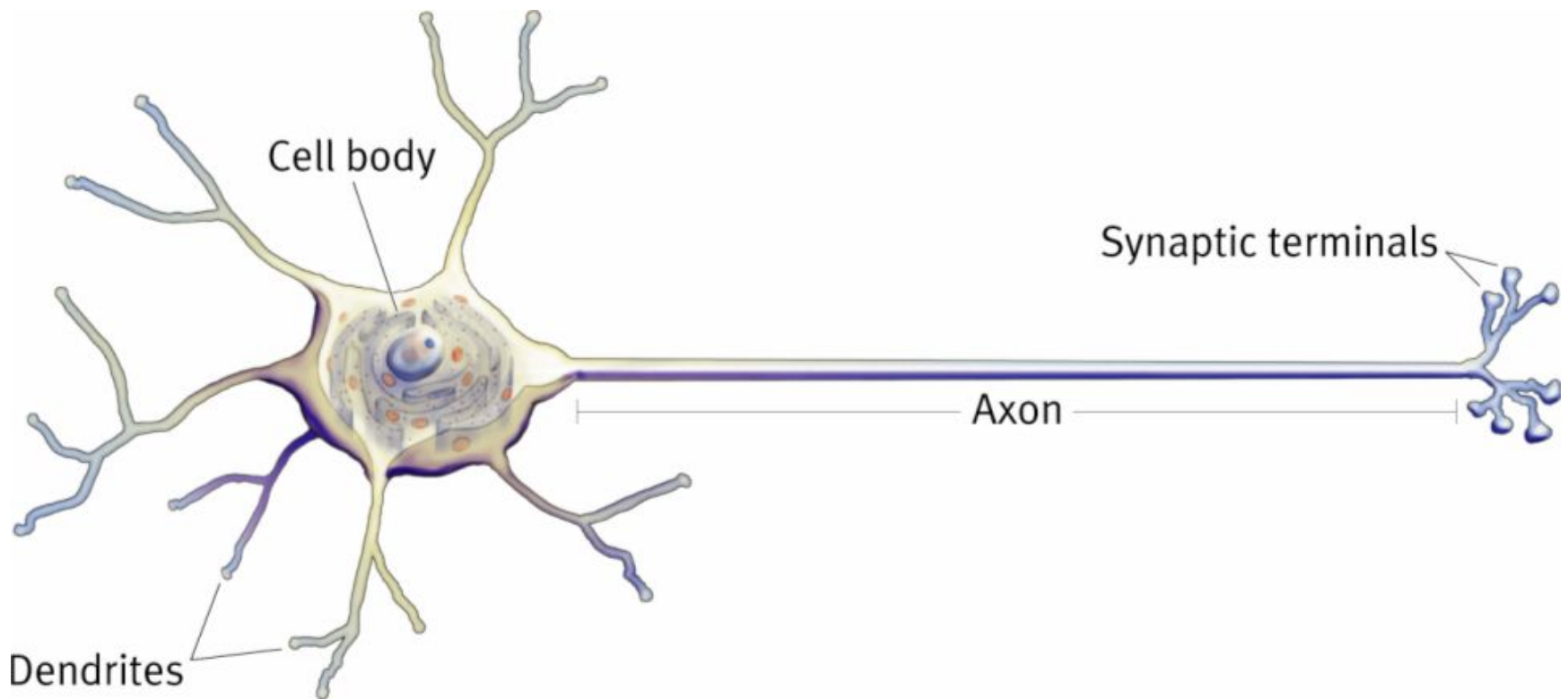


# NEURAL NETWORKS

---

# The Neuron



# The Neuron

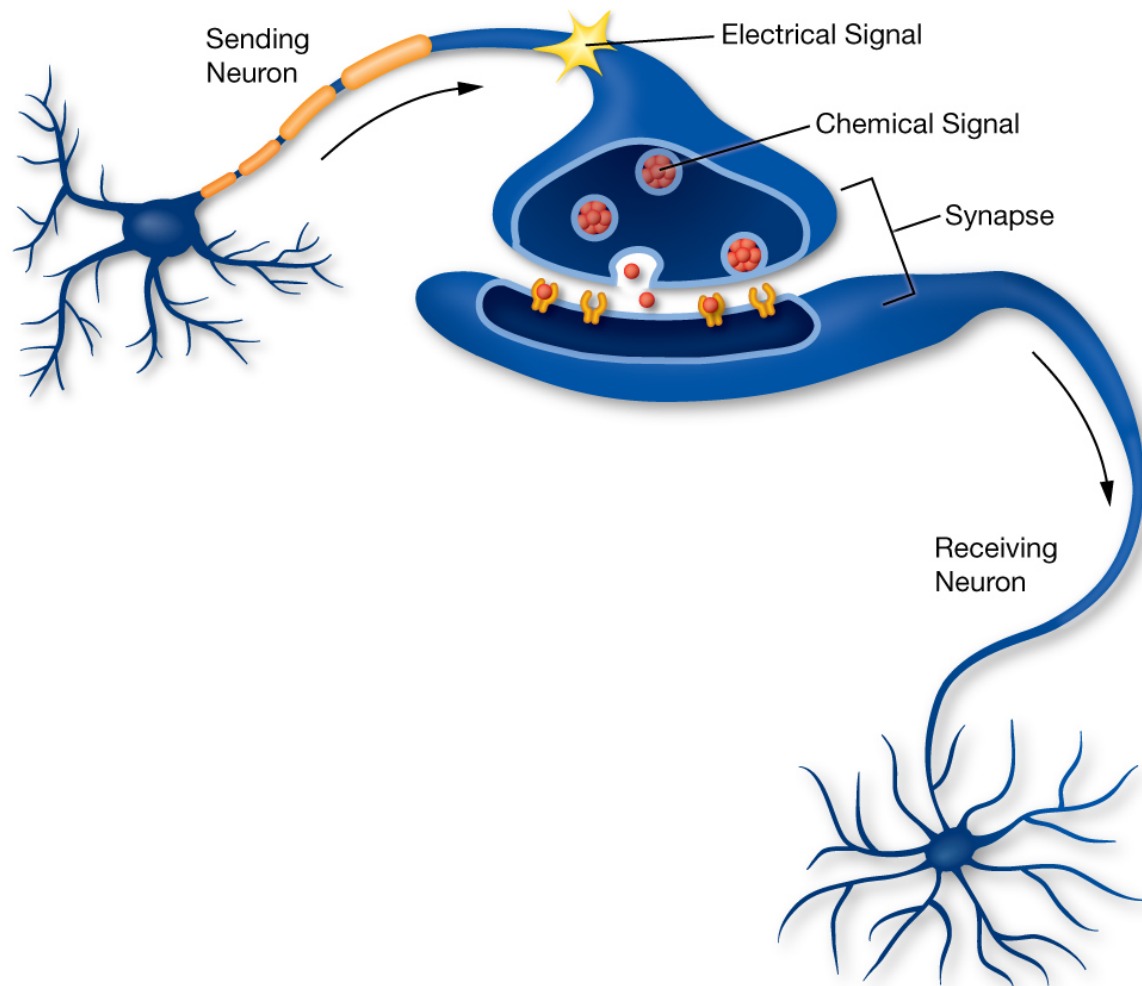
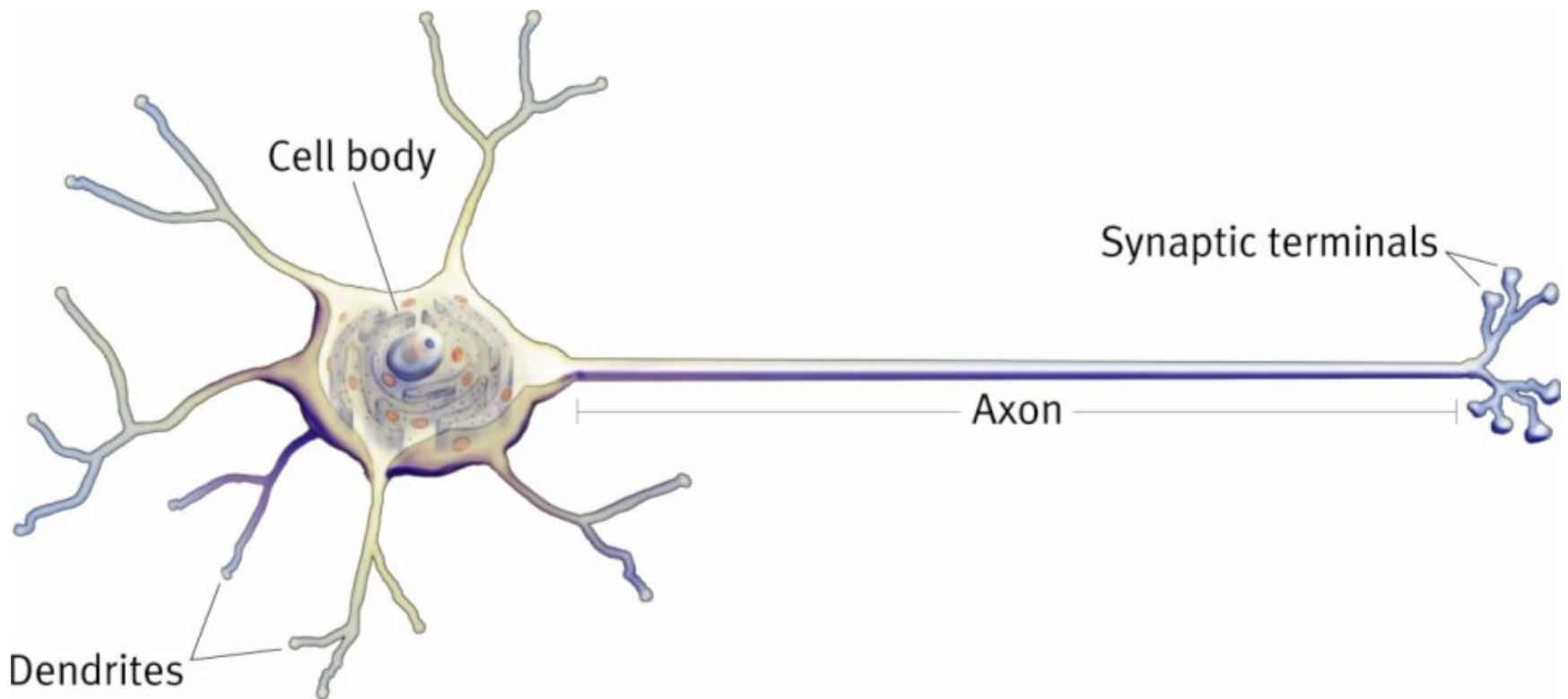


