

# IBM Research Data Science Course

srihari sridharan

# Agenda

## More Neural Network

- Loss
- Gradient Descent
- Backpropagation
- Computational Graphs (Tensorflow 101)

# Training Process (Backpropagation)

Step 1 : Initially weights and bias are randomly assigned.

Step 2 : Vector -> Matrix Multiplication -> Some Output (Vector)

Step 3 : Compare with the actual value (Label)

Step 4 : See the difference and pass it back to all the layers behind it.

Step 5 : Update the weights and repeat from Step 2

# Training Process (Finding the values of W,b )

Step 1 : Initialise the weights to something

(4 x 3)


(1 x 4)

4 . 9	3 . 0	1 . 4	0 . 2
(0,1)	(0,2)	(0,3)	(0,4)

**x**

1	1	1
1	1	1
1	1	1
1	1	1
1	1	1

# Training Process (Finding the values of W,b )

## Step 2 : Do the matrix multiplication with initialised weights.

(4 x 3)

(1 x 4)

4 . 9	3 . 0	1 . 4	0 . 2
(0,1)	(0,2)	(0,3)	(0,4)

x

1	1	1
1	1	1
1	1	1
1	1	1

2

**9.5**

# (Prediction)



# Training Process (Finding the values of W,b )

Step 3 : Compare with actual value (Label)

$$\begin{matrix} & \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix} \\ \boxed{4.9 \quad 3.0 \quad 1.4 \quad 0.2} \quad \mathbf{x} & = \quad \boxed{9.5 \quad 9.5 \quad 9.5} \quad (\text{Prediction}) \end{matrix}$$

↓

Label : 1 | 0 | 0

## Training Process (Finding the values of W,b )

Step 4 : See the difference and pass back the changes (gradients) to previous layers

$$\begin{matrix} & \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix} \\ \boxed{4.9 \quad 3.0 \quad 1.4 \quad 0.2} \quad \mathbf{x} & = \quad \boxed{9.5 \quad 9.5 \quad 9.5} \end{matrix}$$

**Square Error : Square of difference.**

$$E_{total} = \sum \frac{1}{2}(target - output)^2 = 84.25$$

Label : 

1	0	0
---	---	---

# Training Process (Finding the values of W,b )

Step 4 : See the difference and pass back the changes (gradients) to previous layers

$$\begin{matrix} & \begin{matrix} 1 & 1 & 1 \end{matrix} \\ \begin{matrix} 4.9 & 3.0 & 1.4 & 0.2 \end{matrix} & \times & \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix} & = & \begin{matrix} 9.5 & 9.5 & 9.5 \end{matrix} \end{matrix}$$

Error = 84.25

Label : 

1	0	0
---	---	---

$$E_{total} = \sum \frac{1}{2}(target - output)^2$$

$$\frac{\partial E_{total}}{\partial out_{o1}} = 2 * \frac{1}{2}(target_{o1} - out_{o1})^{2-1} * -1 + 0$$

$$\frac{\partial E_{total}}{\partial out_{o1}} = -(target_{o1} - out_{o1})$$



Gradients (Derivatives)

# Training Process (Finding the values of W,b )

## Step 5 : Update the Weights

$$\begin{matrix} 4.9 & 3.0 & 1.4 & 0.2 \end{matrix} \boxed{\mathbf{x}} \begin{matrix} 1-\delta & 1-\delta & 1-\delta \\ 1-\delta & 1-\delta & 1-\delta \\ 1-\delta & 1-\delta & 1-\delta \\ 1-\delta & 1-\delta & 1-\delta \end{matrix} = \begin{matrix} 9.5 & 9.5 & 9.5 \end{matrix}$$

Error = 84.25

Label :  $\begin{matrix} 1 & 0 & 0 \end{matrix}$

$$E_{total} = \sum \frac{1}{2}(target - output)^2$$

$$\frac{\partial E_{total}}{\partial out_{o1}} = 2 * \frac{1}{2}(target_{o1} - out_{o1})^{2-1} * -1 + 0$$

