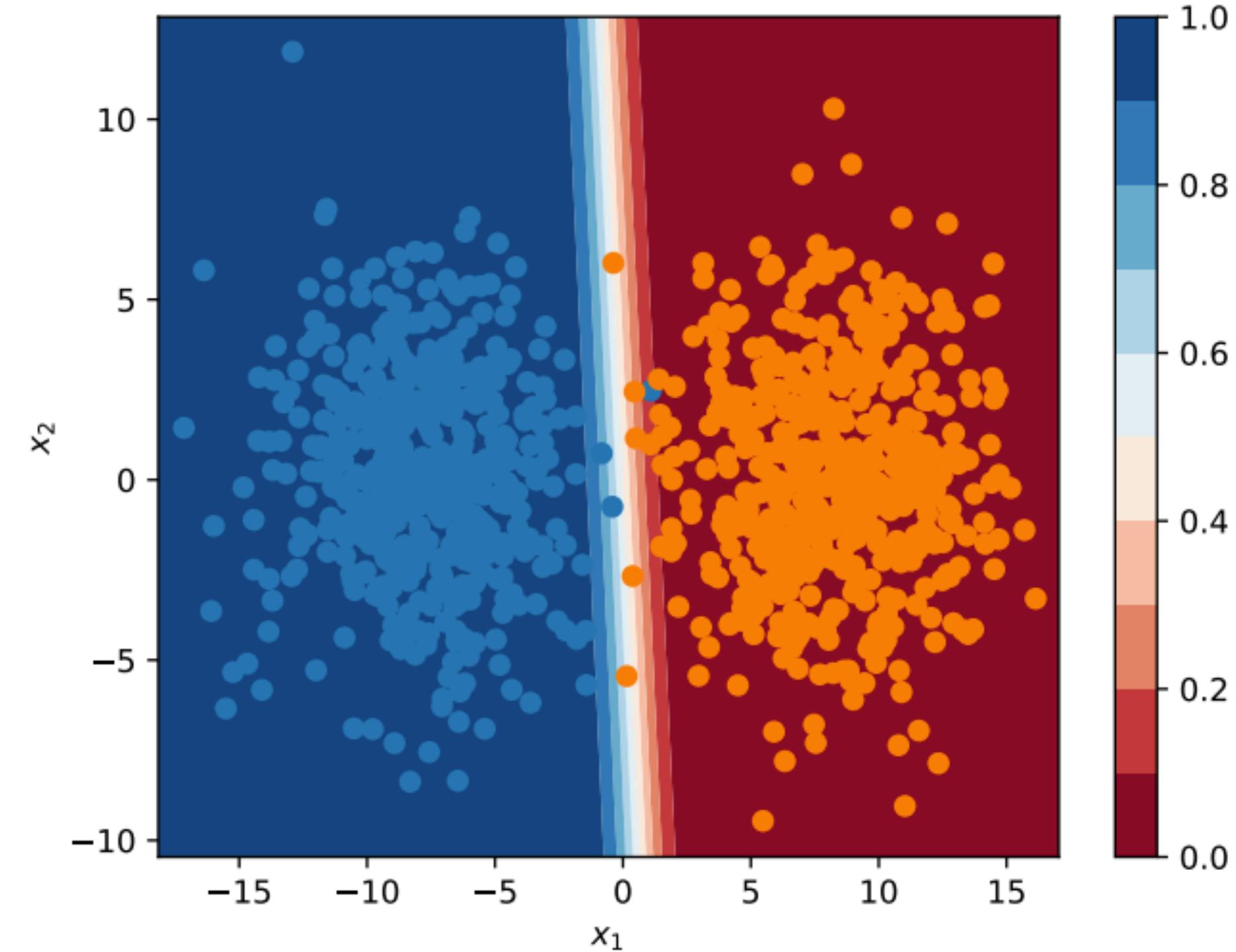
## Variants:

- ◆ Line search (green)
- ◆ Conjugate gradients (red)
- ◆ L-BFGS



# Logistic Regression

**Examples with two-dimensional inputs  $(x_1, x_2) \in \mathcal{R}^2$ :**



- ◆ Logistic regression model:  $f_w(x_1, x_2) = \sigma(w_0 + w_1 x_1 + w_2 x_2)$

# Logistic Regression

Maximizing the **Log-Likelihood** is equivalent to minimizing **Cross Entropy** or **KL Divergence**:

$$\hat{\mathbf{w}}_{ML} = \arg \max_{\mathbf{w}} \sum_{i=1}^N \log p_{model}(y_i | \mathbf{x}_i, \mathbf{w})$$

Log-likelihood

$$= \arg \max_{\mathbf{w}} \mathcal{E}_{p_{data}} [\log p_{model}(y | \mathbf{x}, \mathbf{w})]$$

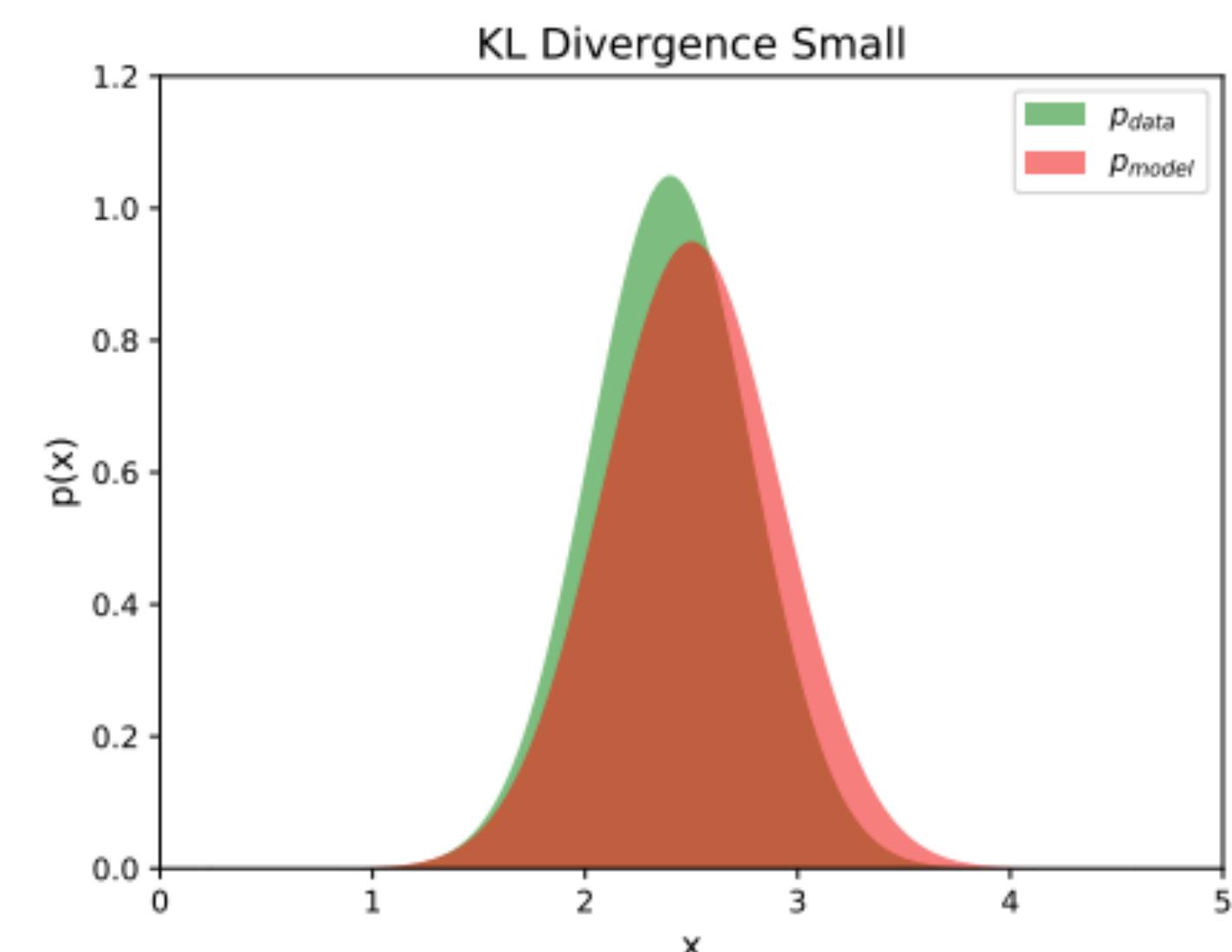
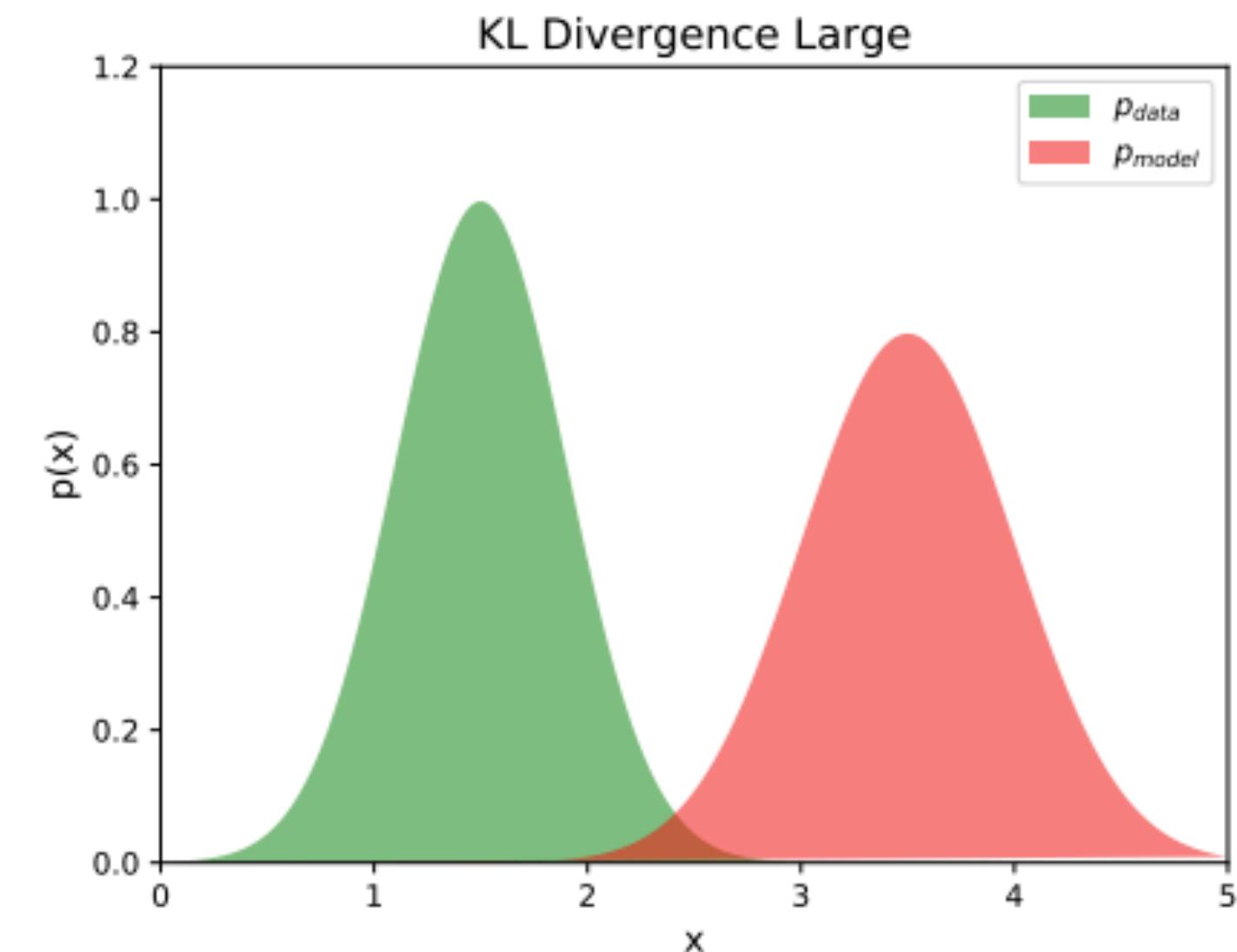
$$= \arg \min_{\mathbf{w}} -\mathcal{E}_{p_{data}} [\log p_{model}(y | \mathbf{x}, \mathbf{w})]$$

Cross Entropy  
 $H(p_{data}, p_{model})$

$$= \arg \min_{\mathbf{w}} \mathcal{E}_{p_{data}} [\log p_{data}(y | \mathbf{x}) - p_{model}(y | \mathbf{x}, \mathbf{w})]$$

$$= \arg \min_{\mathbf{w}} D_{KL}(p_{data} || p_{model})$$

KL Divergence



# Computation Graphs

# Logistic Regression

## Maximum Likelihood for Logistic Regression:

$$\hat{\mathbf{w}}_{ML} = \arg \min_{\mathbf{w}} \sum_{i=1}^N -y_i \log \hat{y}_i - (1 - y_i) \log (1 - \hat{y}_i)$$

Binary Cross  
Entropy Loss  
 $\mathcal{L}(\hat{y}_i, y_i)$

$$\text{where } \hat{y} = f_{\mathbf{w}}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}) \text{ and } \sigma(x) = \frac{1}{1 + e^{-x}}$$

- ◆ Minimization of a **non-linear objective** requires the calculation of gradients  $\nabla_{\mathbf{w}}$
- ◆ Luckily, in the above case the gradient is simple:  $\nabla_{\mathbf{w}} \mathcal{L}(\hat{y}_i, y_i) = (\hat{y}_i - y_i) \mathbf{x}_i$
- ◆ But this is not true for more complex models such as deep neural networks
- ◆ How can we **efficiently** compute gradients in the general case?