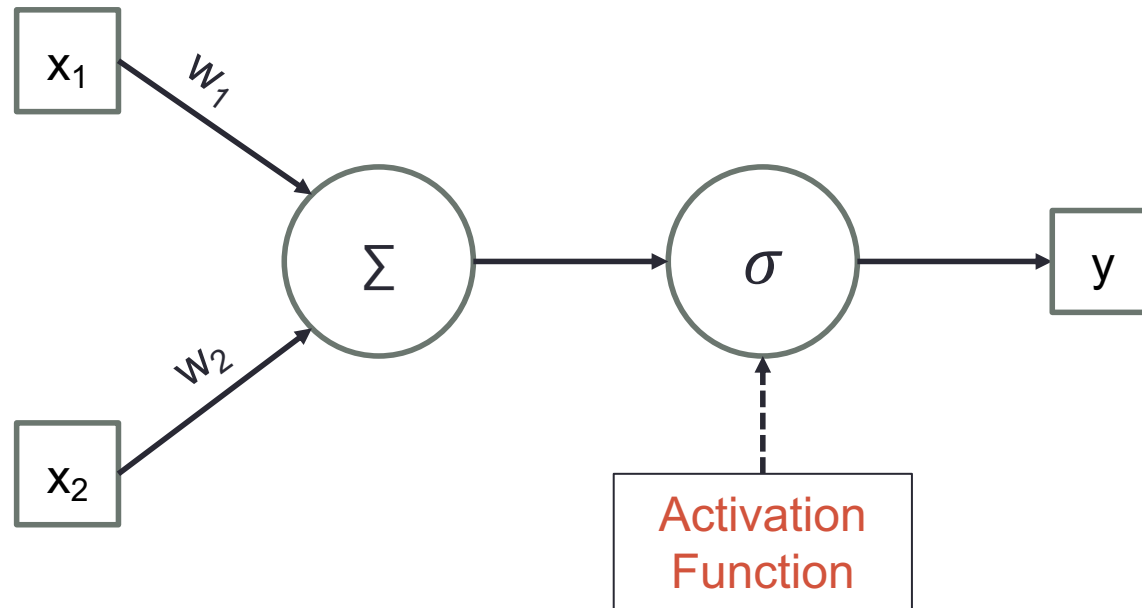
Image courtesy:

<http://learn.genetics.utah.edu/content/addiction/neurons/>

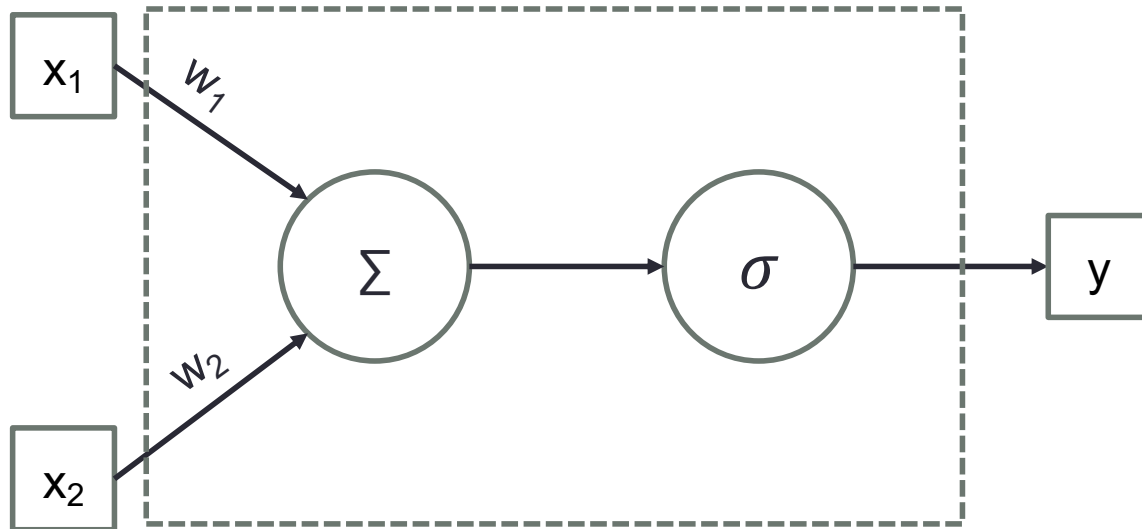
# “Real” Neuron



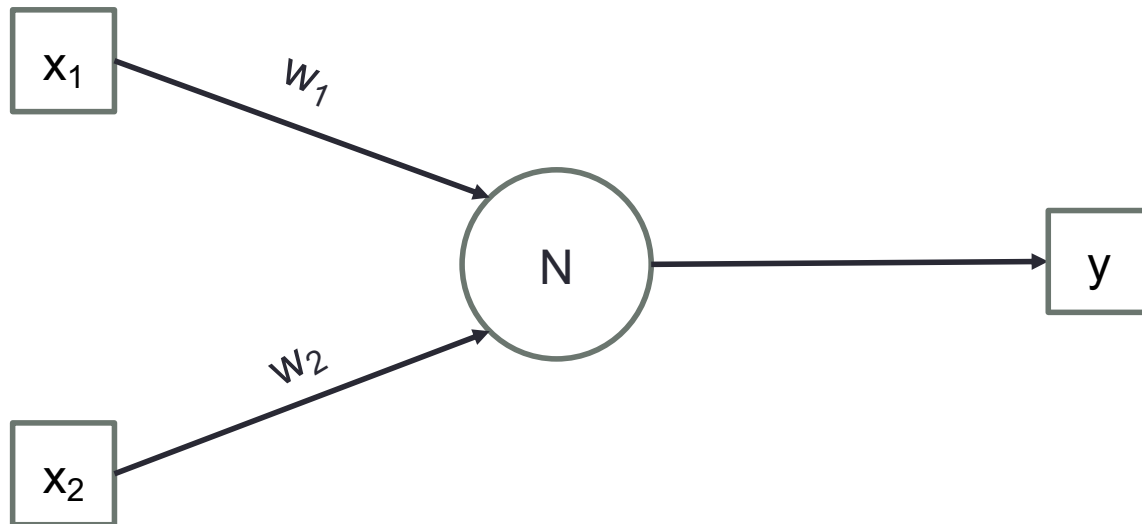
# Artificial Neurons



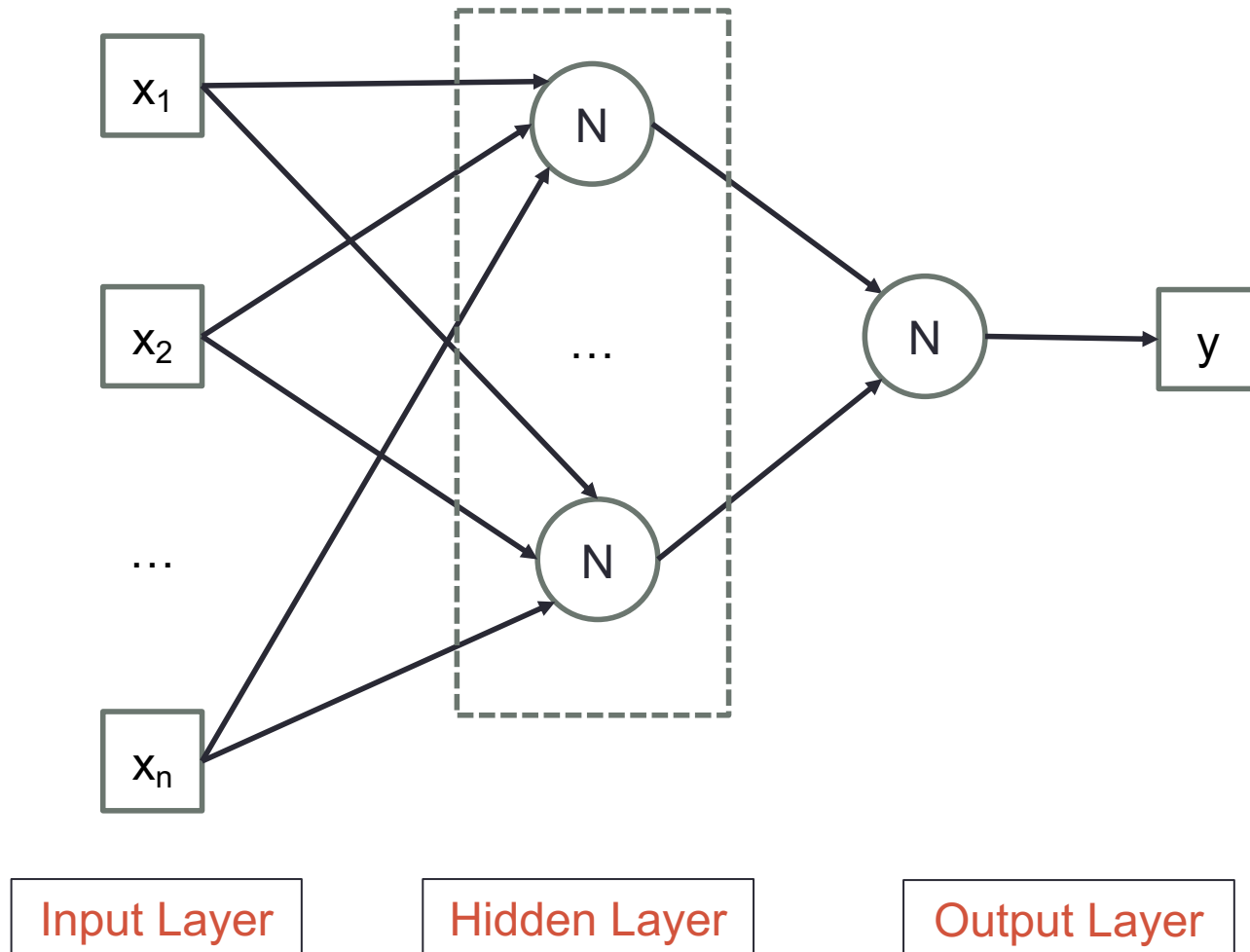
# Artificial Neurons



# Artificial Neurons



# Fully-Connected Neural Network





## BACKPROPAGATION:

Initialize all weights in the network to small, random numbers.