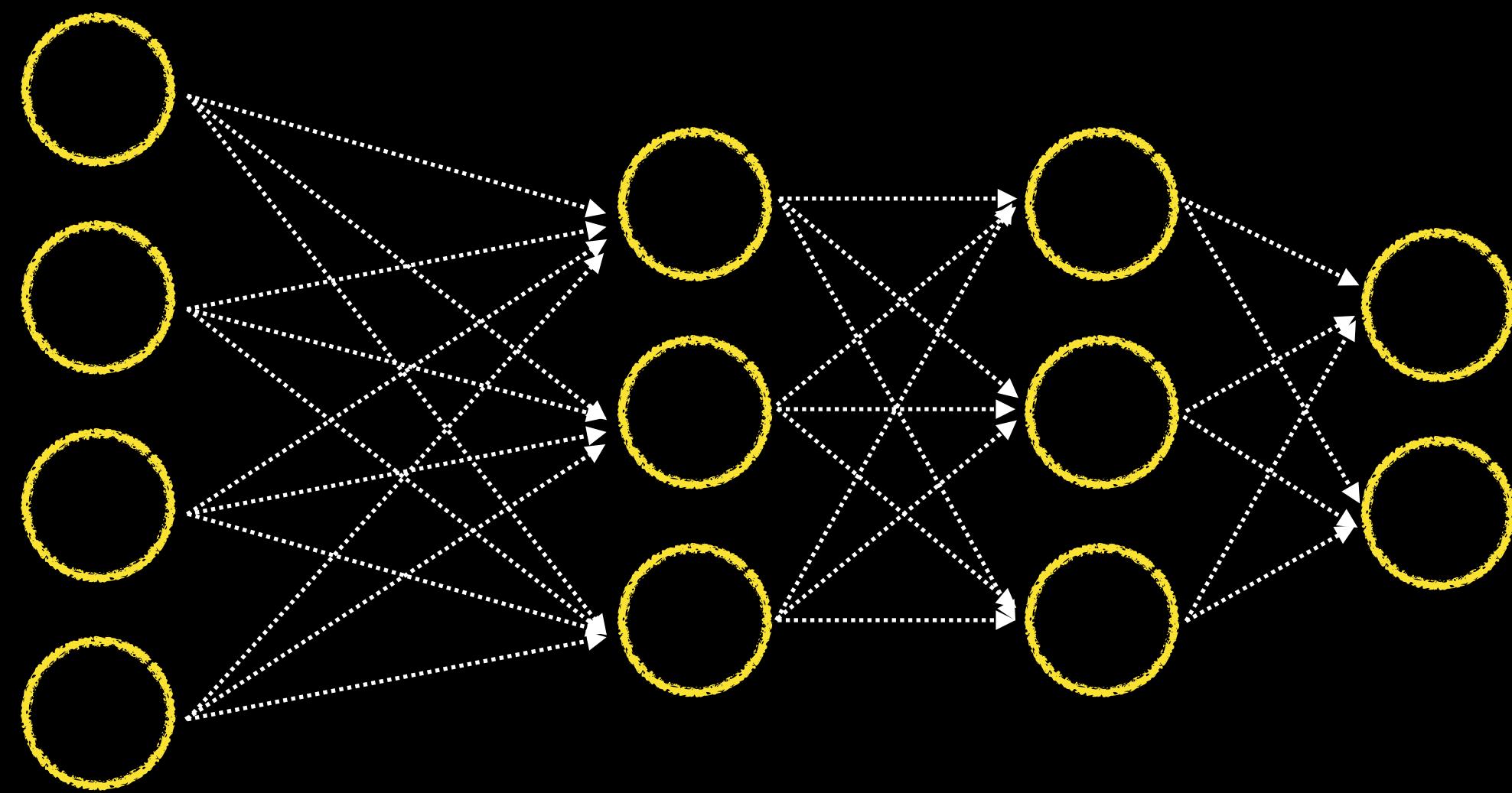
$$\frac{\partial E_{total}}{\partial out_{o1}} = -(target_{o1} - out_{o1})$$