# Computation Graphs

## Key Idea:

- ◆ **Decompose** complex computations into sequence of atomic assignments
- ◆ We call this sequence of assignments a **computation graph** or **source code**
- ◆ The **forward pass** takes a training point  $(\mathbf{x}, y)$  as input and computes a loss, e.g.:

$$\mathcal{L} = - \log p_{model}(y \mid \mathbf{x}, \mathbf{w})$$

- ◆ Gradients  $\nabla_{\mathbf{w}} \mathcal{L}$  can be computed using a **backward pass**
- ◆ Both forward and backward pass are **efficient** due to the use of dynamic programming, i.e., storing and reusing intermediate results
- ◆ This decomposition and reuse of computation is key to the success of the **backpropagation algorithm**, the primary workhorse of deep learning

# Computation Graphs

A **computation graph** has three kinds of nodes:

- Input nodes
- Parameter nodes
- Compute nodes

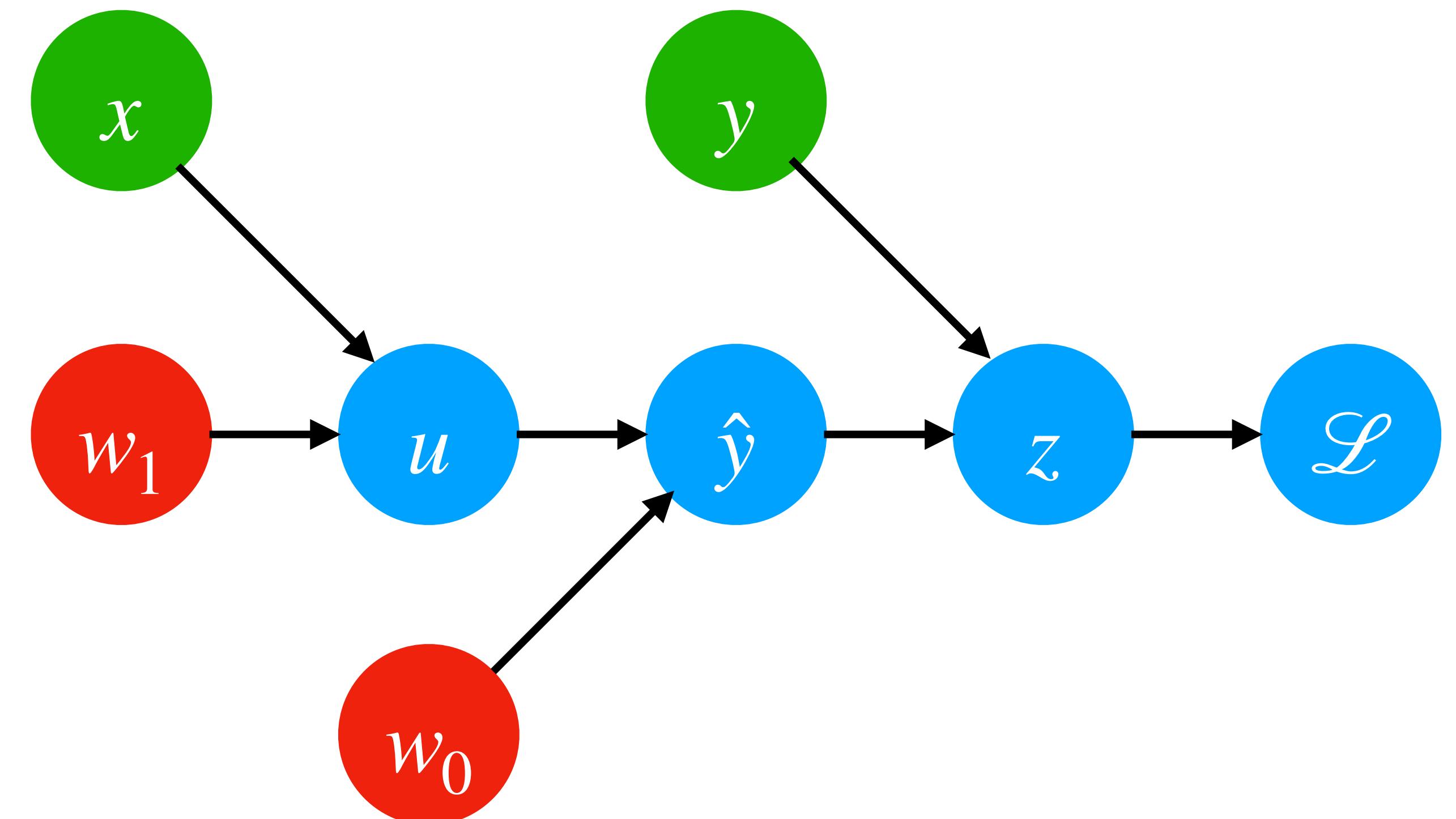
**Example:** Linear Regression

$$(1) \quad u = w_1 x$$

$$(2) \quad \hat{y} = w_0 + u$$

$$(3) \quad z = \hat{y} - y$$

$$(4) \quad \mathcal{L} = z^2$$



# Computation Graphs

A **computation graph** has three kinds of nodes:

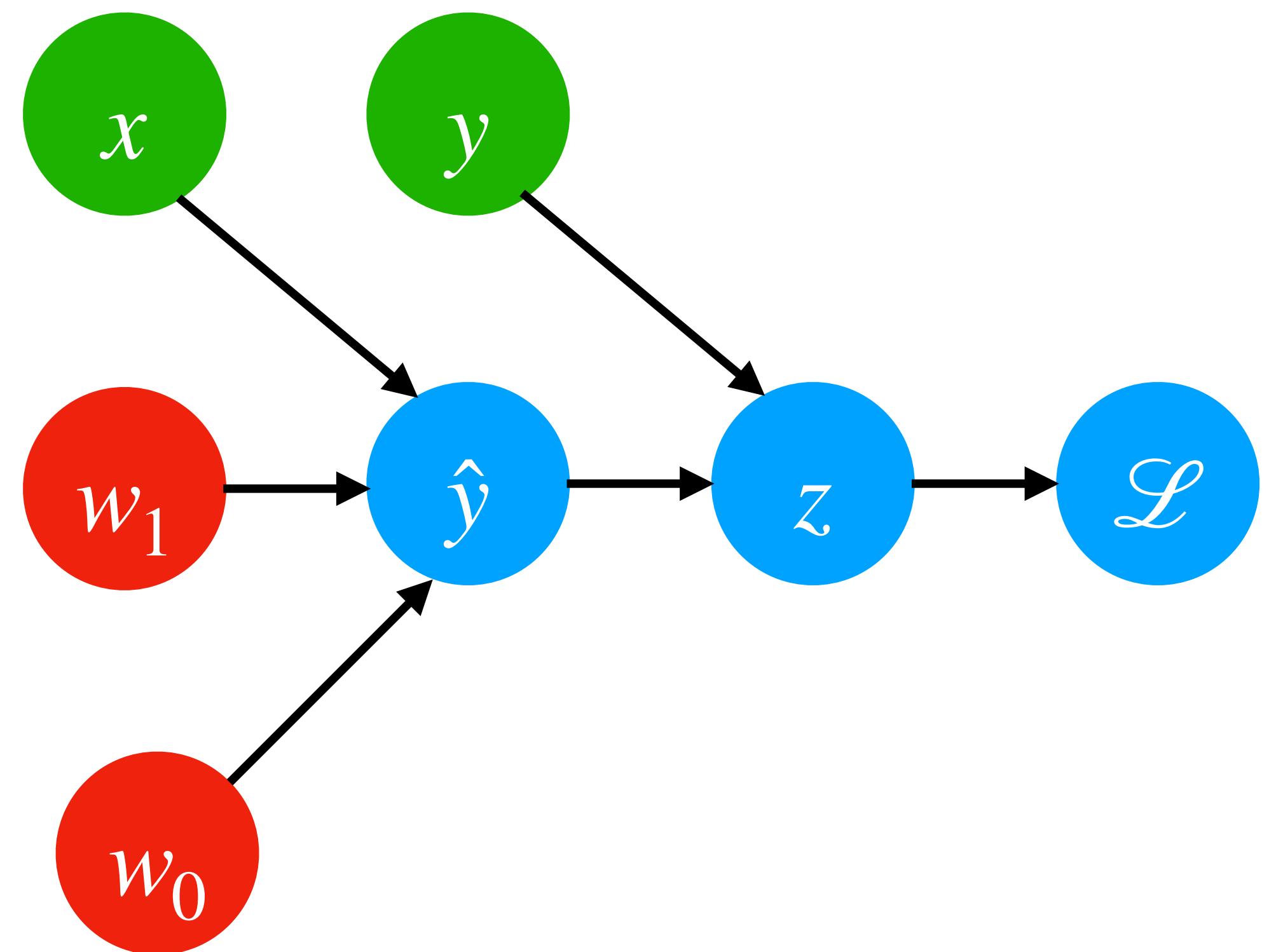
- Input nodes
- Parameter nodes
- Compute nodes

**Example:** Linear Regression

$$(1) \quad \hat{y} = w_0 + w_1 x$$

$$(2) \quad z = \hat{y} - y$$

$$(3) \quad \mathcal{L} = z^2$$



# Computation Graphs

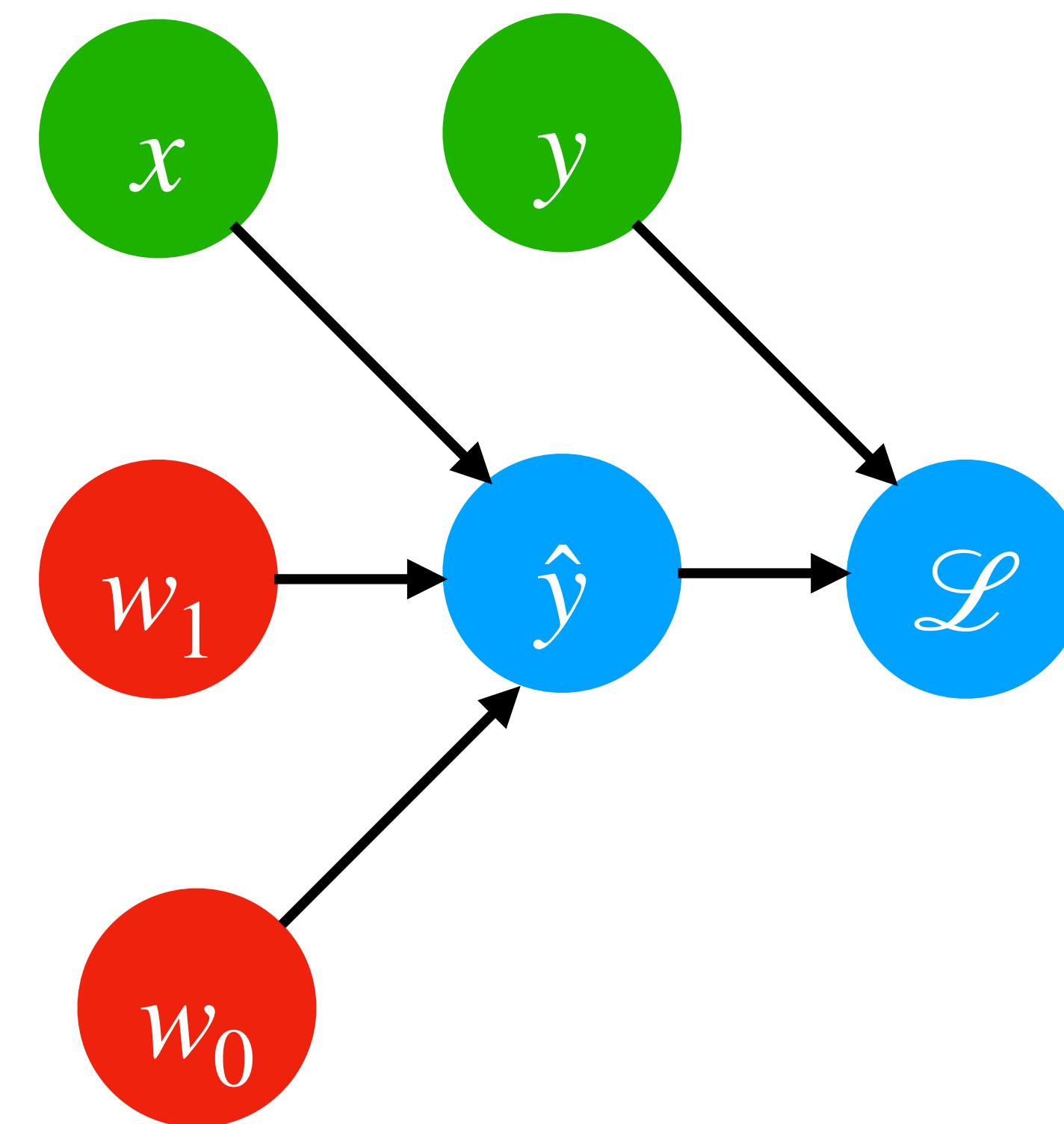
A **computation graph** has three kinds of nodes:

- Input nodes
- Parameter nodes
- Compute nodes

**Example:** Linear Regression

$$(1) \quad \hat{y} = w_0 + w_1 x$$

$$(2) \quad \mathcal{L} = (\hat{y} - y)^2$$



# Computation Graphs

A **computation graph** has three kinds of nodes:

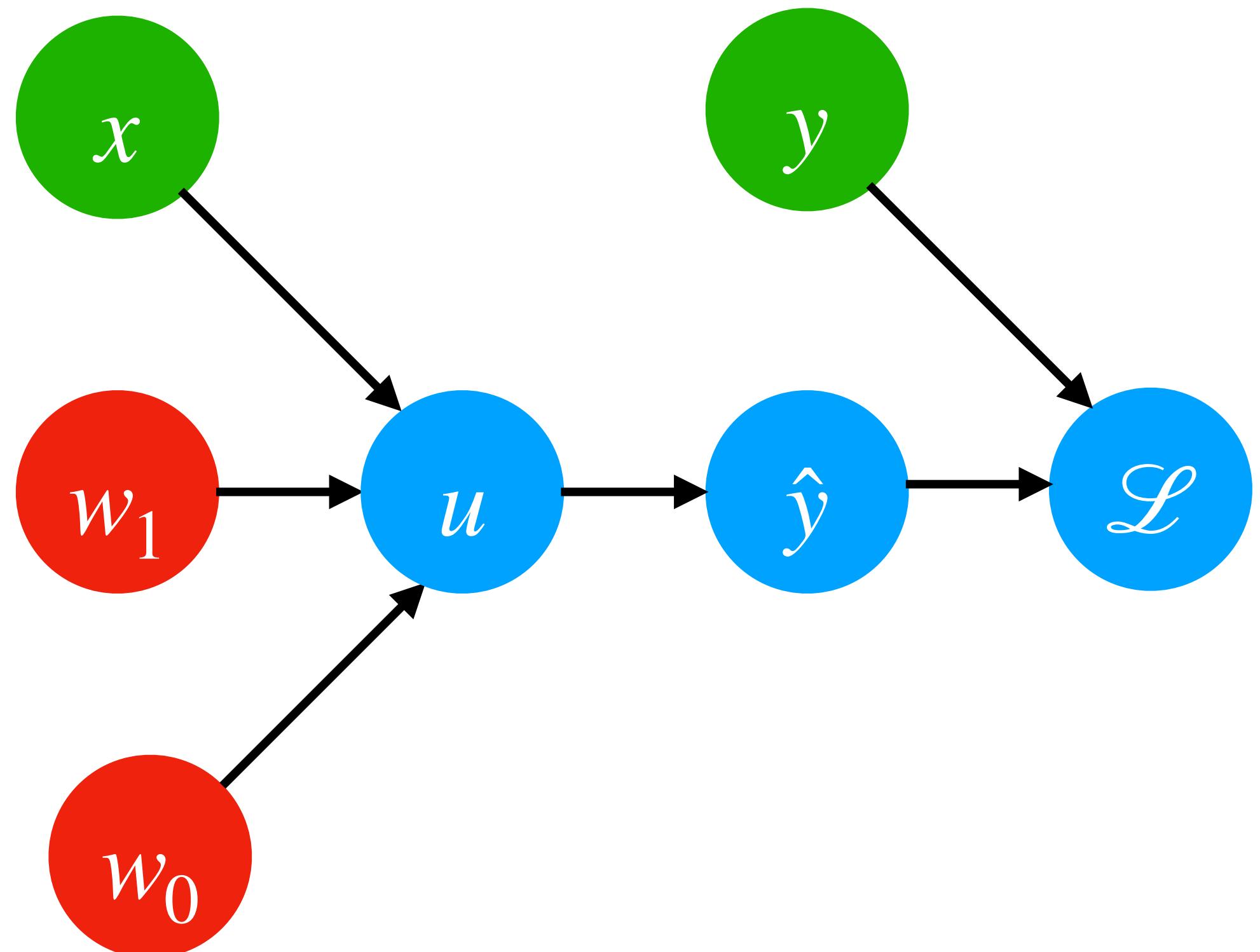
- Input nodes
- Parameter nodes
- Compute nodes

**Example:** Logistic Regression

$$(1) \quad u = w_0 + w_1 x$$

$$(2) \quad \hat{y} = \sigma(u)$$

$$(3) \quad \mathcal{L} = -y \log \hat{y} - (1 - y) \log (1 - \hat{y})$$



# Computation Graphs

A **computation graph** has three kinds of nodes:

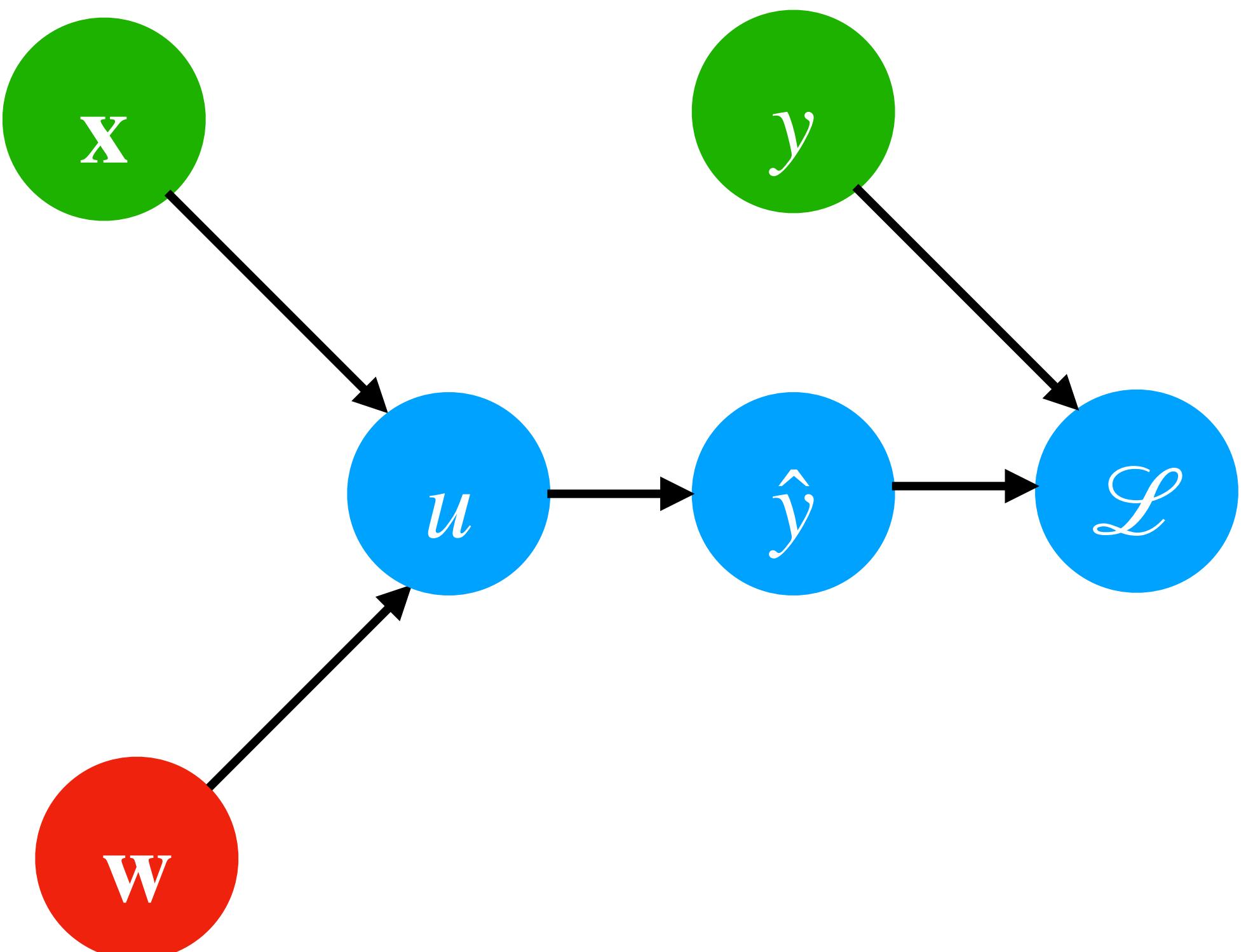
- Input nodes
- Parameter nodes
- Compute nodes

**Example:** Logistic Regression

$$(1) \quad u = \mathbf{w}^T \mathbf{x}$$

$$(2) \quad \hat{y} = \sigma(u)$$

$$(3) \quad \mathcal{L} = -y \log \hat{y} - (1 - y) \log (1 - \hat{y})$$



# Computation Graphs

A **computation graph** has three kinds of nodes:

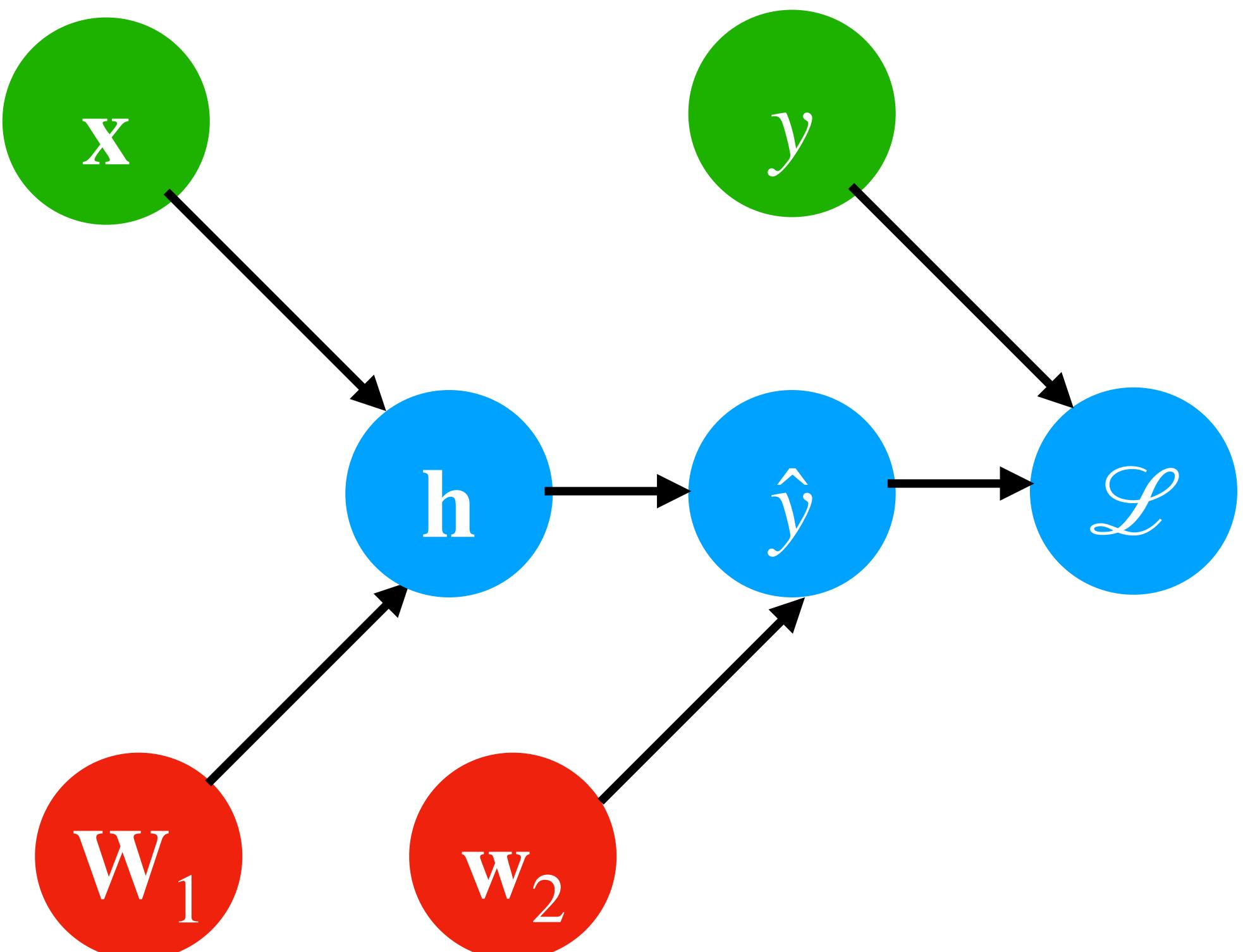
- Input nodes
- Parameter nodes
- Compute nodes

**Example:** Multi-Layer Perceptron

$$(1) \quad \mathbf{h} = \sigma(\mathbf{W}_1^T \mathbf{x})$$

$$(2) \quad \hat{y} = \sigma(\mathbf{w}_2^T \mathbf{h})$$

$$(3) \quad \mathcal{L} = -y \log \hat{y} - (1 - y) \log (1 - \hat{y})$$



# Backpropagation

# Backpropagation

**Goal:** Find gradients of negative log likelihood

$$\nabla_{\mathbf{w}} \sum_{i=1}^N -\log p_{model}(y_i | \mathbf{x}_i, \mathbf{w})$$

$\mathcal{L}(y_i, \mathbf{x}_i, \mathbf{w})$

or more generally of a loss function

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{y}, \mathbf{X}, \mathbf{w}) = \nabla_{\mathbf{w}} \sum_{i=1}^N \mathcal{L}(y_i, \mathbf{x}_i, \mathbf{w}) = \sum_{i=1}^N \nabla_{\mathbf{w}} \mathcal{L}(y_i, \mathbf{x}_i, \mathbf{w})$$

given a dataset  $\mathcal{X} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$  with  $N$  elements. In the following, we consider the computation of gradients wrt. a single data point:  $\nabla_{\mathbf{w}} \mathcal{L}(y_i, \mathbf{x}_i, \mathbf{w})$ . The gradient with respect to the entire dataset  $\mathcal{X}$  is obtained by summing up all individual gradients.

# Chain Rule

**Chain Rule:**

$$\frac{d}{dx} f(g(x)) = \frac{df}{dg}(g) \frac{dg}{dx}(x) = \frac{df}{dg} \frac{dg}{dx}$$

**Multivariate Chain Rule:**

$$\frac{d}{dx} f(g_1(x), \dots, g_M(x)) = \sum_{i=1}^M \frac{df}{dg_i} (g_1(x), \dots, g_M(x)) \frac{dg_i}{dx}(x) = \sum_{i=1}^M \frac{df}{dg_i} \frac{dg_i}{dx}$$

# Backpropagation

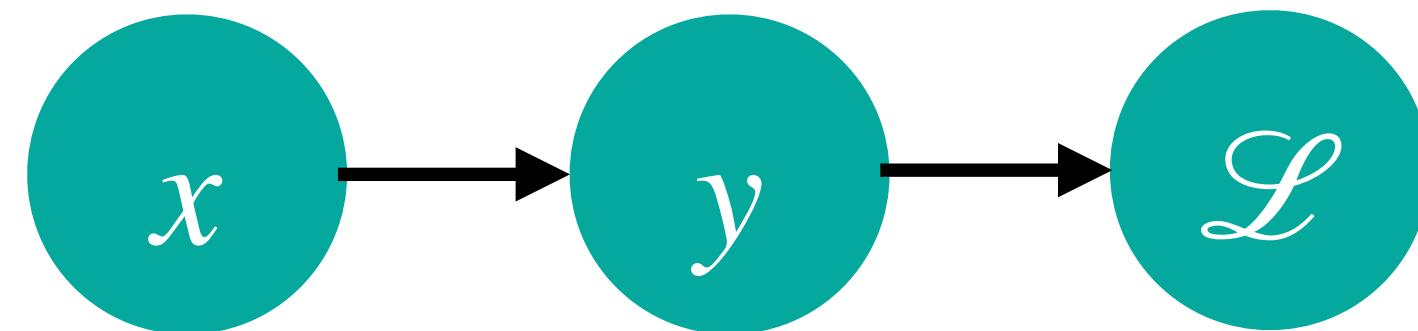
For now: no distinction between node types (input, parameter, compute)

**Forward Pass:**

$$(1) \quad y = x^2$$

$$(2) \quad \mathcal{L} = 2y$$

**Loss:**  $\mathcal{L} = 2x^2$



# Backpropagation

For now: no distinction between node types (input, parameter, compute)

**Forward Pass:**

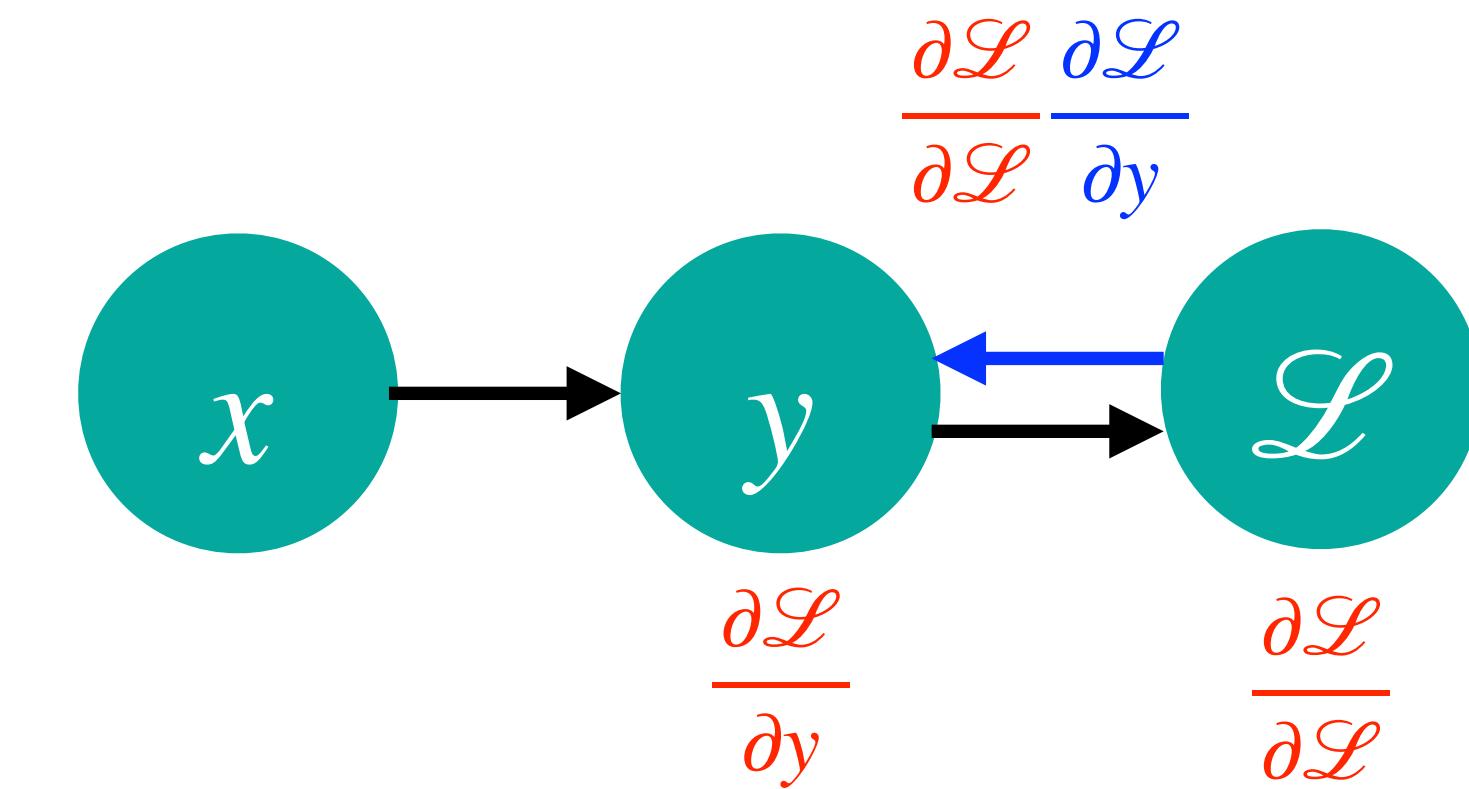
$$(1) \quad y = x^2$$

$$(2) \quad \mathcal{L} = 2y$$

**Backward Pass:**

$$(2) \quad \frac{\partial \mathcal{L}}{\partial y} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial y} = 2$$

**Loss:**  $\mathcal{L} = 2x^2$



◆ **Red:** back-propagated gradients

◆ **Blue:** local gradients

# Backpropagation

For now: no distinction between node types (input, parameter, compute)