**loop**

**for** each training example  $(\mathbf{x}, y)$  **do**

FORWARDPROP:

For each hidden unit  $h$ ,  $a_h = \sigma(\text{net}_h) = \sigma(\sum_i w_{ih}x_i)$

$\hat{y} = a_k = \sigma(\text{net}_k) = \sigma(\sum_h w_{hk}a_h)$

BACKPROP:

$\delta_k = \frac{\partial J}{\partial \text{net}_k} = (y - \hat{y})\hat{y}(1 - \hat{y})$

For each weight  $w_h$ ,  $w_h \leftarrow w_h - \eta\delta_k a_h$

For each hidden unit  $h$ ,  $\delta_h = \delta_k w_{hk} a_h (1 - a_h)$

For each weight  $w_{ih}$ ,  $w_{ih} \leftarrow w_{ih} - \eta\delta_h x_i$

**end for**

**end loop**

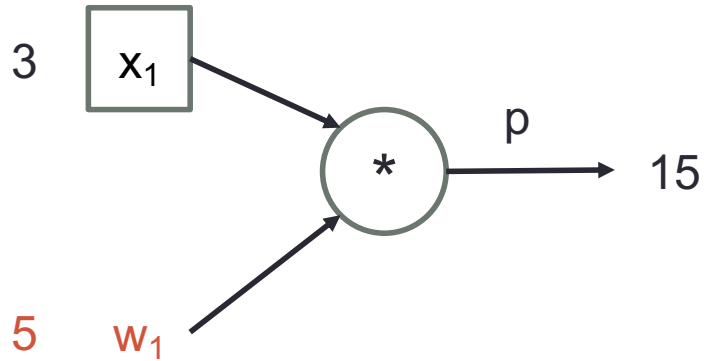
The diagram illustrates the architecture of the proposed model, which is a deep learning network for multi-class classification. The architecture is divided into two main branches, each ending with a Softmax layer. The top branch includes a 'Deficient' layer and a 'Softmax' layer. The bottom branch includes a 'Softmax' layer and a 'Softmax' layer. The architecture is designed to handle multi-class classification tasks.

# Automatic Differentiation

- Use the abstraction of a **computational graph**
- Define your computation and let engine worry about optimization



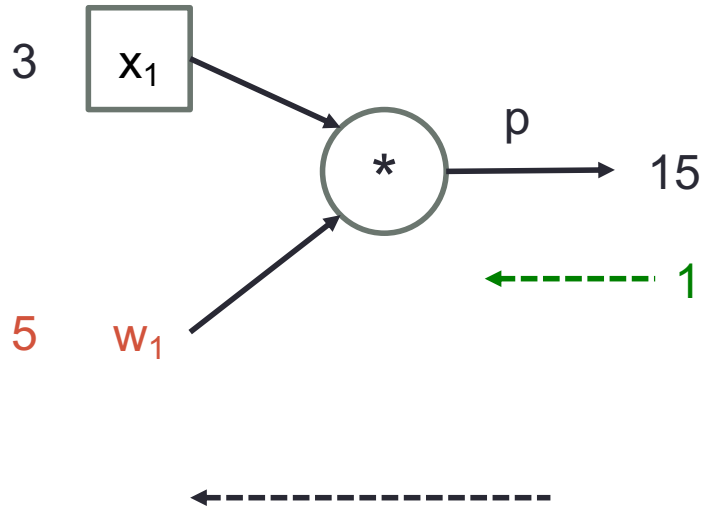
# Computational Graph



## Forward Pass

- Apply the operator

# Computational Graph

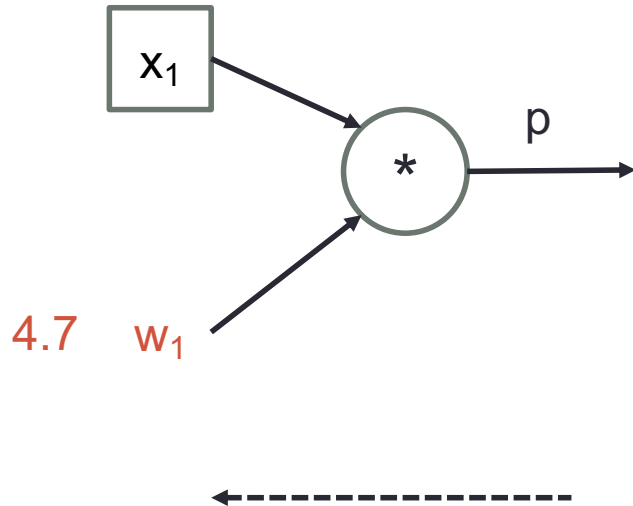


$$\frac{\partial p}{\partial w_1} = x_1$$

## Backward Pass

- Adjust parameter using local gradient 3 (scaled by a learning rate)

# Computational Graph

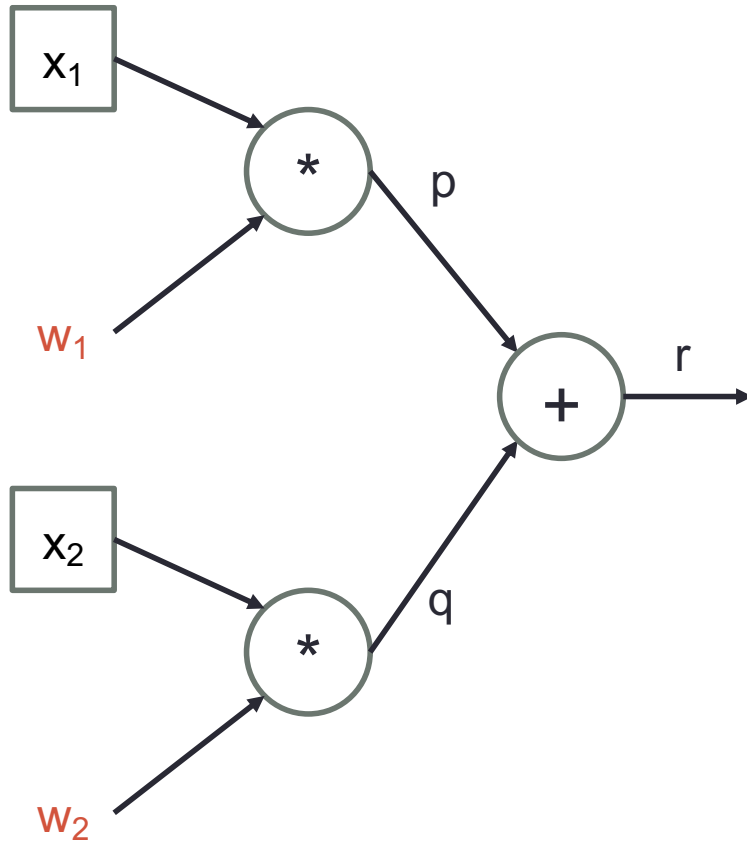


$$\frac{\partial p}{\partial w_1} = x_1$$

## Backward Pass

- Adjust parameter using local gradient 3 (scaled by a learning rate)