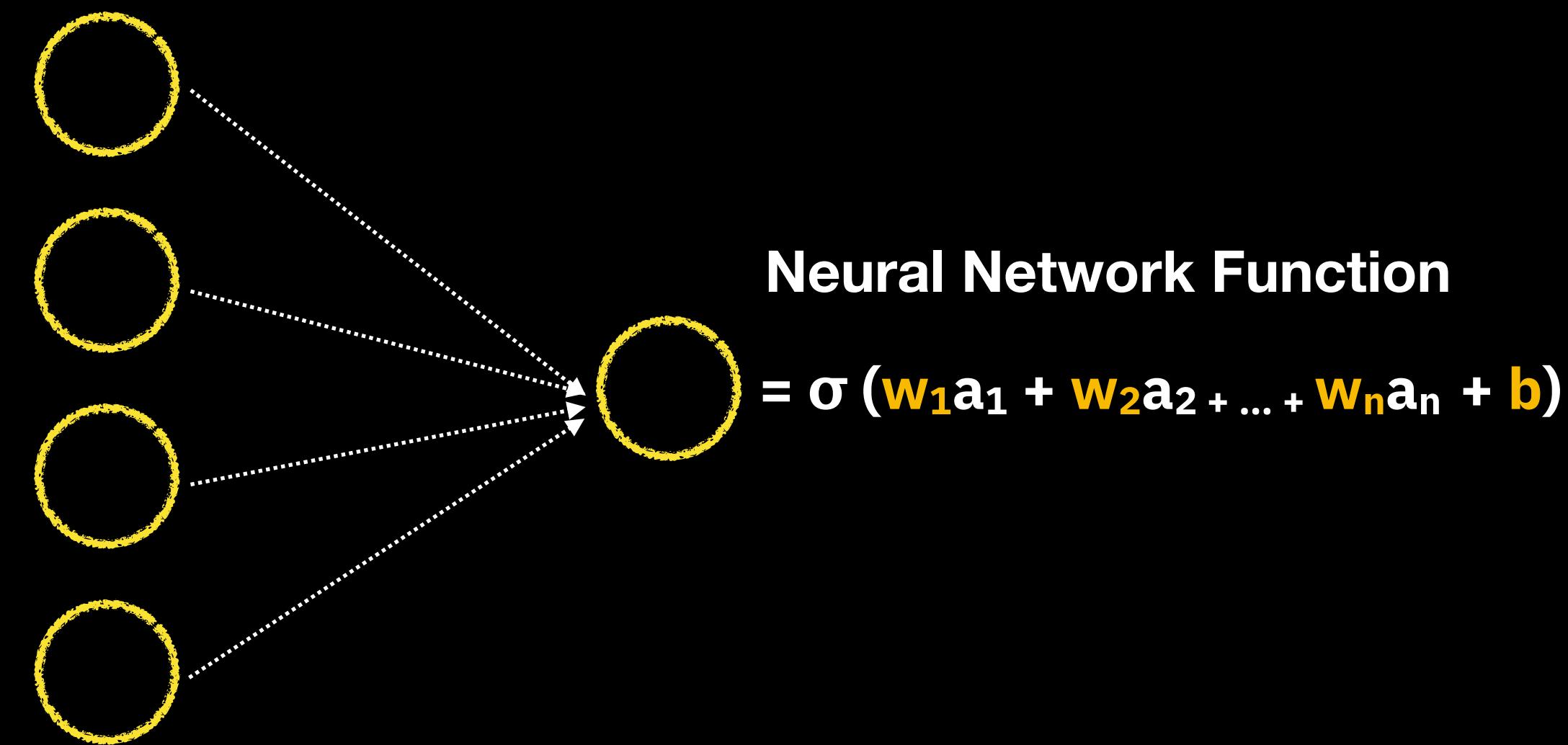
Gradients (Derivatives)

Subtract from each weight a small fraction of its corresponding derivative (Delta Rule)