**Forward Pass:**

$$(1) \quad y = x^2$$

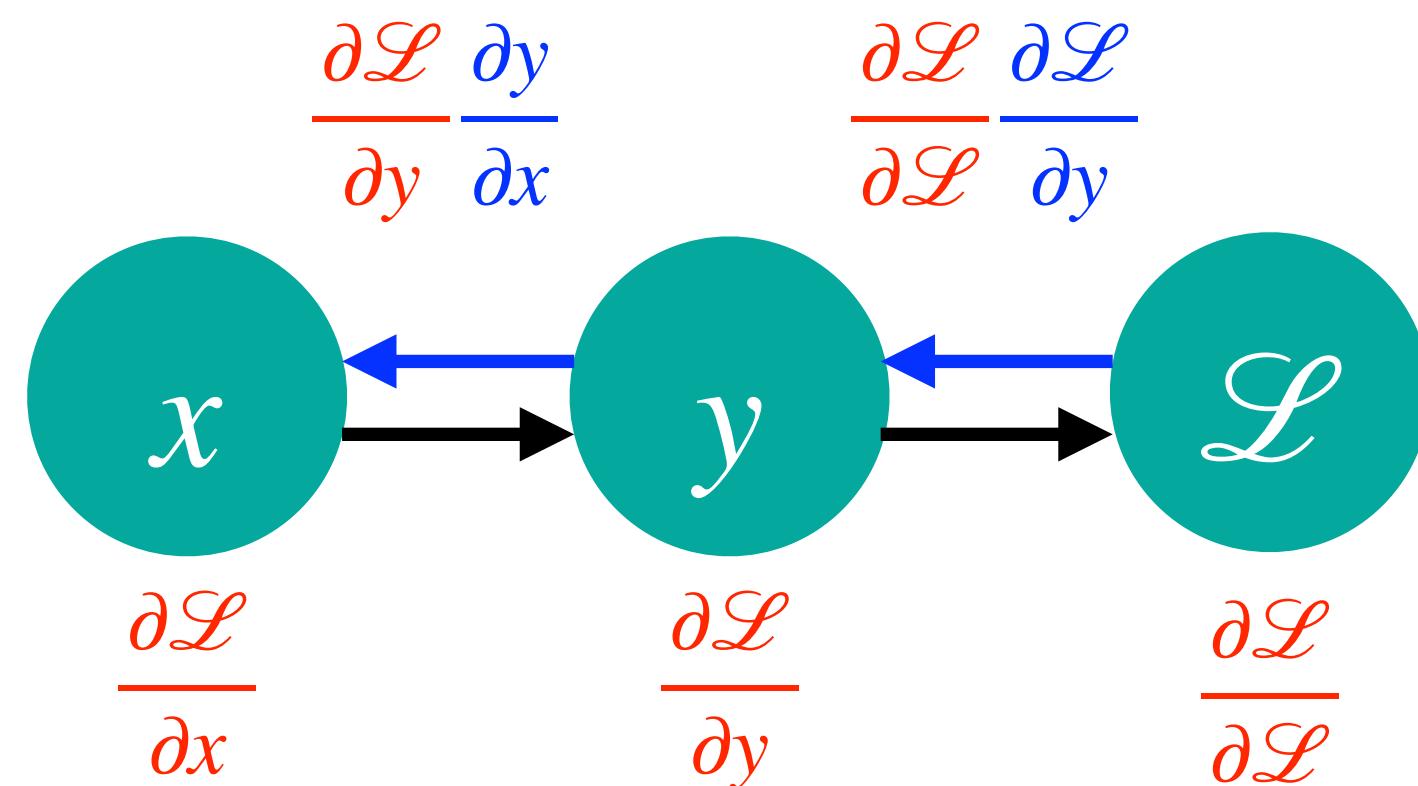
$$(2) \quad \mathcal{L} = 2y$$

**Backward Pass:**

$$(2) \quad \frac{\partial \mathcal{L}}{\partial y} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial y} = 2$$

$$(1) \quad \frac{\partial \mathcal{L}}{\partial x} = \frac{\partial \mathcal{L}}{\partial y} \frac{\partial y}{\partial x} = \frac{\partial \mathcal{L}}{\partial y} 2x$$

**Loss:**  $\mathcal{L} = 2x^2$



♦ **Red:** back-propagated gradients

♦ **Blue:** local gradients

# Backpropagation

For now: no distinction between node types (input, parameter, compute)

**Forward Pass:**

$$(1) \quad y = x^2$$

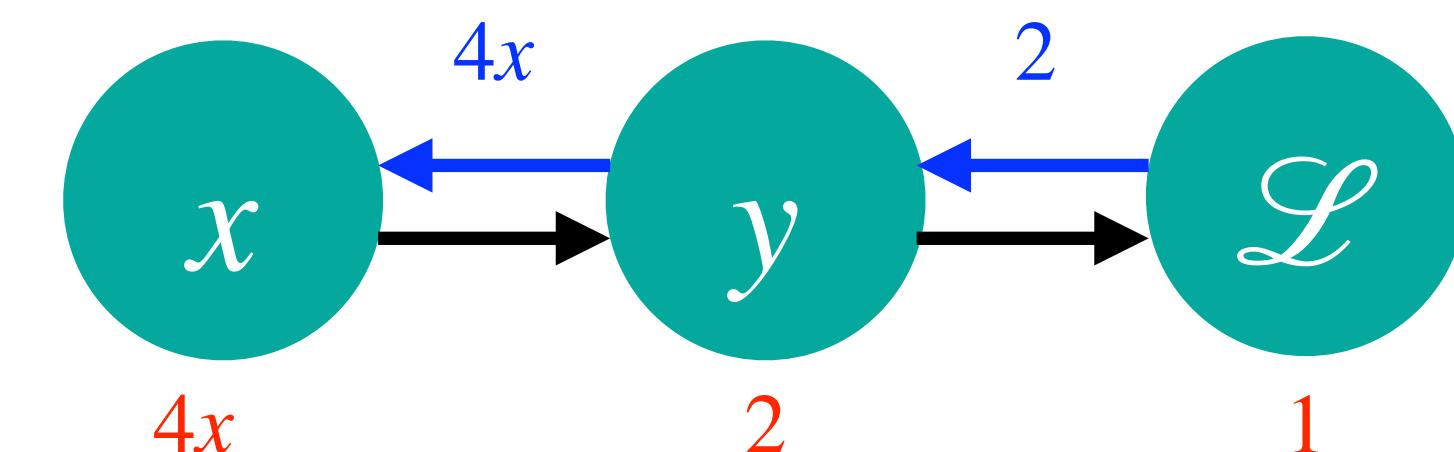
$$(2) \quad \mathcal{L} = 2y$$

**Backward Pass:**

$$(2) \quad \frac{\partial \mathcal{L}}{\partial y} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial y} = 2$$

$$(1) \quad \frac{\partial \mathcal{L}}{\partial x} = \frac{\partial \mathcal{L}}{\partial y} \frac{\partial y}{\partial x} = \frac{\partial \mathcal{L}}{\partial y} 2x$$

**Loss:**  $\mathcal{L} = 2x^2$



♦ **Red:** back-propagated gradients

♦ **Blue:** local gradients

# Backpropagation: Abstract Case

For now: no distinction between node types (input, parameter, compute)

**Forward Pass:**

$$(1) \quad y = y(x)$$

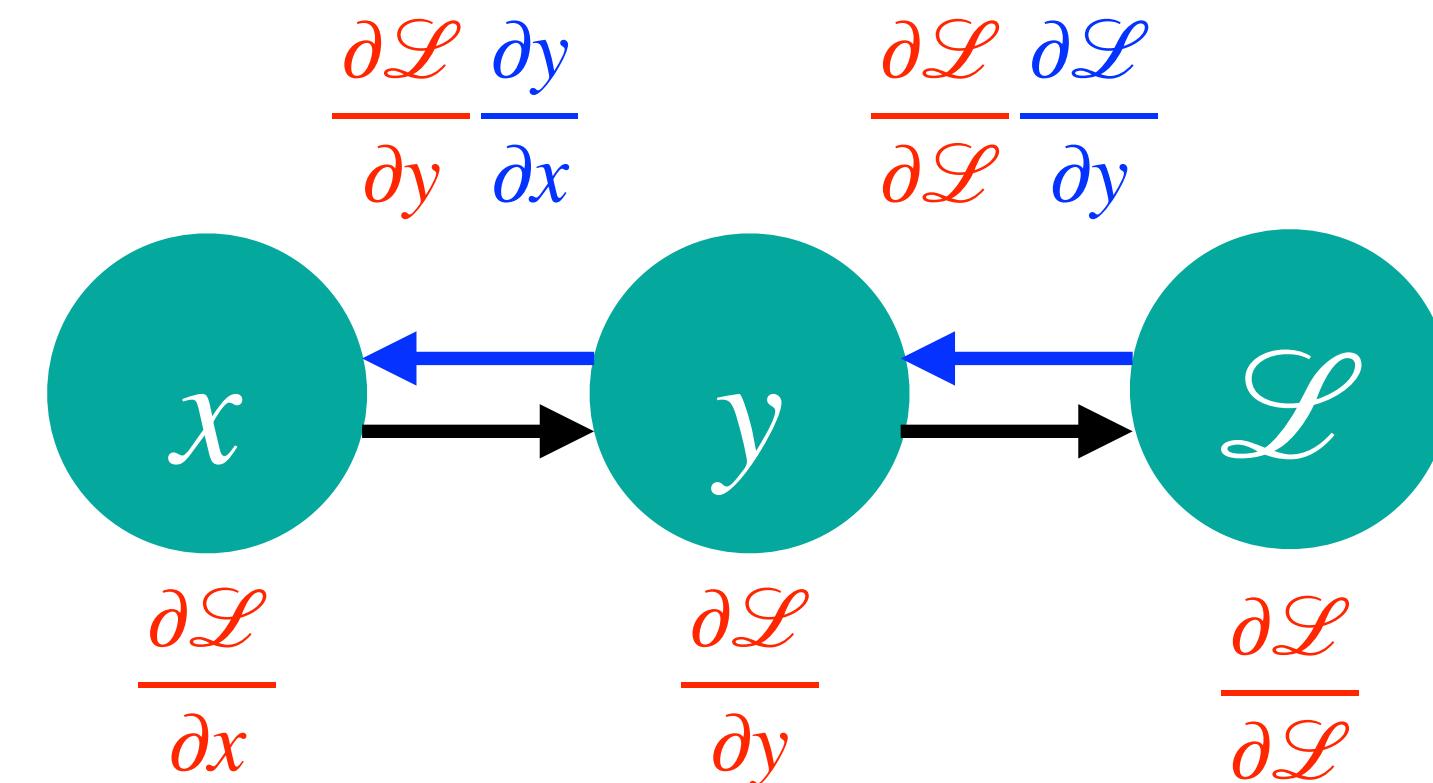
$$(2) \quad \mathcal{L} = \mathcal{L}(y)$$

**Backward Pass:**

$$(2) \quad \frac{\partial \mathcal{L}}{\partial y} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial y} = \frac{\partial \mathcal{L}}{\partial y}$$

$$(1) \quad \frac{\partial \mathcal{L}}{\partial x} = \frac{\partial \mathcal{L}}{\partial y} \frac{\partial y}{\partial x}$$