# Computational Graph

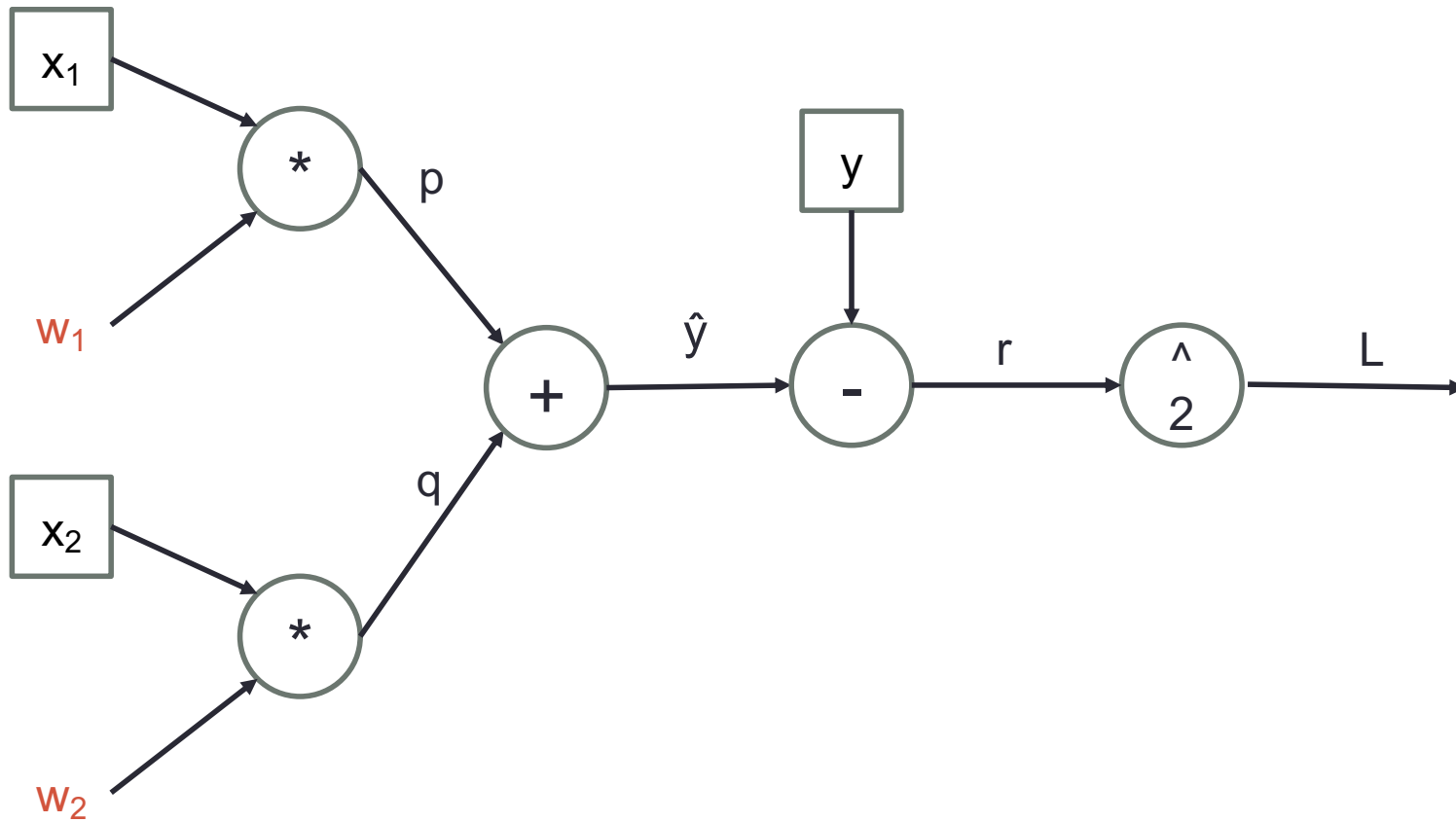


$$\frac{\partial r}{\partial p} = 1$$

$$\frac{\partial p}{\partial w_1} = x_1$$

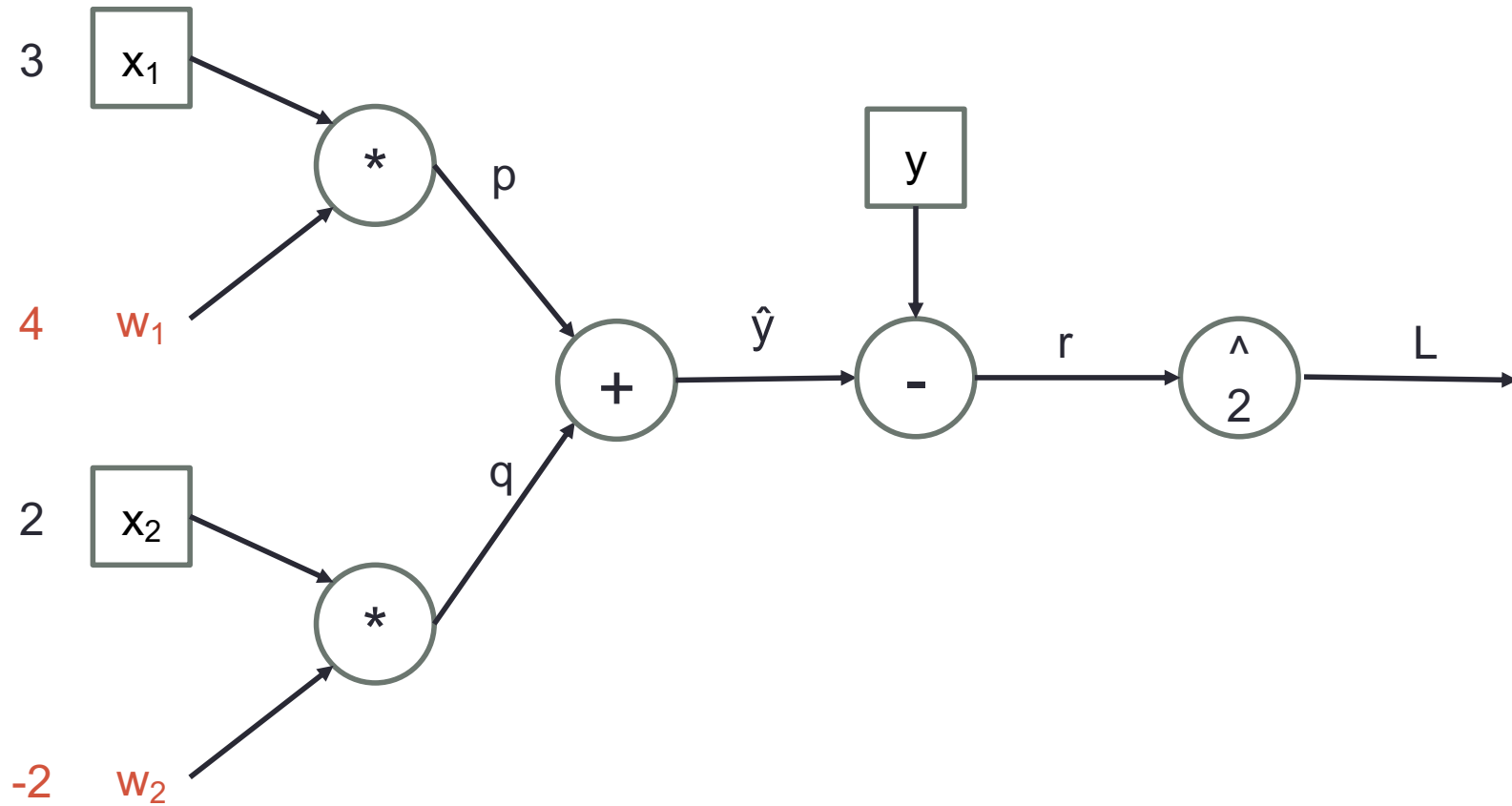
$$\frac{\partial r}{\partial w_1} = \frac{\partial r}{\partial p} \frac{\partial p}{\partial w_1}$$