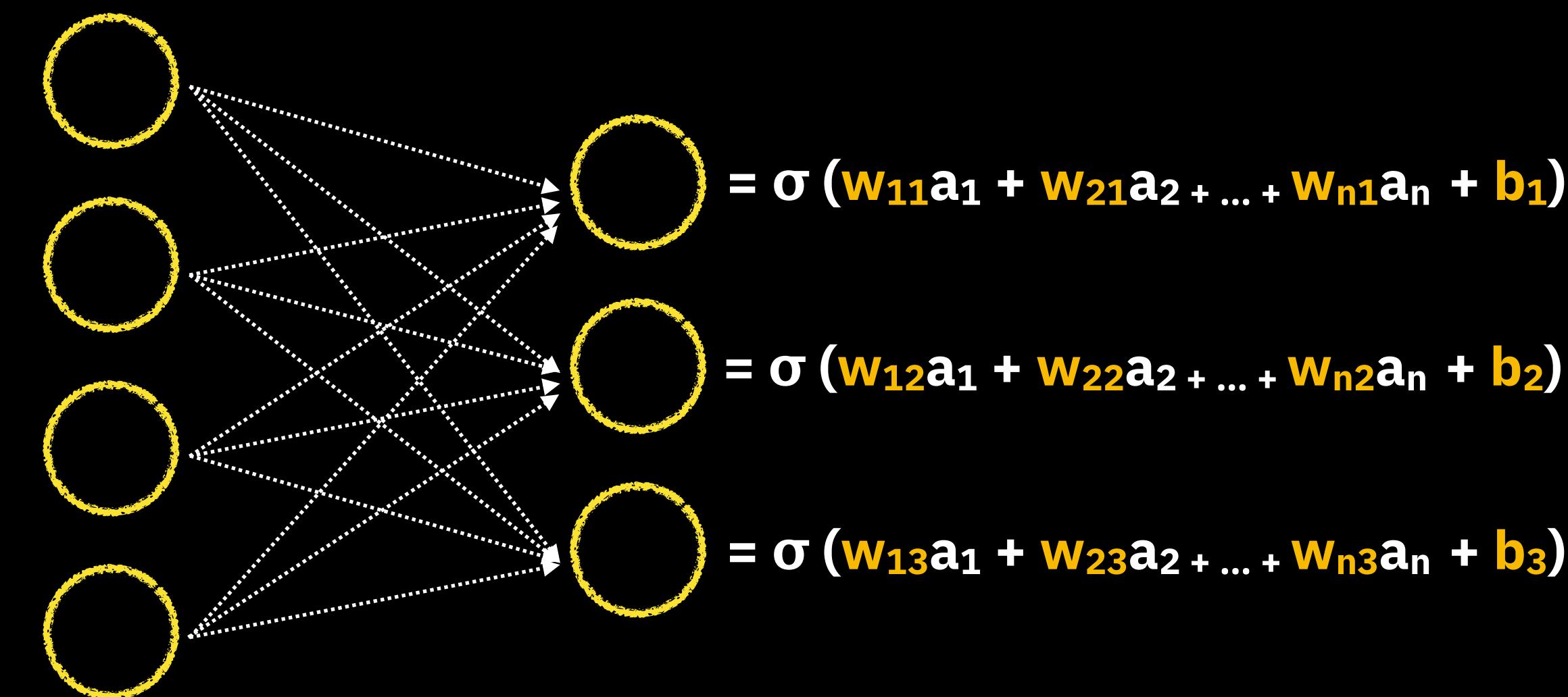
# Fully Connected Layer



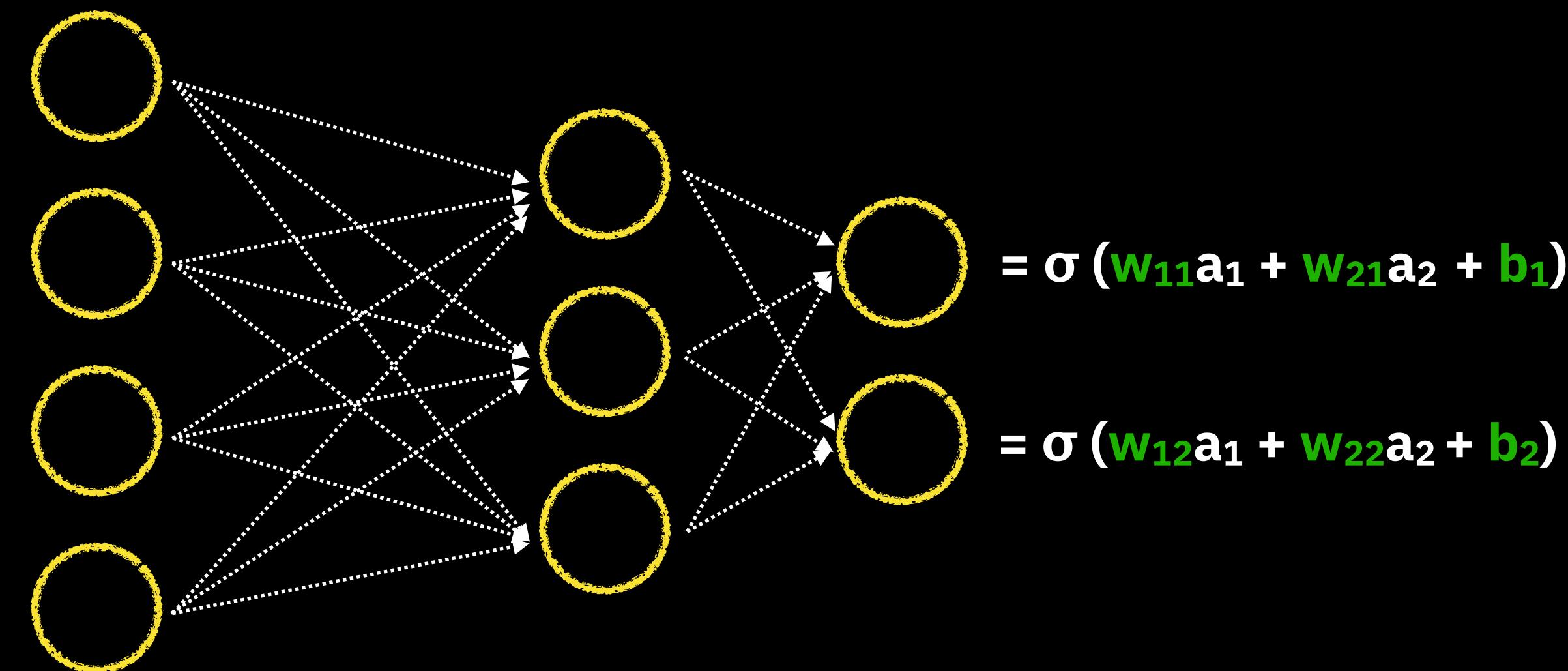