**Loss:**  $\mathcal{L}(y(x))$



♦ **Red:** back-propagated gradients

♦ **Blue:** local gradients

# Backpropagation: Fan-Out > 1

**Forward Pass:**

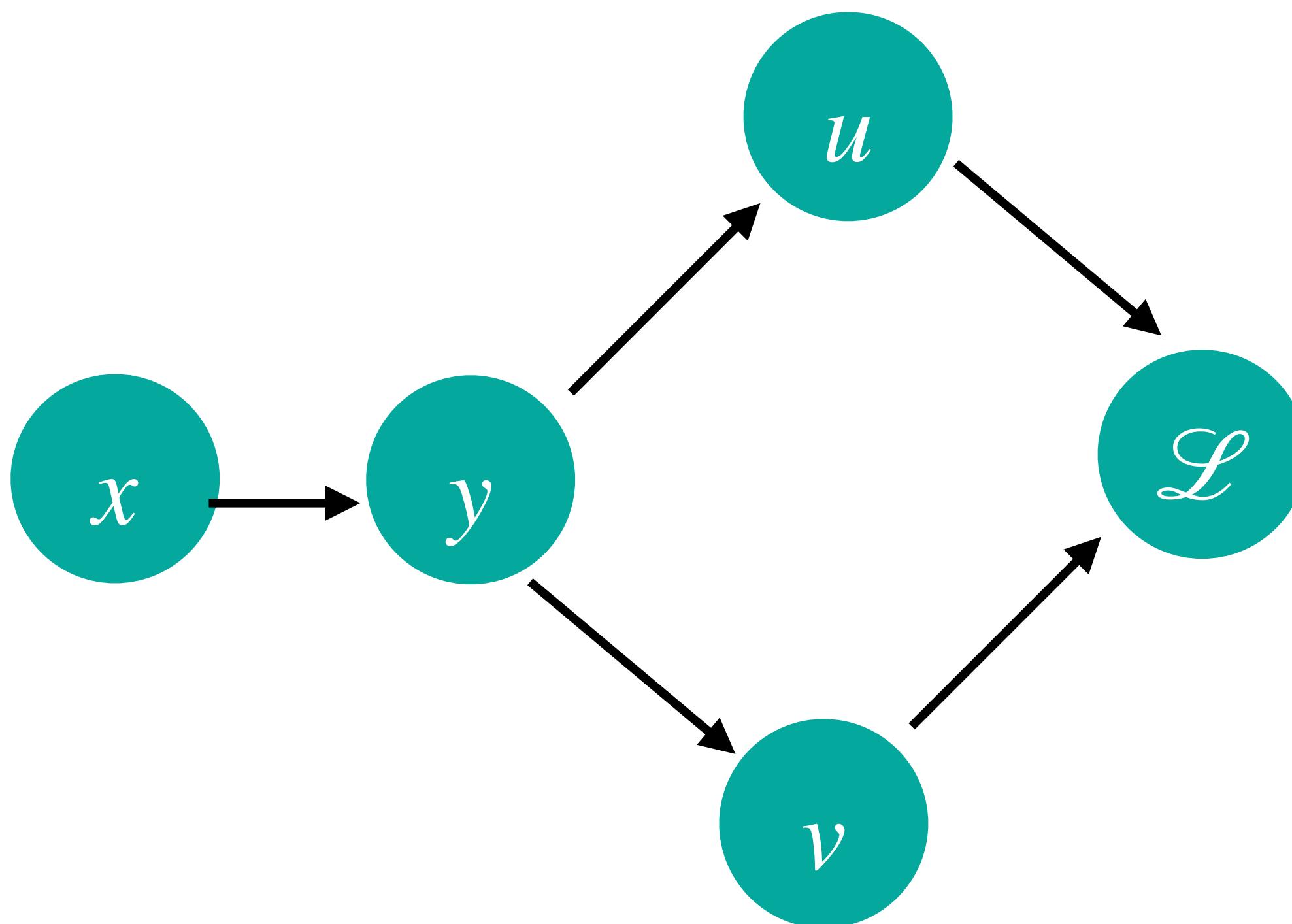
$$(1) \quad y = y(x)$$

$$(2) \quad u = u(y)$$

$$(3) \quad v = v(y)$$

$$(4) \quad \mathcal{L} = \mathcal{L}(u, v)$$

**Loss:**  $\mathcal{L} \left( u \left( y(x) \right), v \left( y(x) \right) \right)$



# Backpropagation: Fan-Out > 1

## Forward Pass:

$$(1) \quad y = y(x)$$

$$(2) \quad u = u(y)$$

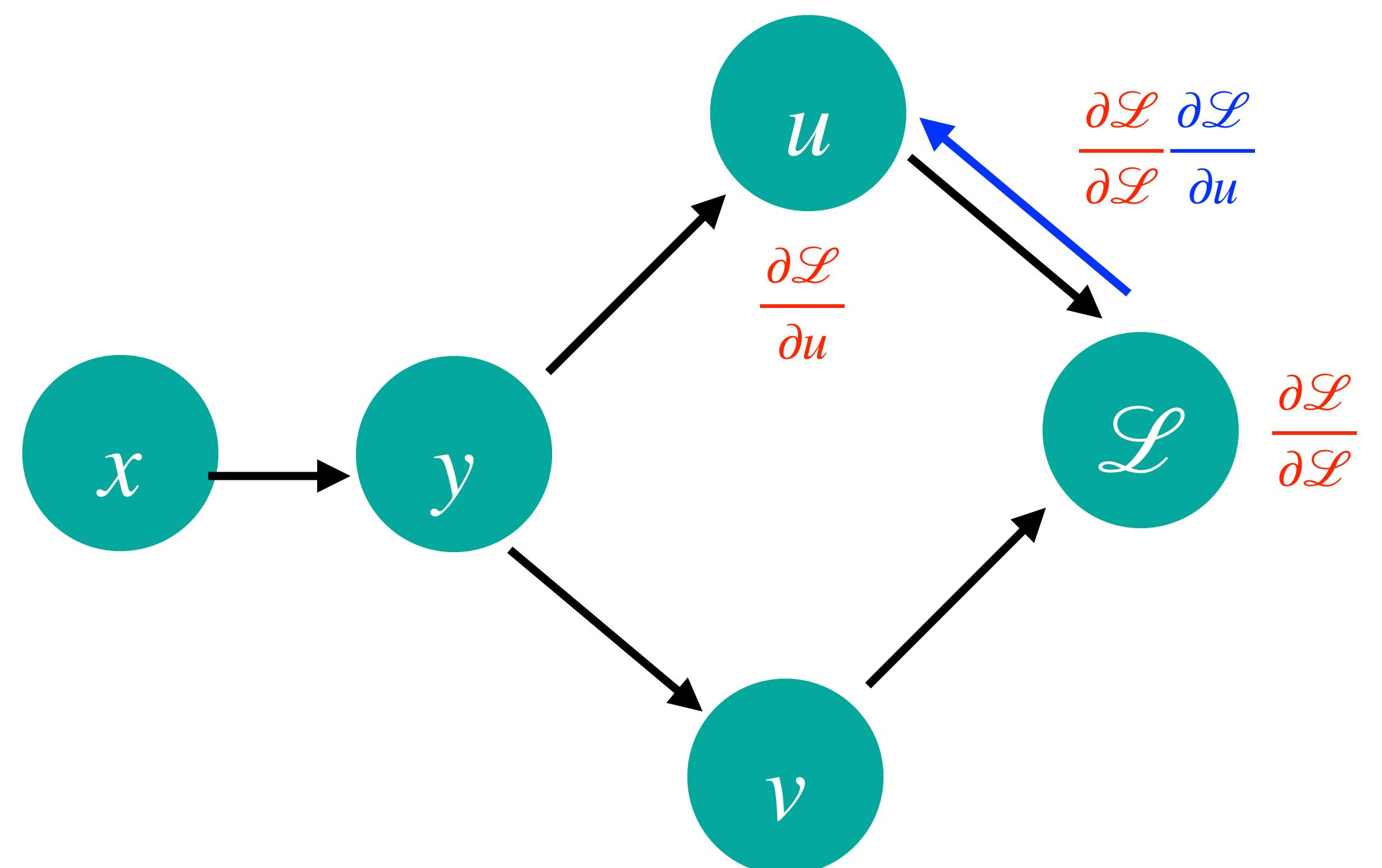
$$(3) \quad v = v(y)$$

$$(4) \quad \mathcal{L} = \mathcal{L}(u, v)$$

## Backward Pass:

$$(4) \quad \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial u}$$

**Loss:**  $\mathcal{L} \left( u(y(x)), v(y(x)) \right)$



# Backpropagation: Fan-Out > 1

## Forward Pass:

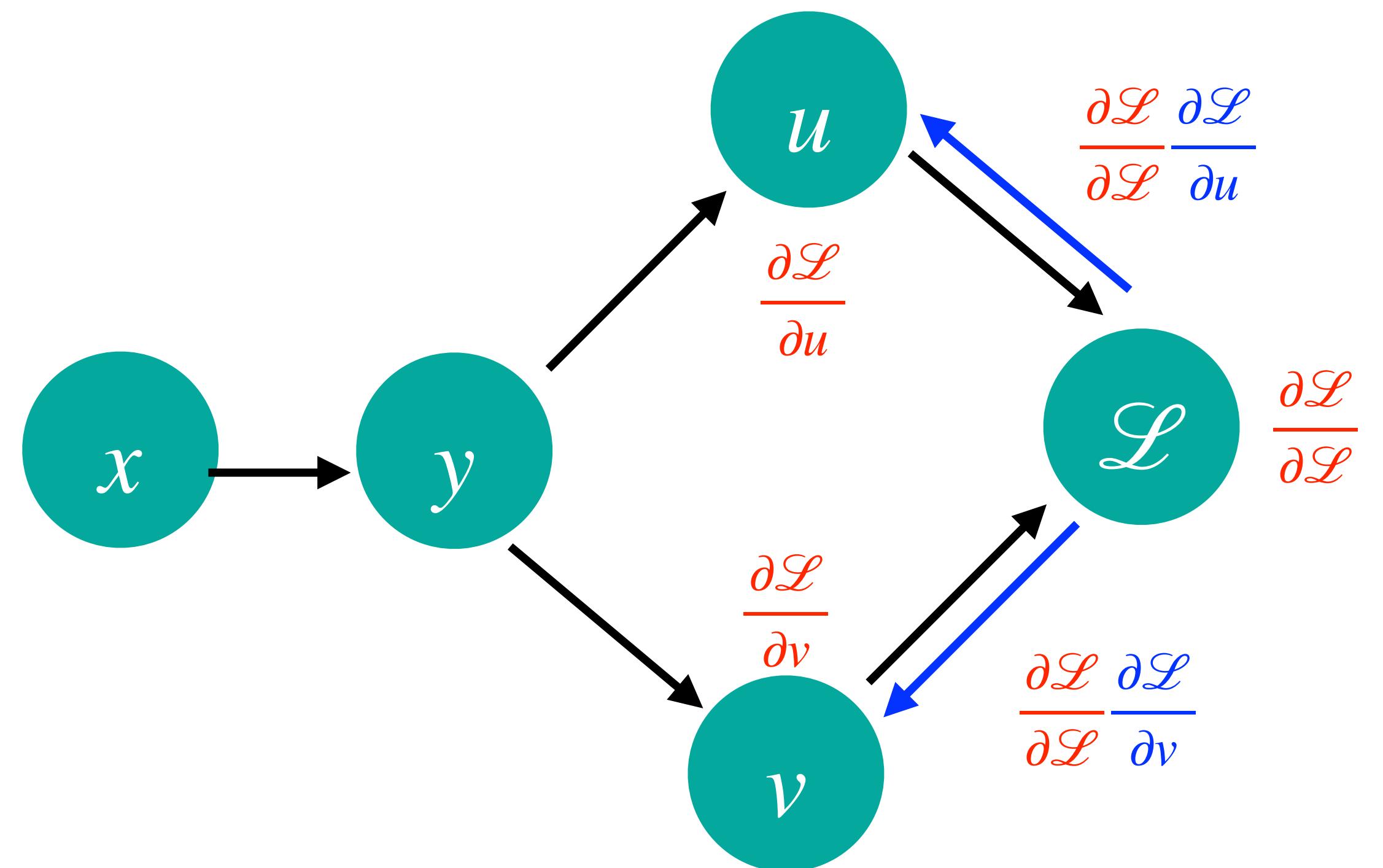
$$(1) \quad y = y(x)$$

$$(2) \quad u = u(y)$$

$$(3) \quad v = v(y)$$

$$(4) \quad \mathcal{L} = \mathcal{L}(u, v)$$

**Loss:**  $\mathcal{L} \left( u(y(x)), v(y(x)) \right)$



## Backward Pass:

$$(4) \quad \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial u}$$

$$(3) \quad \frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial v}$$

# Backpropagation: Fan-Out > 1

## Forward Pass:

$$(1) \quad y = y(x)$$

$$(2) \quad u = u(y)$$

$$(3) \quad v = v(y)$$

$$(4) \quad \mathcal{L} = \mathcal{L}(u, v)$$

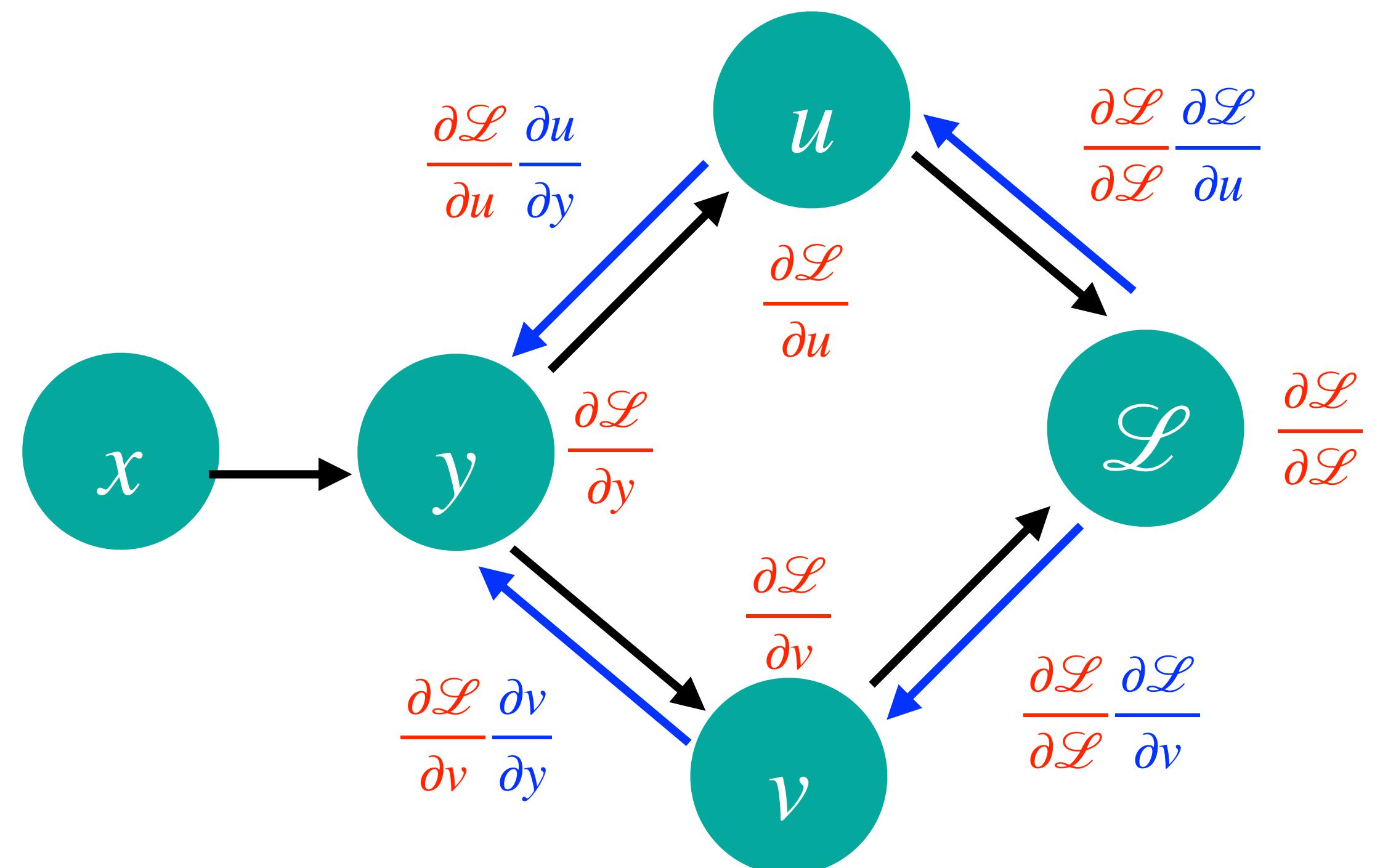
## Backward Pass:

$$(4) \quad \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial u}$$

$$(3) \quad \frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial v}$$

$$(2) \quad \frac{\partial \mathcal{L}}{\partial y} =$$

**Loss:**  $\mathcal{L} \left( u(y(x)), v(y(x)) \right)$



# Backpropagation: Fan-Out > 1

## Forward Pass:

$$(1) \quad y = y(x)$$

$$(2) \quad u = u(y)$$

$$(3) \quad v = v(y)$$

$$(4) \quad \mathcal{L} = \mathcal{L}(u, v)$$

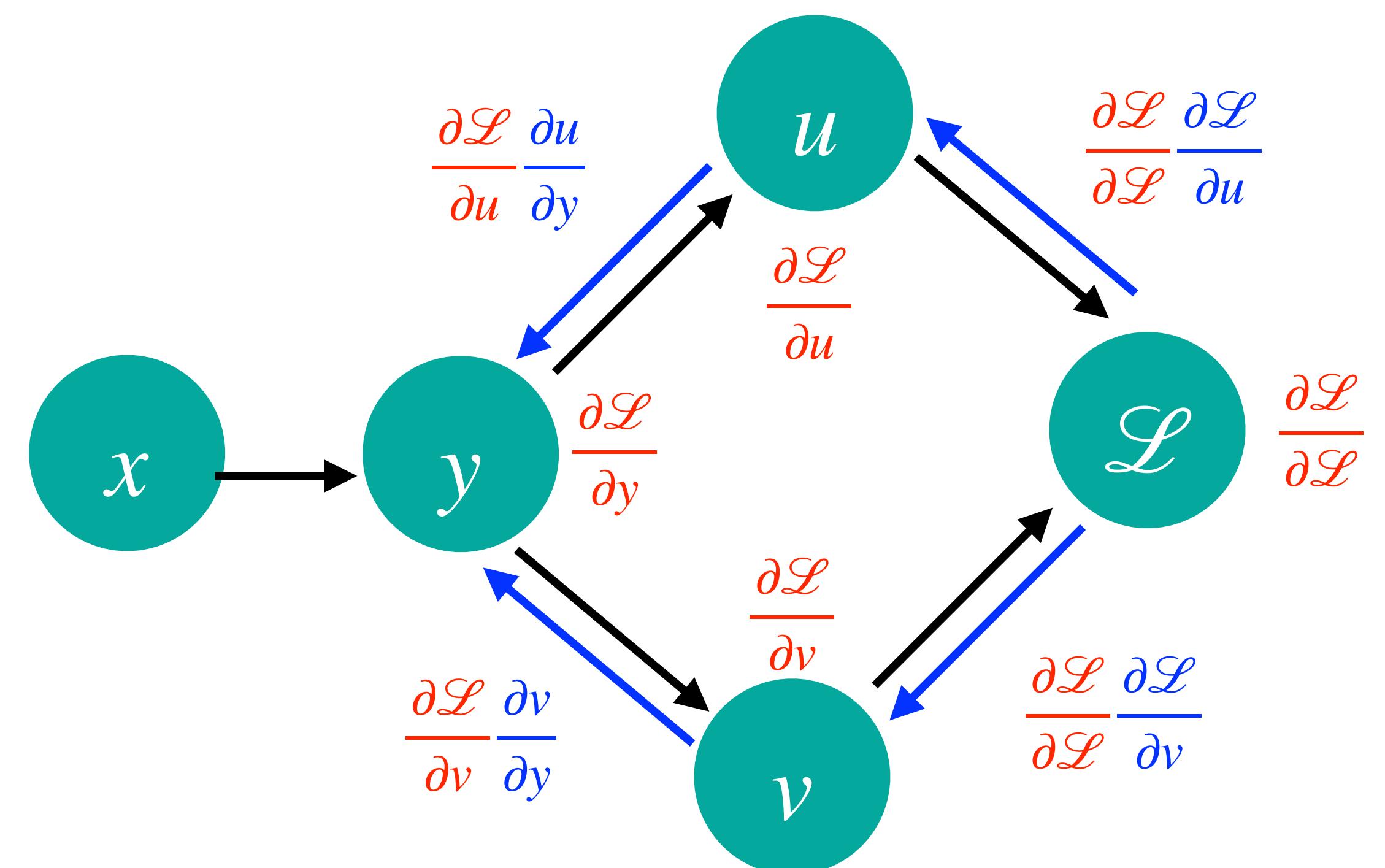
## Backward Pass:

$$(4) \quad \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial u}$$

$$(3) \quad \frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial v}$$

$$(2) \quad \frac{\partial \mathcal{L}}{\partial y} =$$

**Loss:**  $\mathcal{L} (u(y), v(y))$



$$\frac{d}{dy} \mathcal{L} (u(y), v(y)) = ?$$

# Backpropagation: Fan-Out > 1

## Forward Pass:

$$(1) \quad y = y(x)$$

$$(2) \quad u = u(y)$$

$$(3) \quad v = v(y)$$

$$(4) \quad \mathcal{L} = \mathcal{L}(u, v)$$

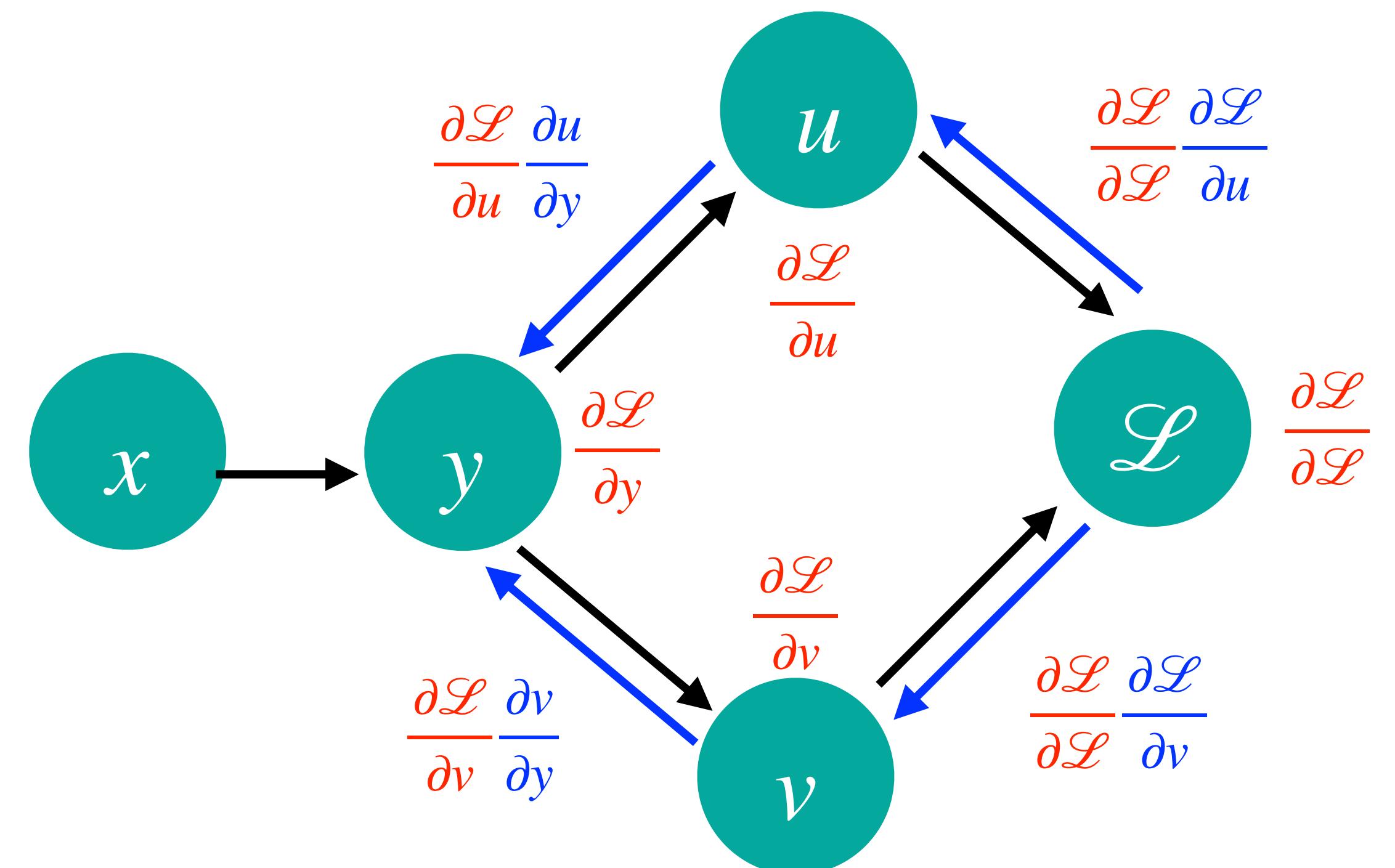
## Backward Pass:

$$(4) \quad \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial u}$$

$$(3) \quad \frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial v}$$

$$(2) \quad \frac{\partial \mathcal{L}}{\partial y} =$$

**Loss:**  $\mathcal{L} (u(y), v(y))$



$$\frac{d}{dy} \mathcal{L} (u(y), v(y)) = ?$$

# Backpropagation: Fan-Out > 1

## Forward Pass:

$$(1) \quad y = y(x)$$

$$(2) \quad u = u(y)$$

$$(3) \quad v = v(y)$$

$$(4) \quad \mathcal{L} = \mathcal{L}(u, v)$$

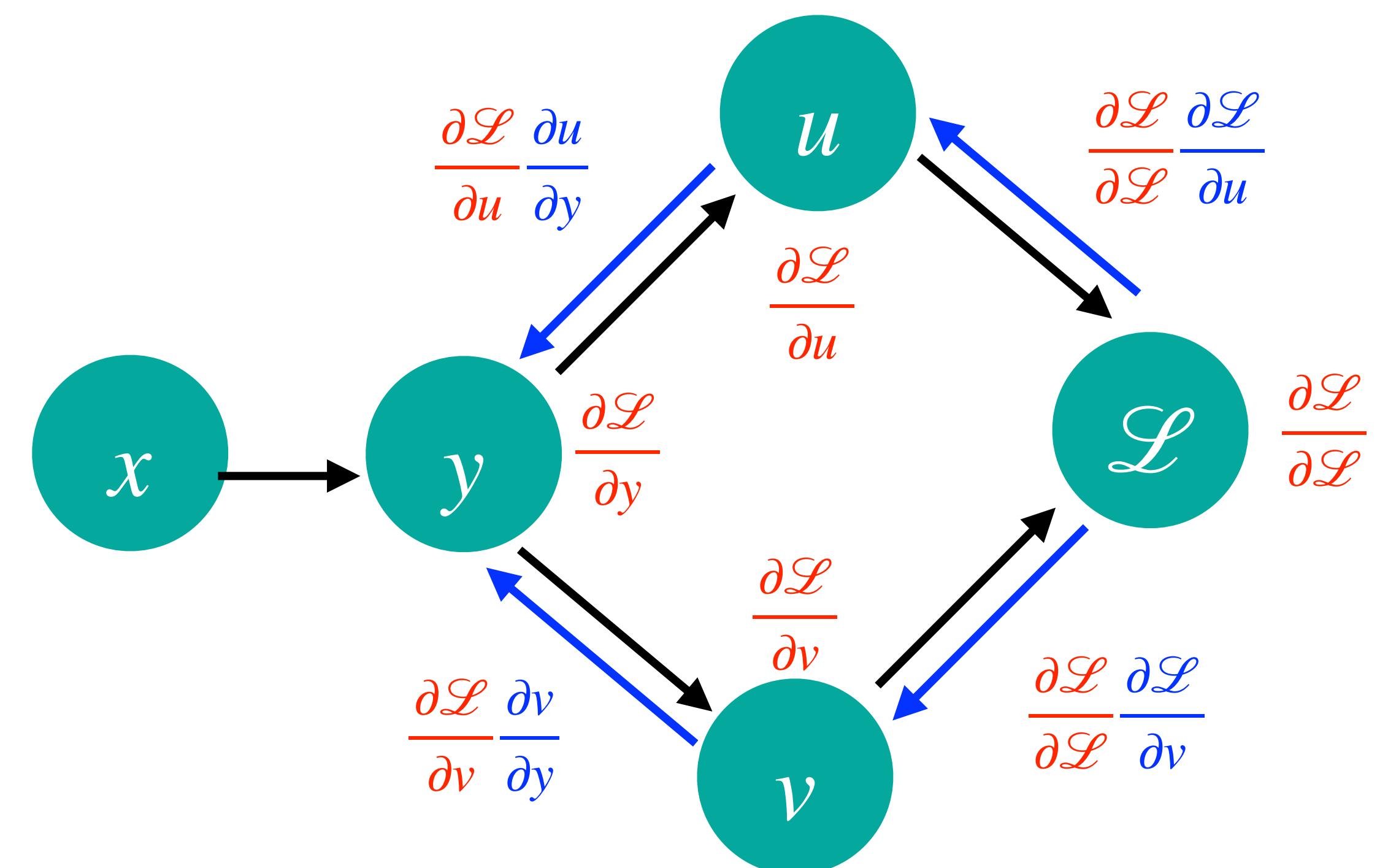
## Backward Pass:

$$(4) \quad \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial u}$$

$$(3) \quad \frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial v}$$

$$(2) \quad \frac{\partial \mathcal{L}}{\partial y} = \frac{\partial \mathcal{L}}{\partial u} \frac{\partial u}{\partial y} + \frac{\partial \mathcal{L}}{\partial v} \frac{\partial v}{\partial y}$$

**Loss:**  $\mathcal{L} (u(y), v(y))$



$$\frac{d}{dy} \mathcal{L} (u(y), v(y)) = \frac{\partial \mathcal{L}}{\partial u} \frac{\partial u}{\partial y} + \frac{\partial \mathcal{L}}{\partial v} \frac{\partial v}{\partial y}$$

All incoming gradients must be summed up!

# Backpropagation: Fan-Out > 1

## Forward Pass:

$$(1) \quad y = y(x)$$

$$(2) \quad u = u(y)$$

$$(3) \quad v = v(y)$$

$$(4) \quad \mathcal{L} = \mathcal{L}(u, v)$$

## Backward Pass:

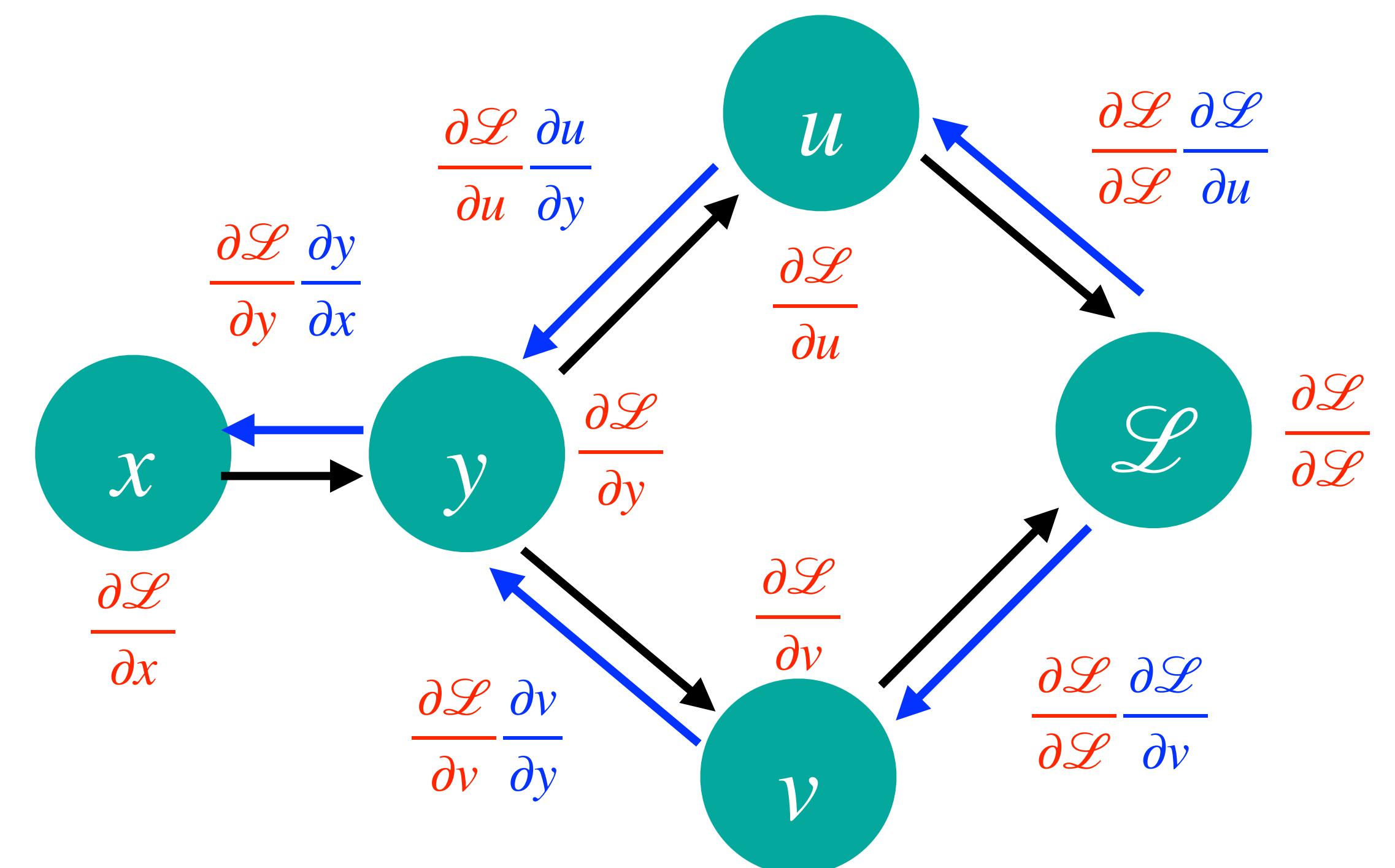
$$(4) \quad \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial u}$$

$$(3) \quad \frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial v}$$

$$(2) \quad \frac{\partial \mathcal{L}}{\partial y} = \frac{\partial \mathcal{L}}{\partial u} \frac{\partial u}{\partial y} + \frac{\partial \mathcal{L}}{\partial v} \frac{\partial v}{\partial y}$$

$$(1) \quad \frac{\partial \mathcal{L}}{\partial x} = \frac{\partial \mathcal{L}}{\partial y} \frac{\partial y}{\partial x}$$

**Loss:**  $\mathcal{L} (u(y(x)), v(y(x)))$



$$\frac{d}{dy} \mathcal{L} (u(y), v(y)) = \frac{\partial \mathcal{L}}{\partial u} \frac{\partial u}{\partial y} + \frac{\partial \mathcal{L}}{\partial v} \frac{\partial v}{\partial y}$$

All incoming gradients must be summed up!