# Computational Graph

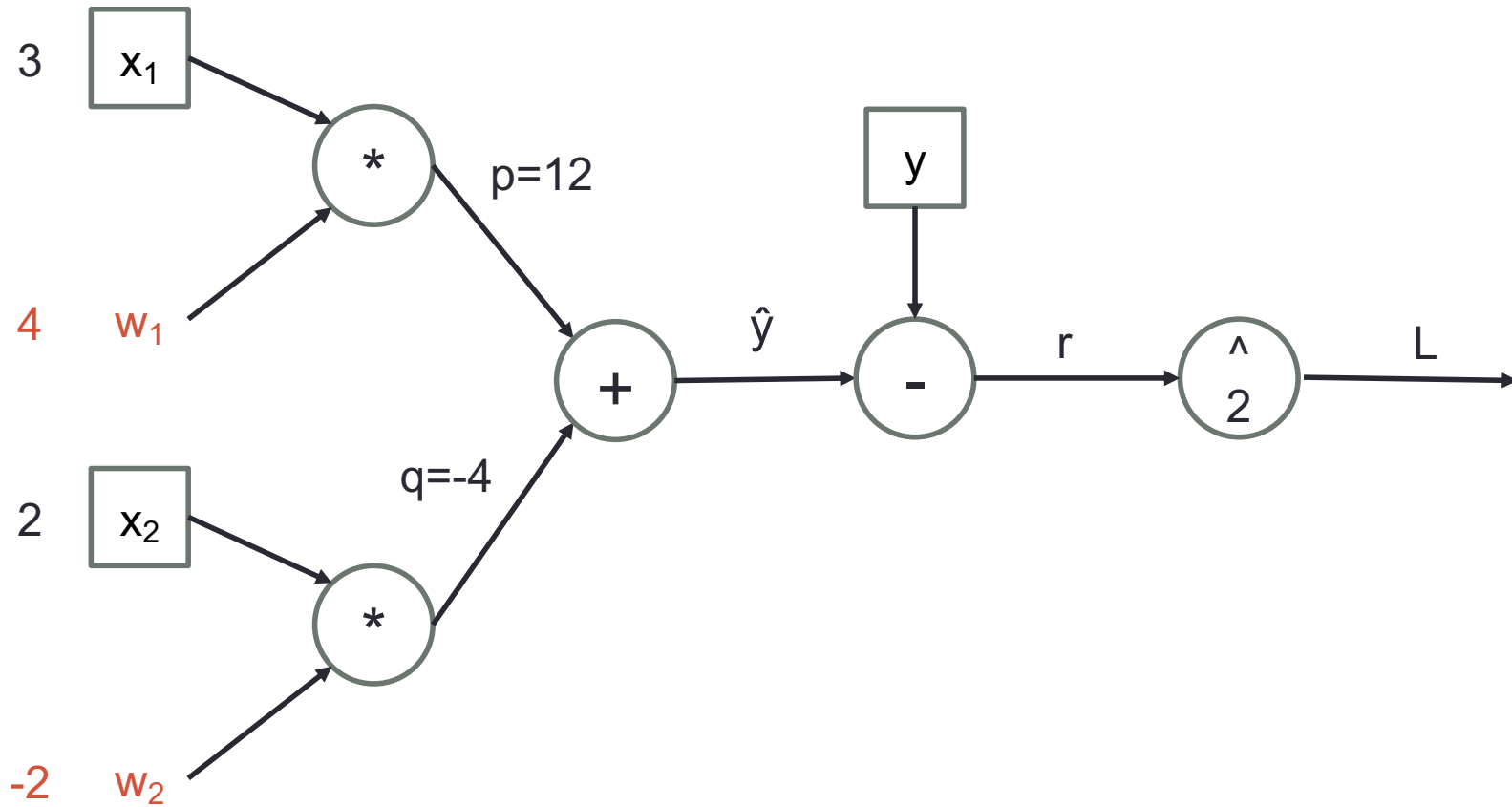




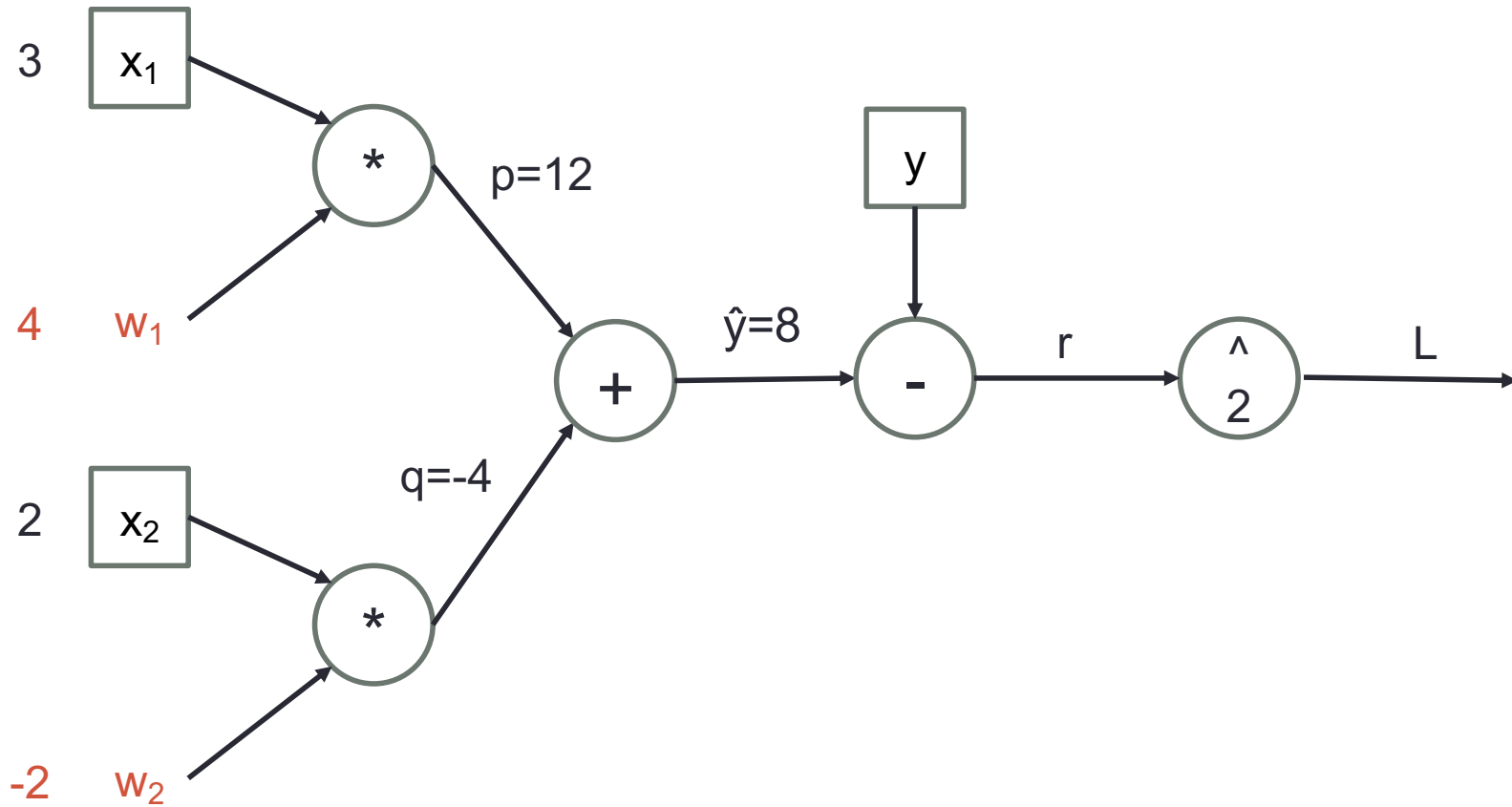
# Computational Graph



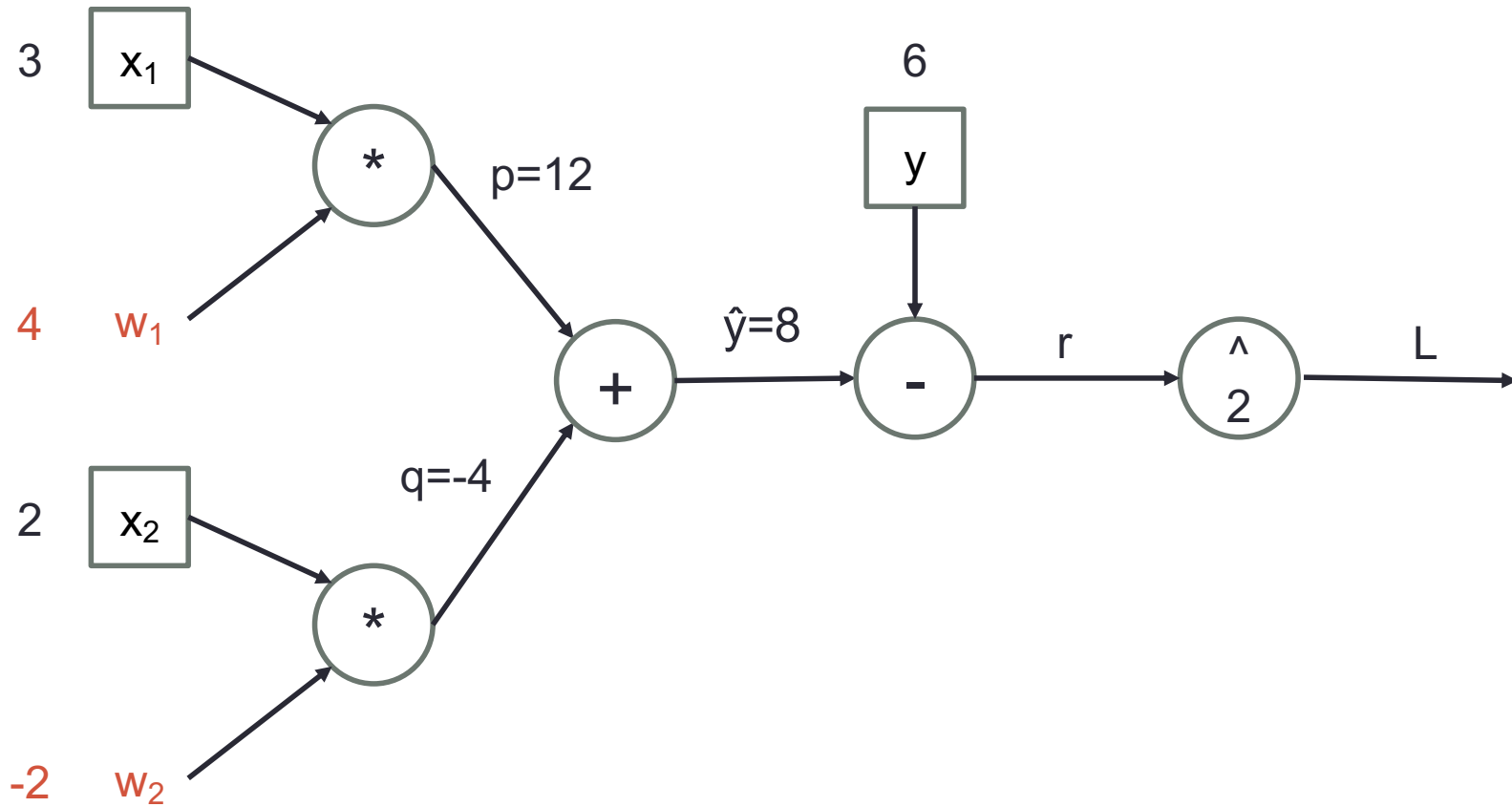
# Computational Graph



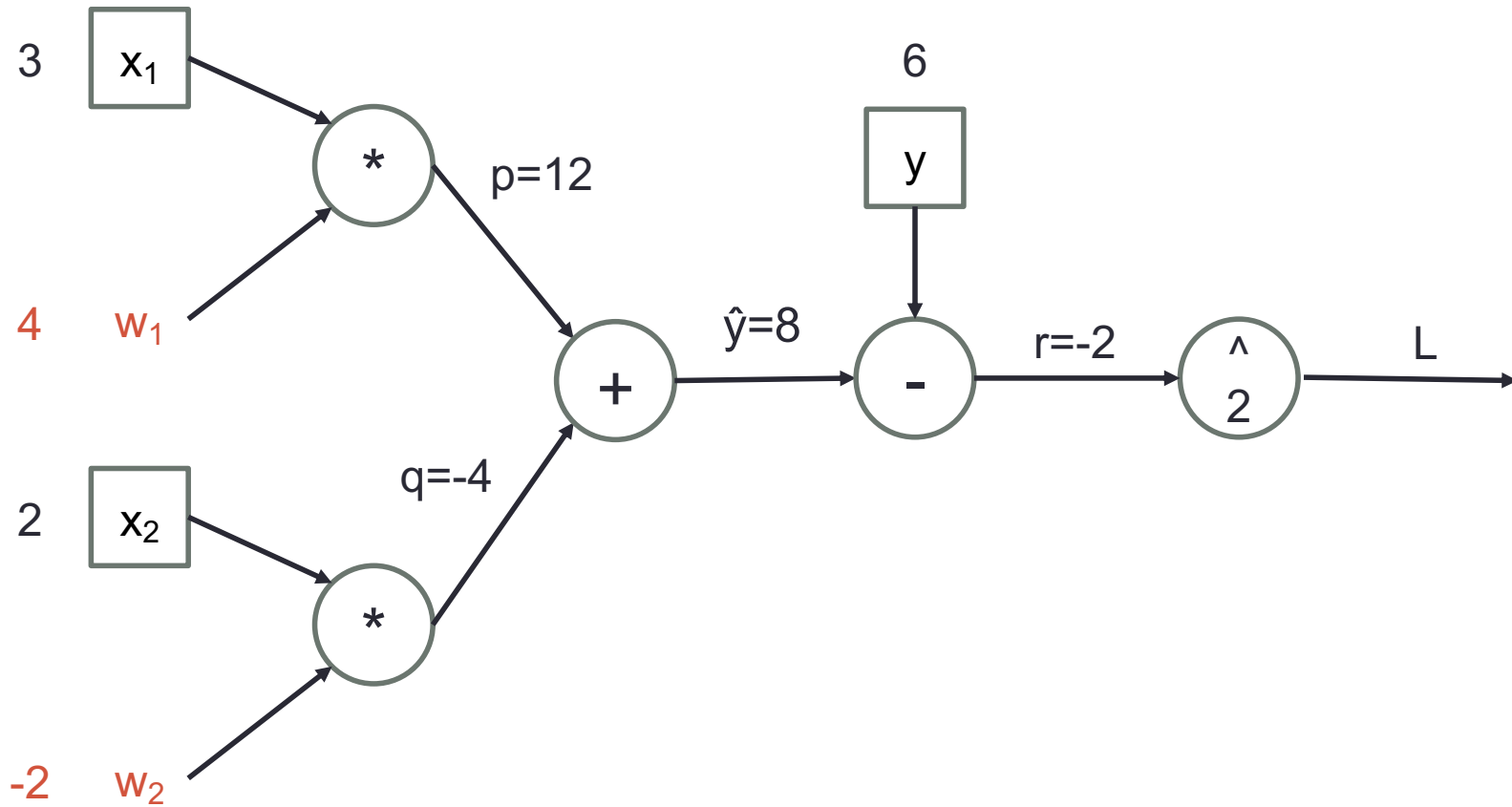
# Computational Graph



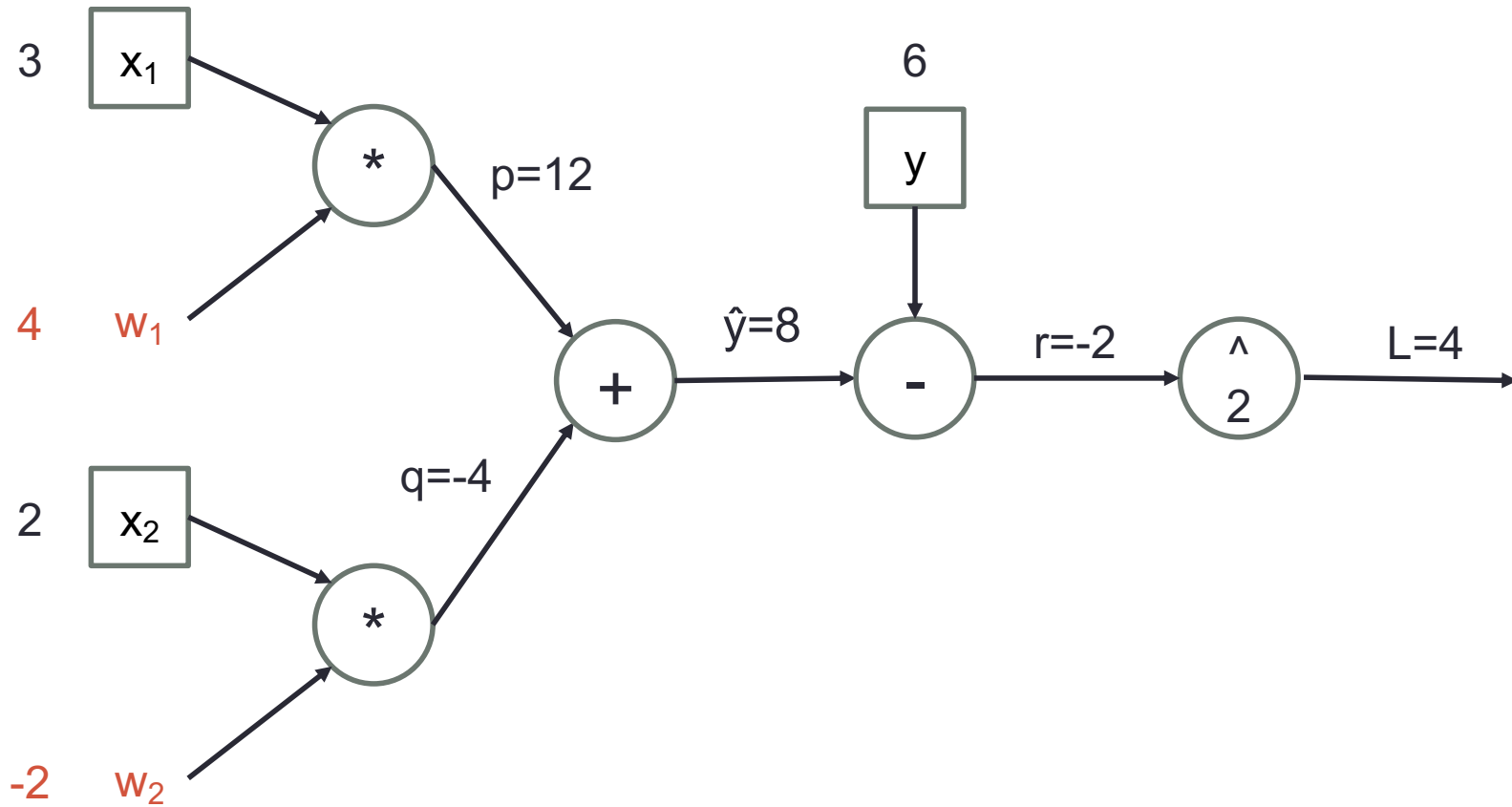
# Computational Graph



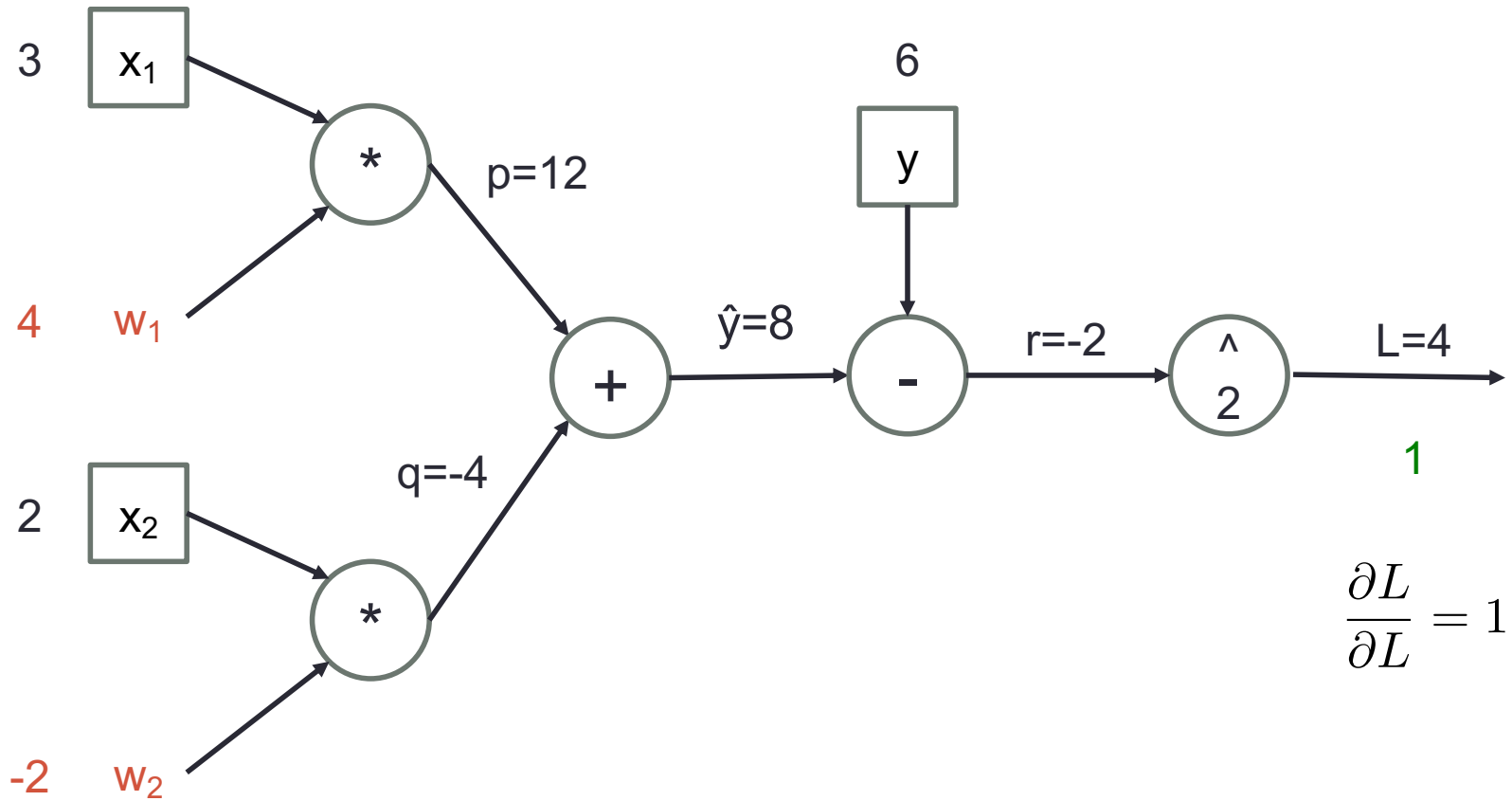
# Computational Graph



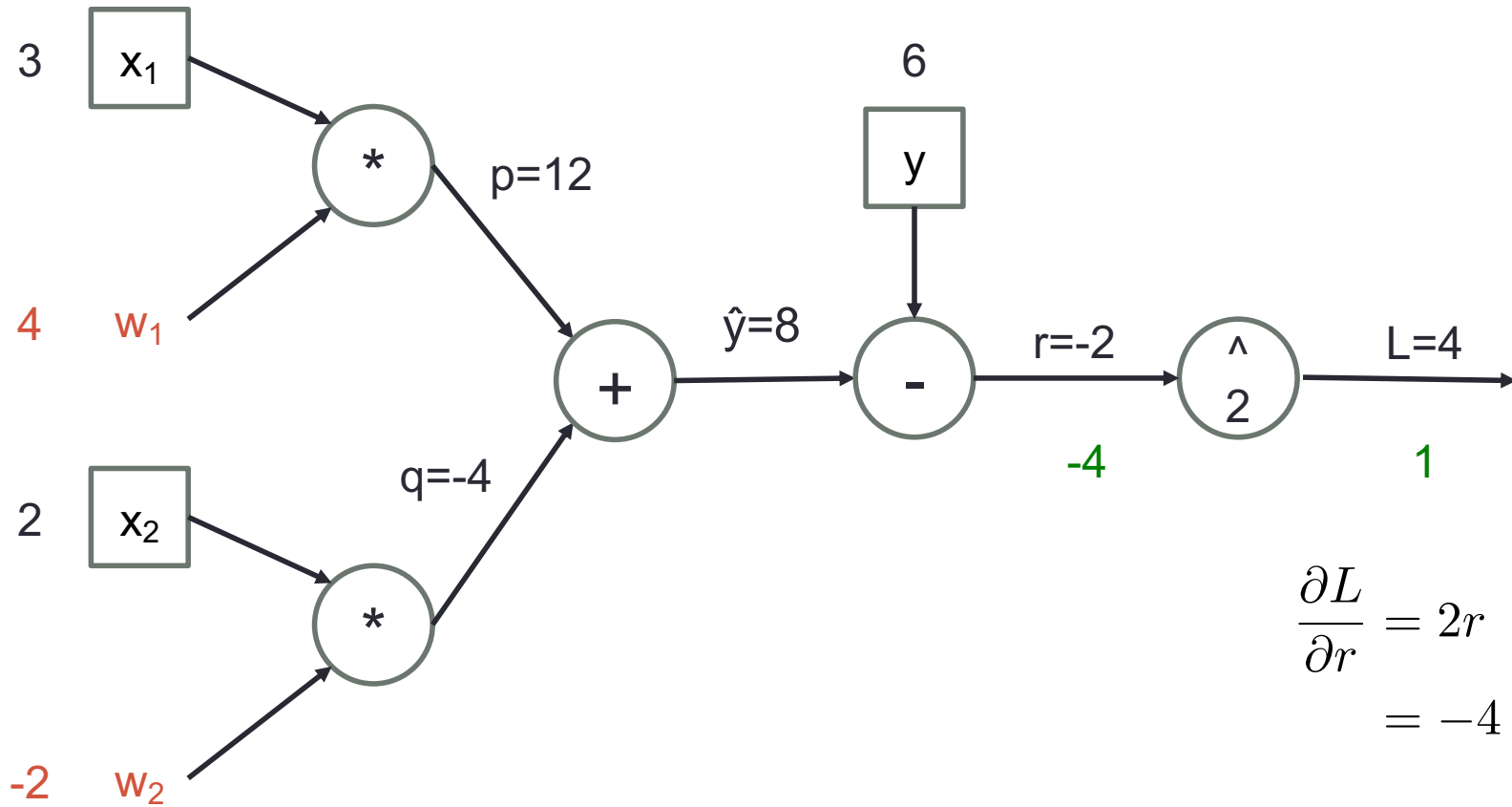
# Computational Graph



# Computational Graph

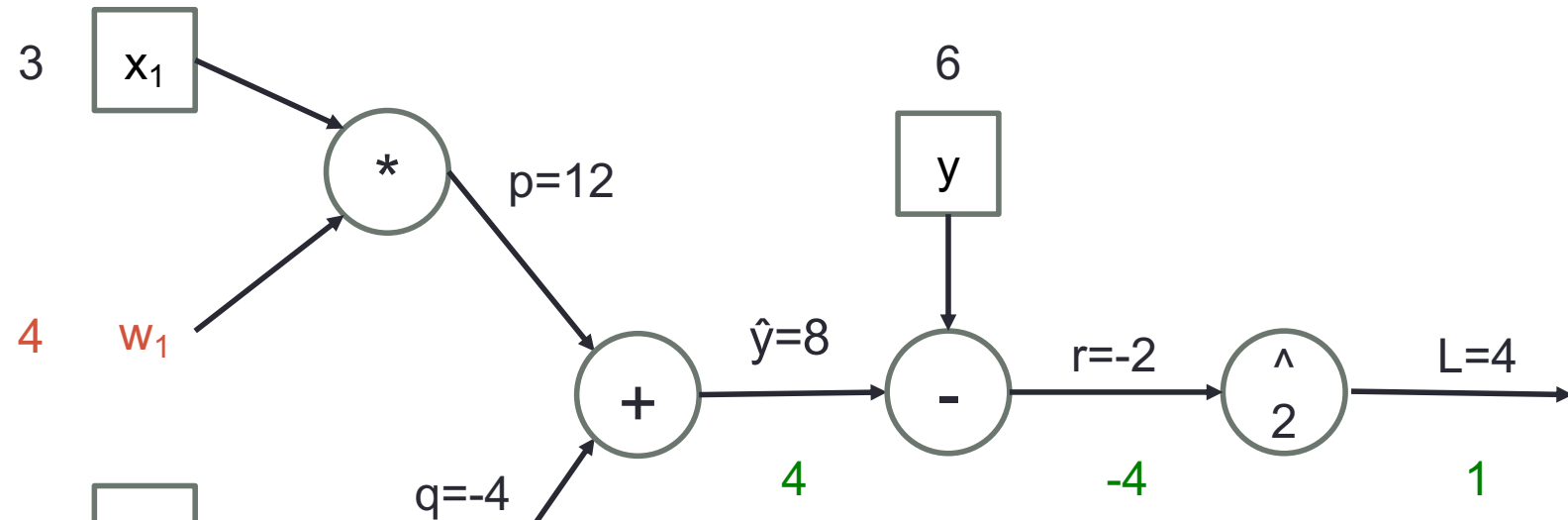