# Backpropagation: Fan-Out > 1

## Forward Pass:

$$(1) \quad y = y(x)$$

$$(2) \quad u = u(y)$$

$$(3) \quad v = v(y)$$

$$(4) \quad \mathcal{L} = \mathcal{L}(u, v)$$

## Backward Pass:

$$(4) \quad \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial u}$$

$$(3) \quad \frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial v}$$

$$(2) \quad \frac{\partial \mathcal{L}}{\partial y} = \frac{\partial \mathcal{L}}{\partial u} \frac{\partial u}{\partial y} + \frac{\partial \mathcal{L}}{\partial v} \frac{\partial v}{\partial y}$$

$$(1) \quad \frac{\partial \mathcal{L}}{\partial x} = \frac{\partial \mathcal{L}}{\partial y} \frac{\partial y}{\partial x}$$

**Implementation:** Each variable/node is an object and has attributes `x.value` and `x.grad`. Values are computed **forward**:

`x.value` = Input

`y.value` =  $y(x.value)$

`u.value` =  $u(y.value)$

`v.value` =  $v(y.value)$

`L.value` =  $\mathcal{L}(u.value, v.value)$

# Backpropagation: Fan-Out > 1

## Forward Pass:

$$(1) \quad y = y(x)$$

$$(2) \quad u = u(y)$$

$$(3) \quad v = v(y)$$

$$(4) \quad \mathcal{L} = \mathcal{L}(u, v)$$

## Backward Pass:

$$(4) \quad \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial u}$$

$$(3) \quad \frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial v}$$

$$(2) \quad \frac{\partial \mathcal{L}}{\partial y} = \frac{\partial \mathcal{L}}{\partial u} \frac{\partial u}{\partial y} + \frac{\partial \mathcal{L}}{\partial v} \frac{\partial v}{\partial y}$$

$$(1) \quad \frac{\partial \mathcal{L}}{\partial x} = \frac{\partial \mathcal{L}}{\partial y} \frac{\partial y}{\partial x}$$

**Implementation:** Each variable/node is an object and has attributes `x.value` and `x.grad`. Values are computed **forward** and gradients **backward**:

`x.grad = y.grad = u.grad = v.grad = 0`

`L.grad = 1`

`u.grad += L.grad * (dL/du) (u.value, v.value)`

`v.grad += L.grad * (dL/dv) (u.value, v.value)`

`y.grad += u.grad * (du/dy) (y.value)`

`y.grad += v.grad * (dv/dy) (y.value)`

`x.grad += y.grad * (dy/dx) (x.value)`

# Backpropagation: Logistic Regression with 1D Inputs

**Forward Pass:**

$$(1) \quad u = w_0 + w_1 x$$

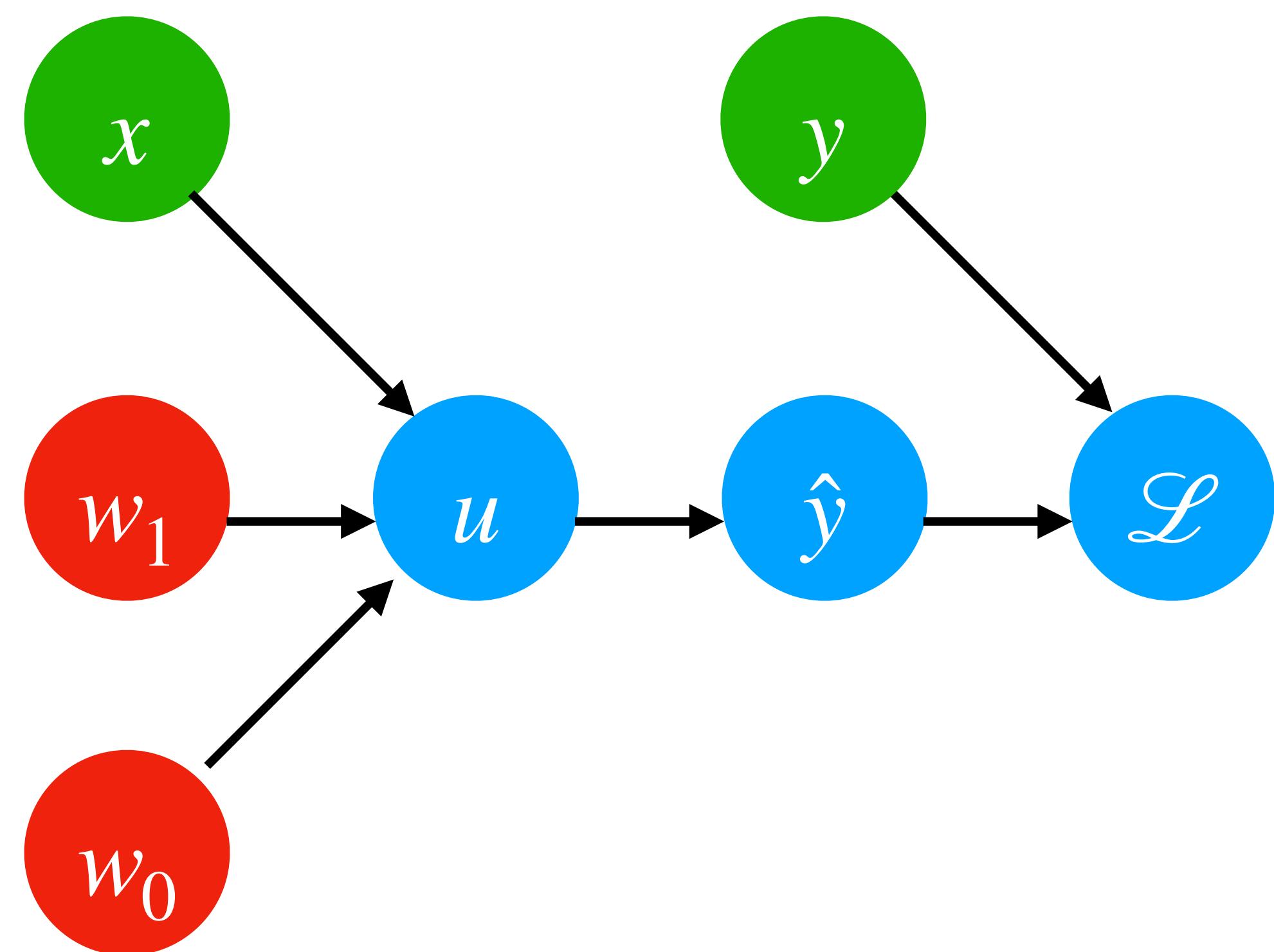
$$(2) \quad \hat{y} = \sigma(u)$$

$$(3) \quad \mathcal{L} = -y \log \hat{y} - (1 - y) \log (1 - \hat{y})$$

BCE( $\hat{y}, y$ )

**Loss:**  $\mathcal{L} = \text{BCE}(\sigma(w_0 + w_1 x), y)$

**Backward Pass:**



# Backpropagation: Logistic Regression with 1D Inputs

**Forward Pass:**

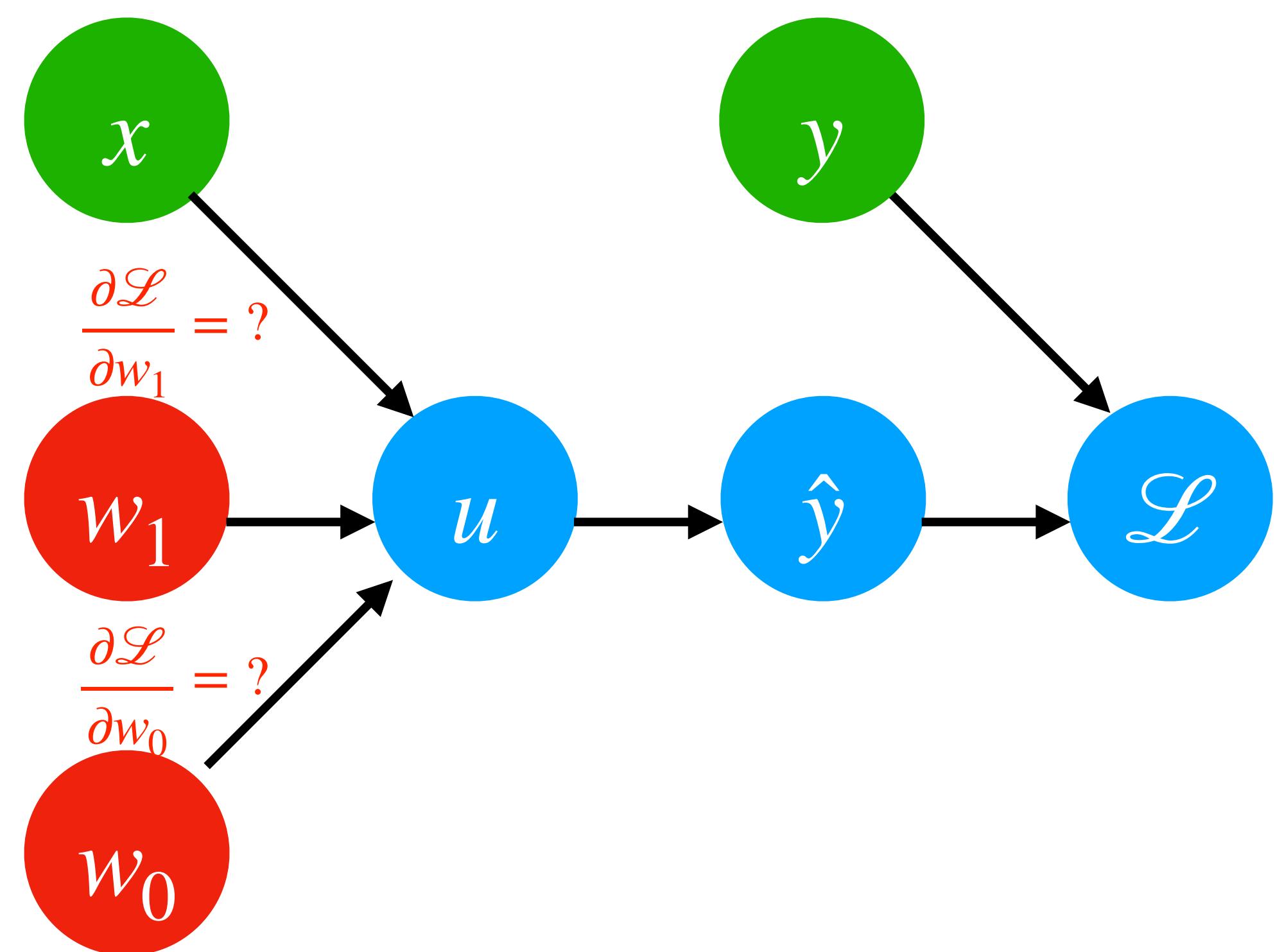
$$(1) \quad u = w_0 + w_1 x$$

$$(2) \quad \hat{y} = \sigma(u)$$

$$(3) \quad \mathcal{L} = -y \log \hat{y} - (1 - y) \log (1 - \hat{y})$$

**Loss:**  $\mathcal{L} = \text{BCE}(\sigma(w_0 + w_1 x), y)$

**Backward Pass:**



# Backpropagation: Logistic Regression with 1D Inputs

**Forward Pass:**

$$(1) \quad u = w_0 + w_1 x$$

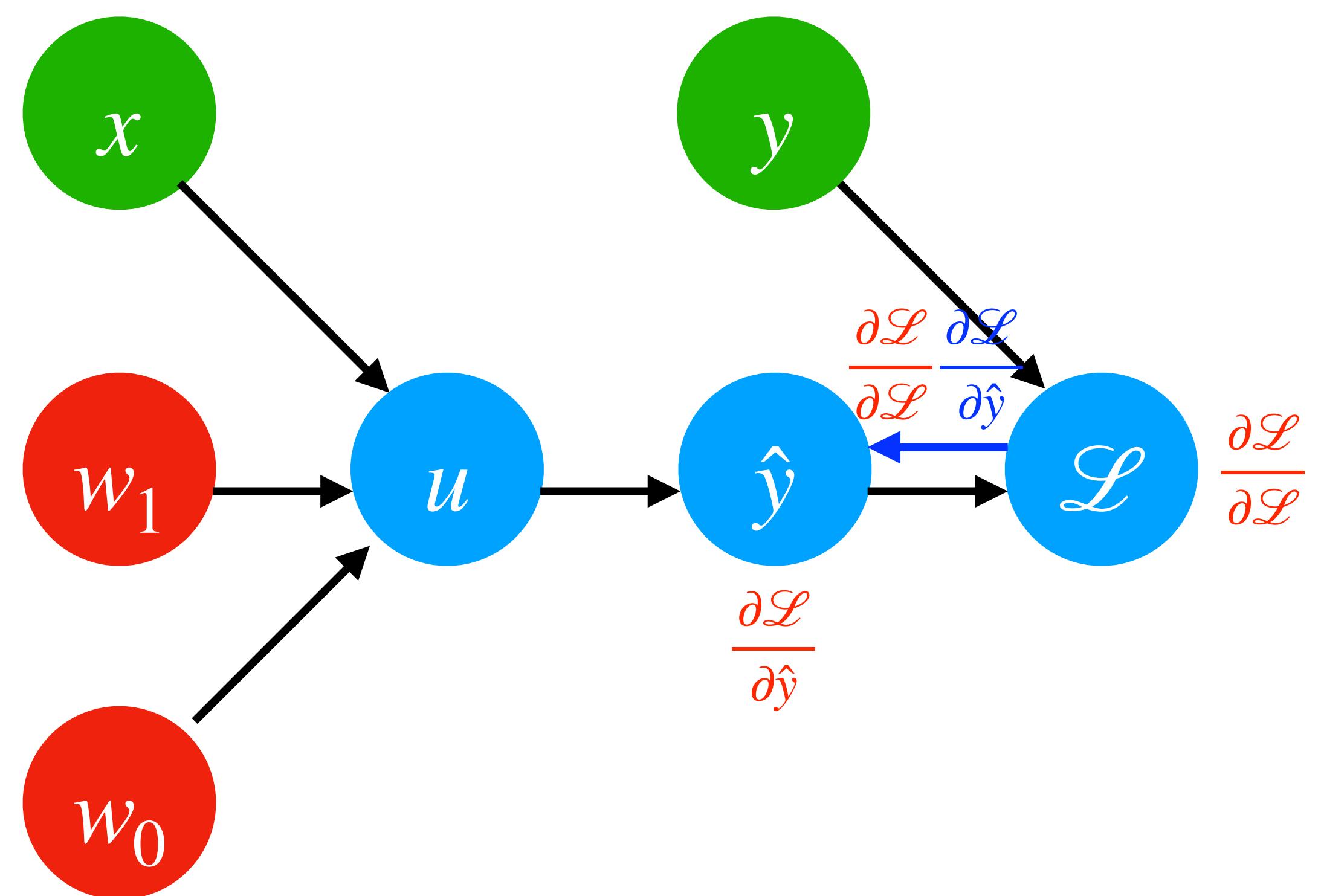
$$(2) \quad \hat{y} = \sigma(u)$$

$$(3) \quad \mathcal{L} = -y \log \hat{y} - (1 - y) \log (1 - \hat{y})$$

**Loss:**  $\mathcal{L} = \text{BCE}(\sigma(w_0 + w_1 x), y)$

**Backward Pass:**

$$(4) \quad \frac{\partial \mathcal{L}}{\partial \hat{y}} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial \hat{y}} = \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})}$$



# Backpropagation: Logistic Regression with 1D Inputs

**Forward Pass:**

$$(1) \quad u = w_0 + w_1 x$$

$$(2) \quad \hat{y} = \sigma(u)$$

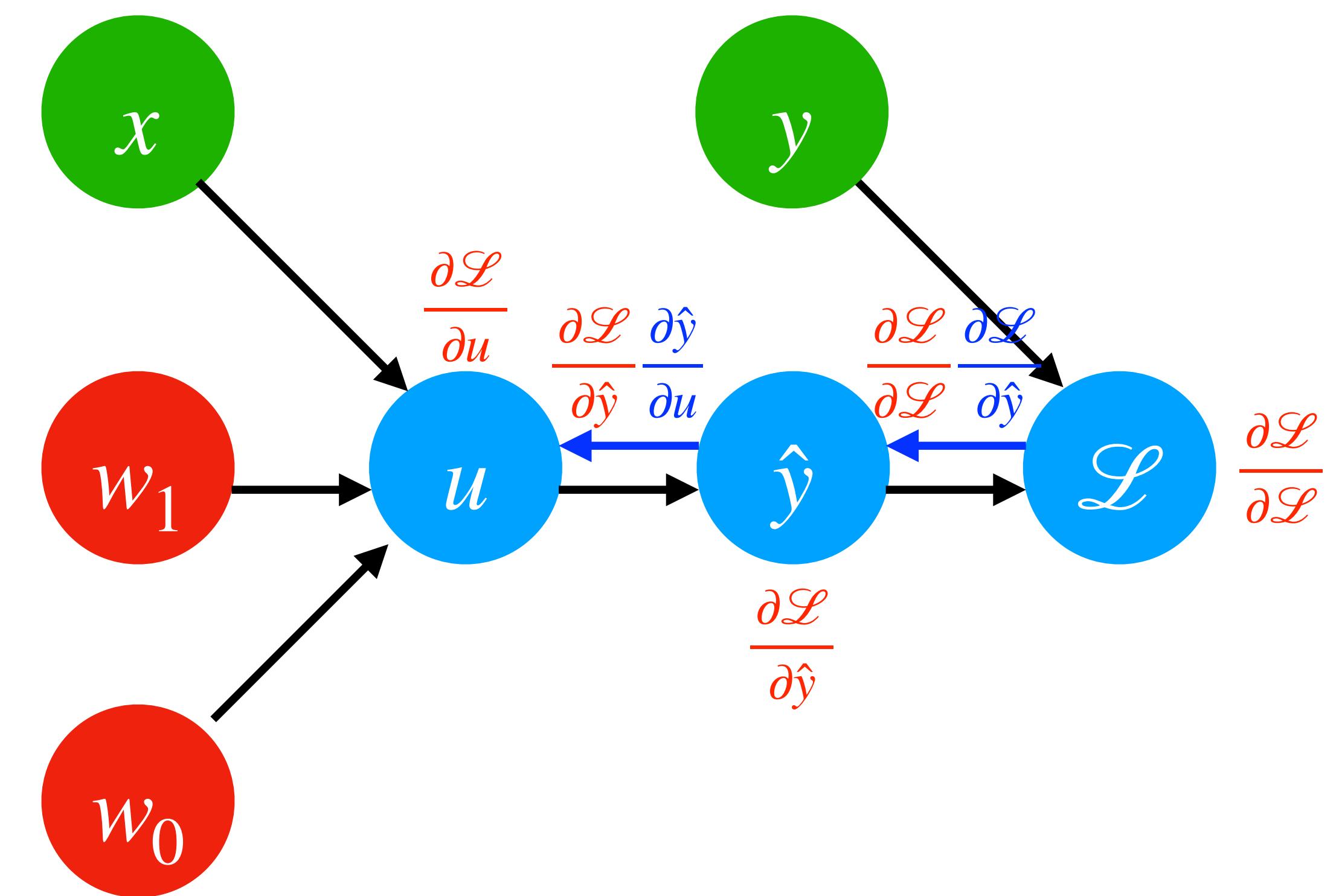
$$(3) \quad \mathcal{L} = -y \log \hat{y} - (1 - y) \log (1 - \hat{y})$$