# Computational Graph



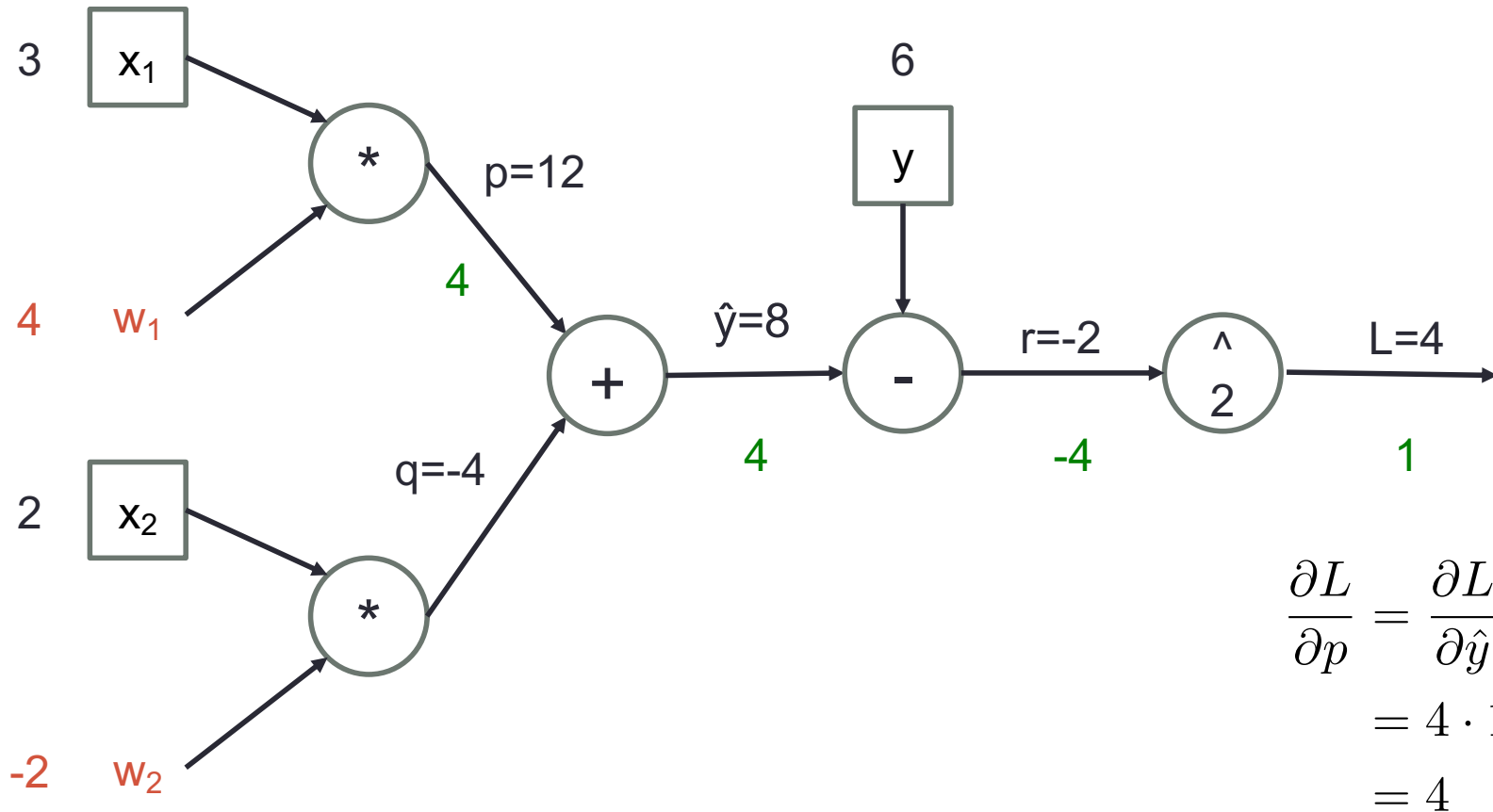


# Computational Graph

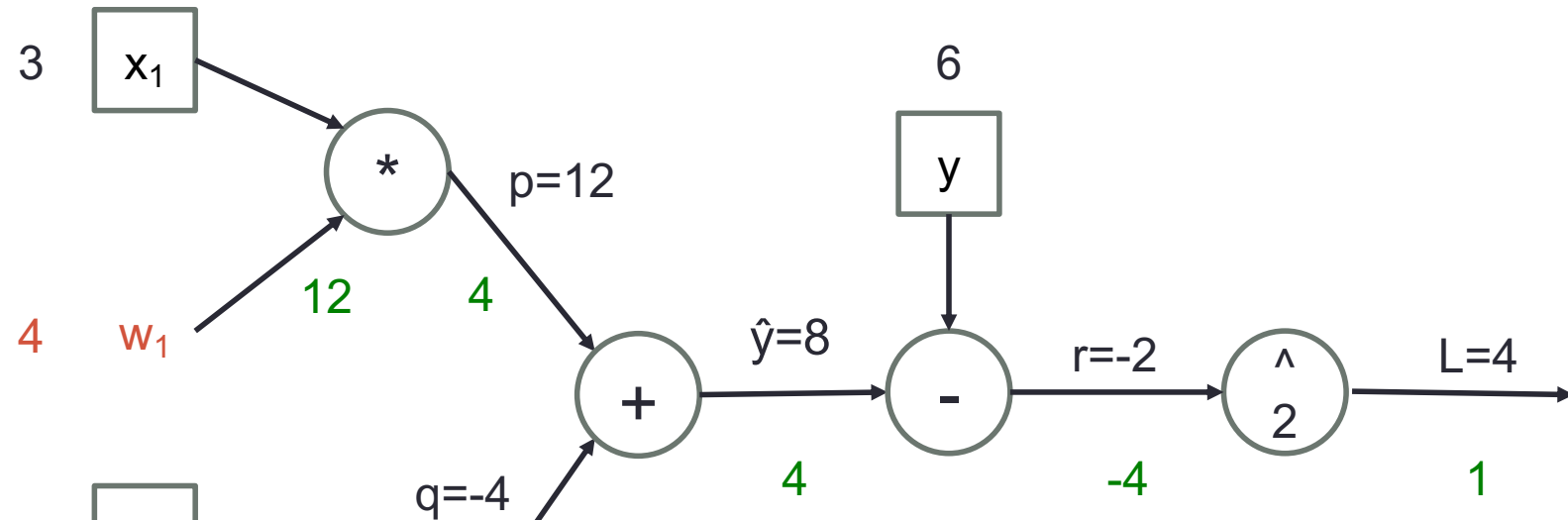


$$\begin{aligned}\frac{\partial L}{\partial \hat{y}} &= \frac{\partial L}{\partial r} \frac{\partial r}{\partial \hat{y}} \\ &= -4 \cdot -1 \\ &= 4\end{aligned}$$