**Loss:**  $\mathcal{L} = \text{BCE}(\sigma(w_0 + w_1 x), y)$

**Backward Pass:**

$$(4) \quad \frac{\partial \mathcal{L}}{\partial \hat{y}} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial \hat{y}} = \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})}$$

$$(5) \quad \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \sigma(u)(1 - \sigma(u))$$



# Backpropagation: Logistic Regression with 1D Inputs

**Forward Pass:**

$$(1) \quad u = w_0 + w_1 x$$

$$(2) \quad \hat{y} = \sigma(u)$$

**BCE( $\hat{y}, y$ )**

$$(3) \quad \mathcal{L} = -y \log \hat{y} - (1 - y) \log (1 - \hat{y})$$

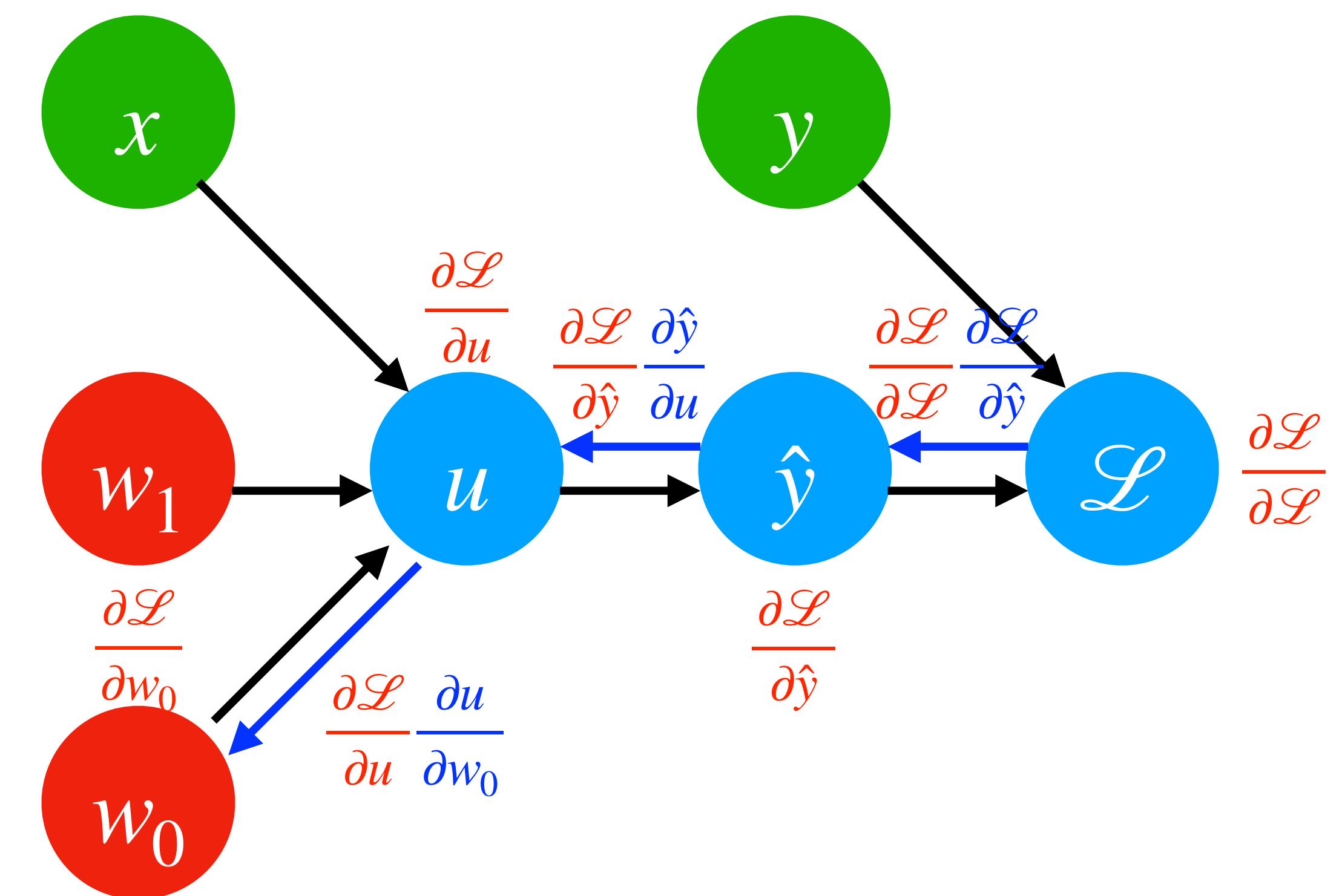
**Loss:**  $\mathcal{L} = \text{BCE}(\sigma(w_0 + w_1 x), y)$

**Backward Pass:**

$$(4) \quad \frac{\partial \mathcal{L}}{\partial \hat{y}} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial \hat{y}} = \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})}$$

$$(3) \quad \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \sigma(u)(1 - \sigma(u))$$

$$(1) \quad \frac{\partial \mathcal{L}}{\partial w_0} = \frac{\partial \mathcal{L}}{\partial u} \frac{\partial u}{\partial w_0} = \frac{\partial \mathcal{L}}{\partial u}$$



# Backpropagation: Logistic Regression with 1D Inputs

**Forward Pass:**

$$(1) \quad u = w_0 + w_1 x$$

$$(2) \quad \hat{y} = \sigma(u)$$

$$(3) \quad \mathcal{L} = -y \log \hat{y} - (1 - y) \log (1 - \hat{y})$$

**Loss:**  $\mathcal{L} = \text{BCE}(\sigma(w_0 + w_1 x), y)$

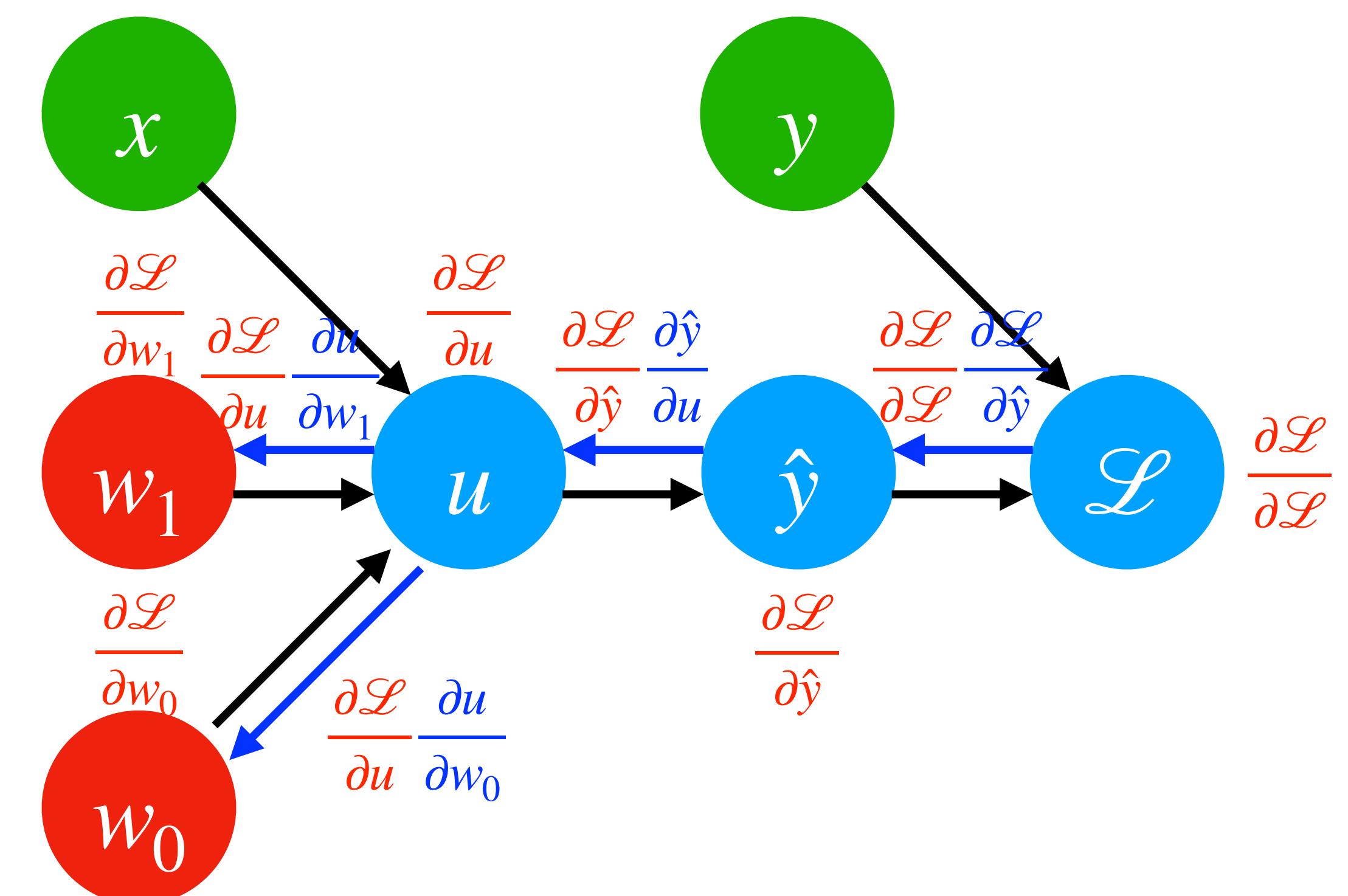
**Backward Pass:**

$$(4) \quad \frac{\partial \mathcal{L}}{\partial \hat{y}} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}} \frac{\partial \mathcal{L}}{\partial \hat{y}} = \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})}$$

$$(3) \quad \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \sigma(u)(1 - \sigma(u))$$

$$(1) \quad \frac{\partial \mathcal{L}}{\partial w_0} = \frac{\partial \mathcal{L}}{\partial u} \frac{\partial u}{\partial w_0} = \frac{\partial \mathcal{L}}{\partial u}$$

$$(1) \quad \frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial u} \frac{\partial u}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial u} x$$

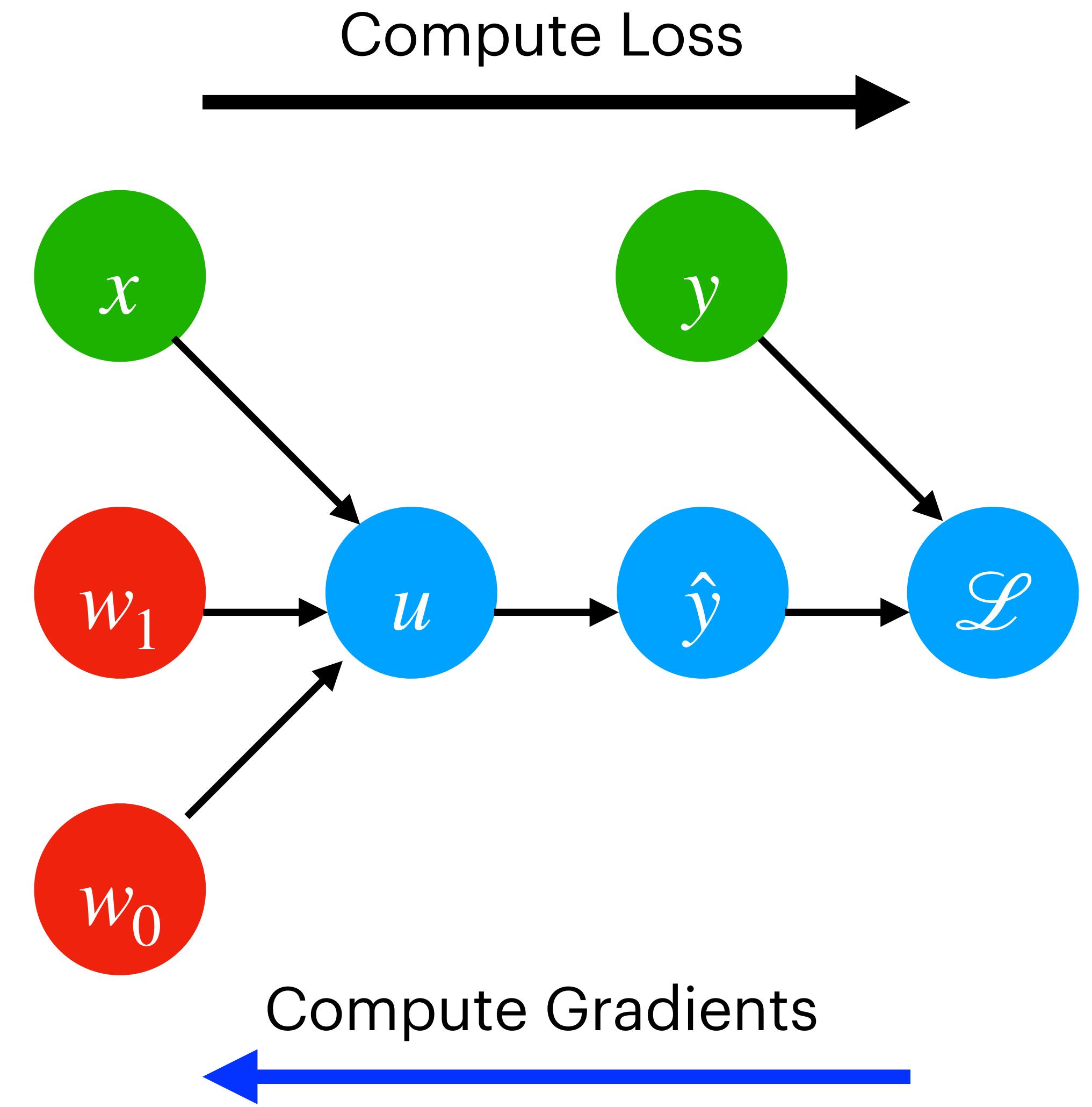


# Summary

- ◆ We can write mathematical expressions as a computation graph
- ◆ Values are efficiently computed forward, gradients backward
- ◆ Multiple incoming gradients are summed up (multivariate chain rule)
- ◆ Modularity: Each node must only “know” how to compute gradients w.r.t. its own arguments
- ◆ One fw/bw pass per data point:

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{y}, \mathbf{X}, \mathbf{w}) = \sum_{i=1}^N \boxed{\nabla_{\mathbf{w}} \mathcal{L}(y_i, \mathbf{x}_i, \mathbf{w})}$$

backprop



**Disclaimer:** So far, we discussed backpropagation only for scalar values. In the next lecture, we will discuss backpropagation with arrays and tensors.