# Computational Graph

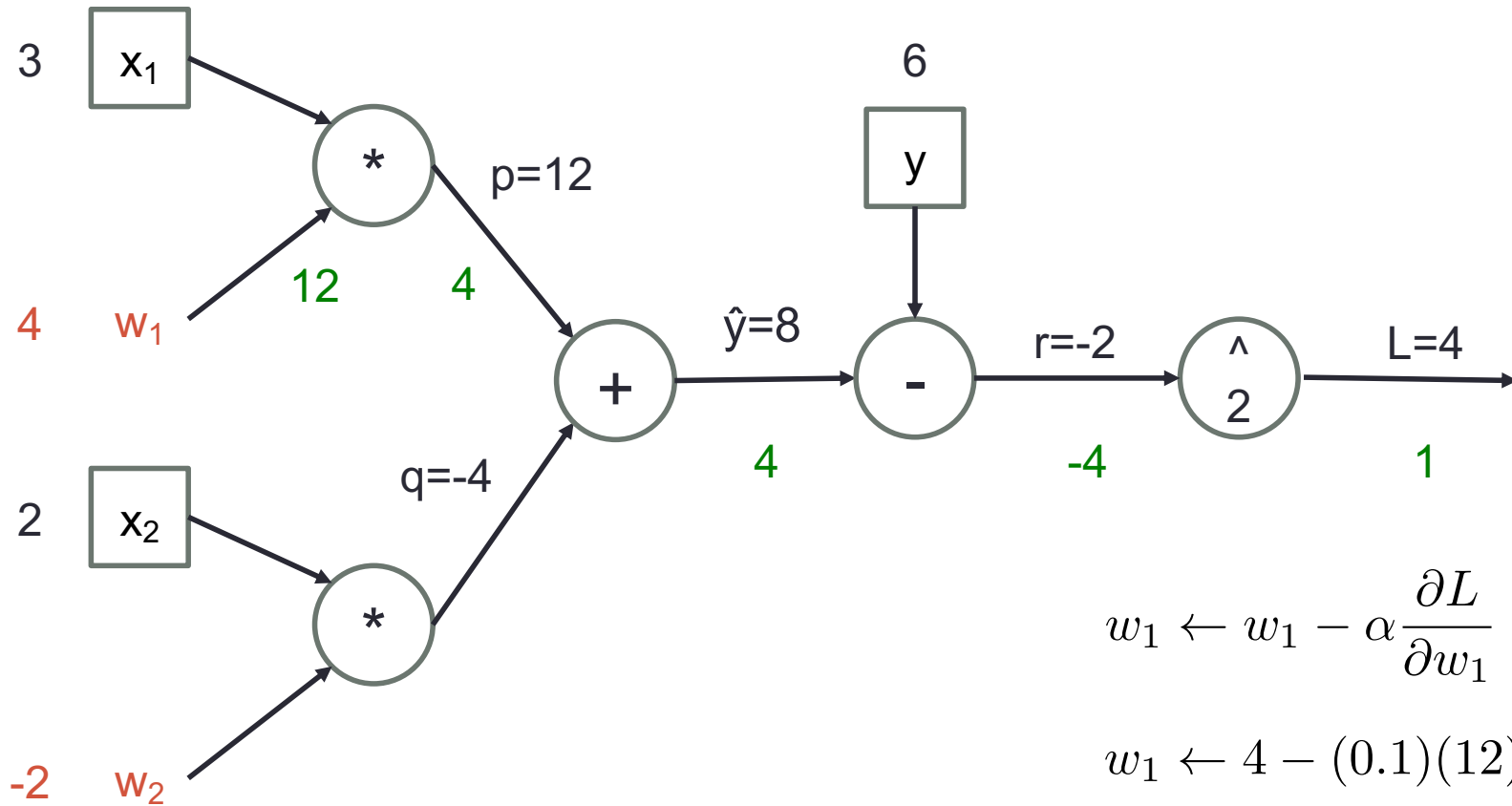


# Computational Graph

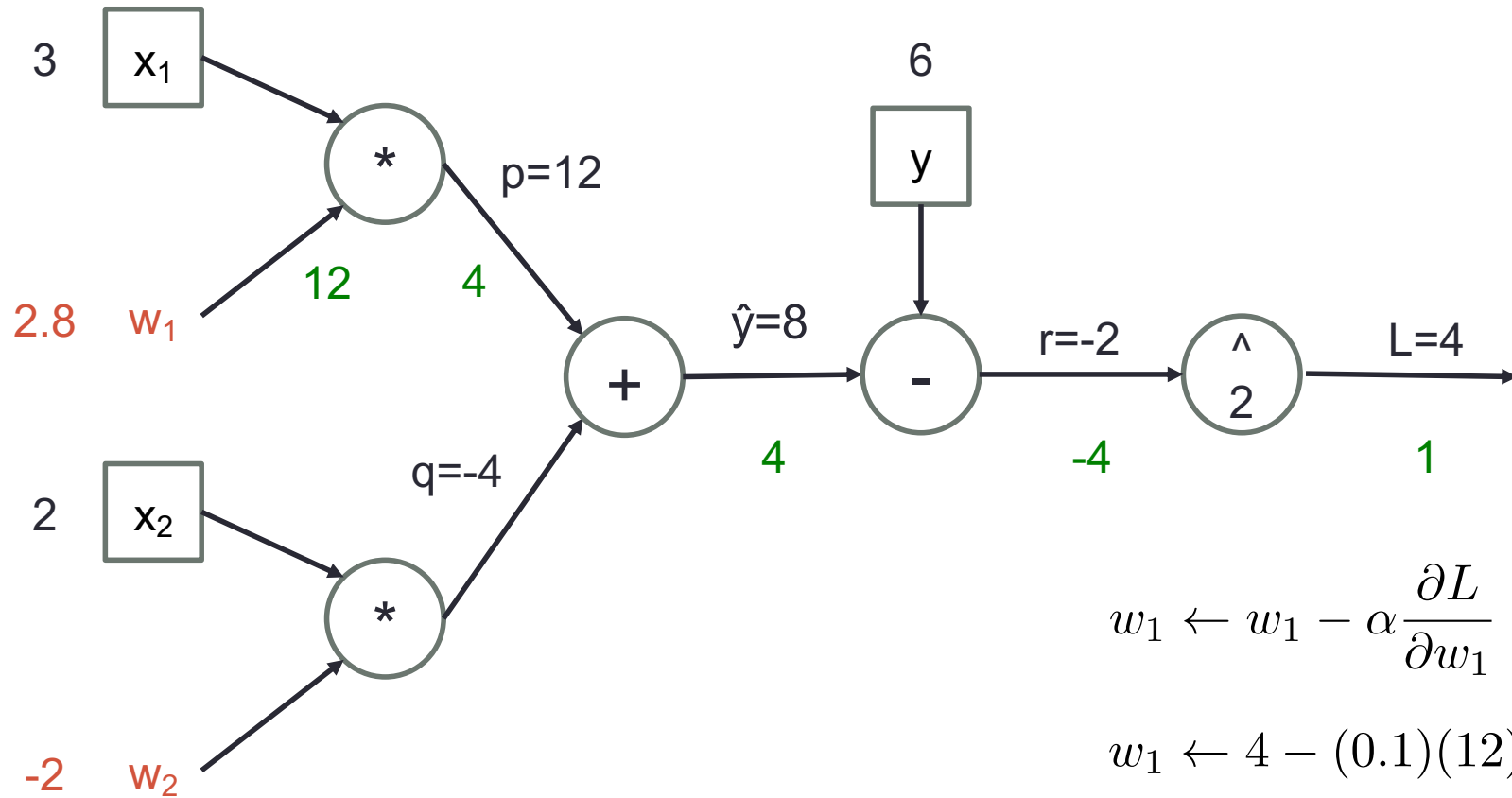


$$\begin{aligned}\frac{\partial L}{\partial w_1} &= \frac{\partial L}{\partial p} \frac{\partial p}{\partial w_1} \\ &= 4x_1 \\ &= 12\end{aligned}$$

# Computational Graph



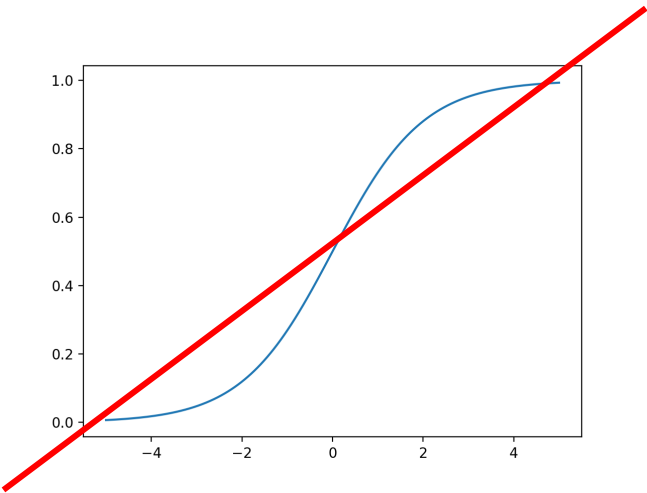
# Computational Graph



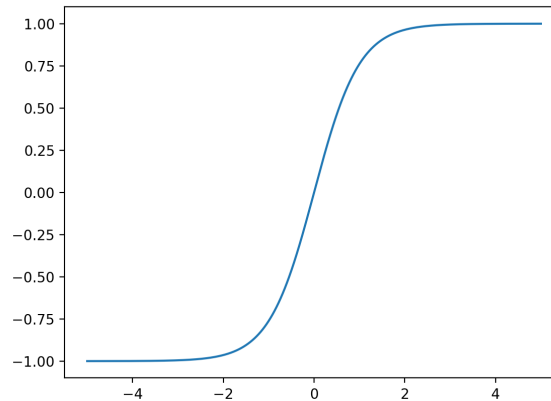
# Preparing Your Data

- Shuffle your data
- Mean center your data
  - Why?
- Normalize the variance
  - Why?
- “Whitening”
  - Decorrelates data
  - Can be hit or miss
- When to do train/test split?

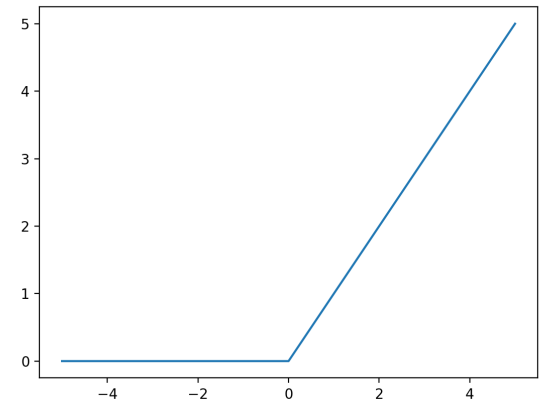
# Choosing an Activation Function



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



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



$$f(x) = \max(0, x)$$

# Initializing Your Network

- Set all weights to 0?
  - Terrible idea
- Set all weights to random values?
  - *Small* random values
- State-of-the-art: Xavier or Glorot initialization
  - Takes into account fan-in/fan-out of a neuron when initializing its weights



# Optimization Methods

- Stochastic gradient descent

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla L(\mathbf{w})$$

- Stochastic gradient descent + momentum

$$\mathbf{z} \leftarrow \beta \mathbf{z} + \nabla L(\mathbf{w})$$

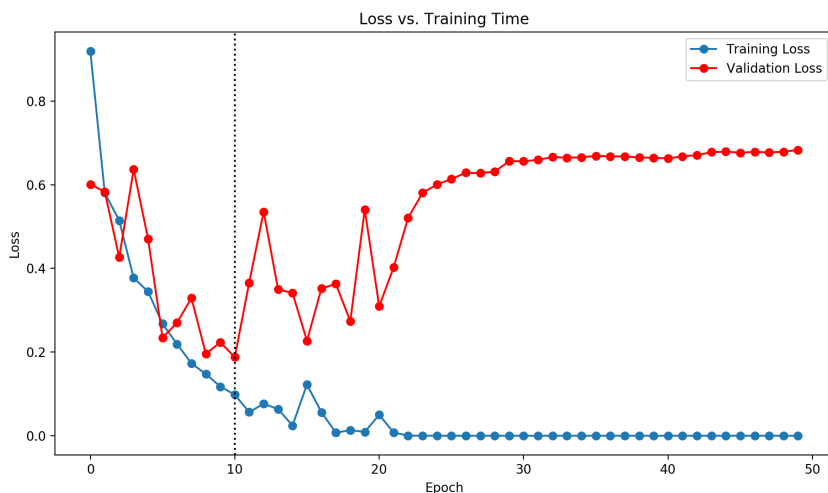
$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \mathbf{z}$$

- State-of-the-art approaches:
  - RMSProp
  - Adam

# Regularization

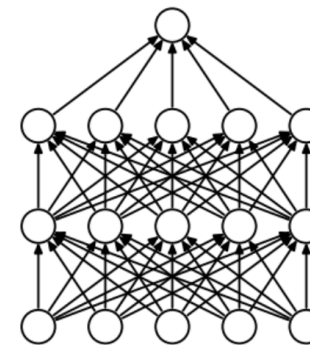
## Classical Approaches

- Weight decay
  - L2 term
- Early stopping

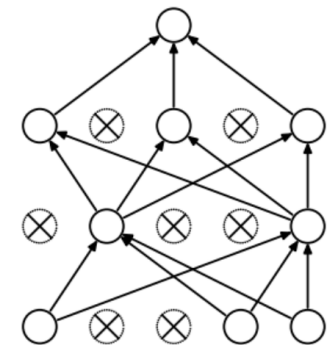


## Modern Approaches

- Dropout



(a) Standard Neural Net



(b) After applying dropout.

- Batch normalization