# Fully Connected Layer



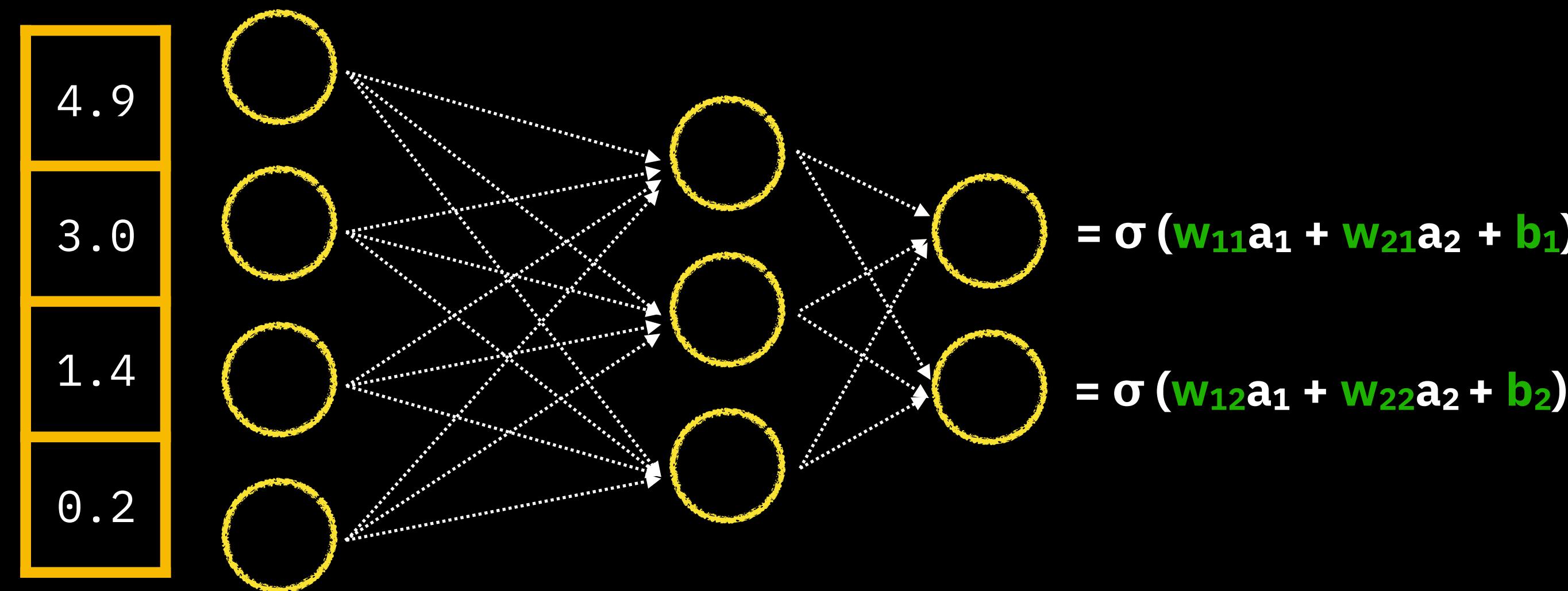
# Fully Connected Layer



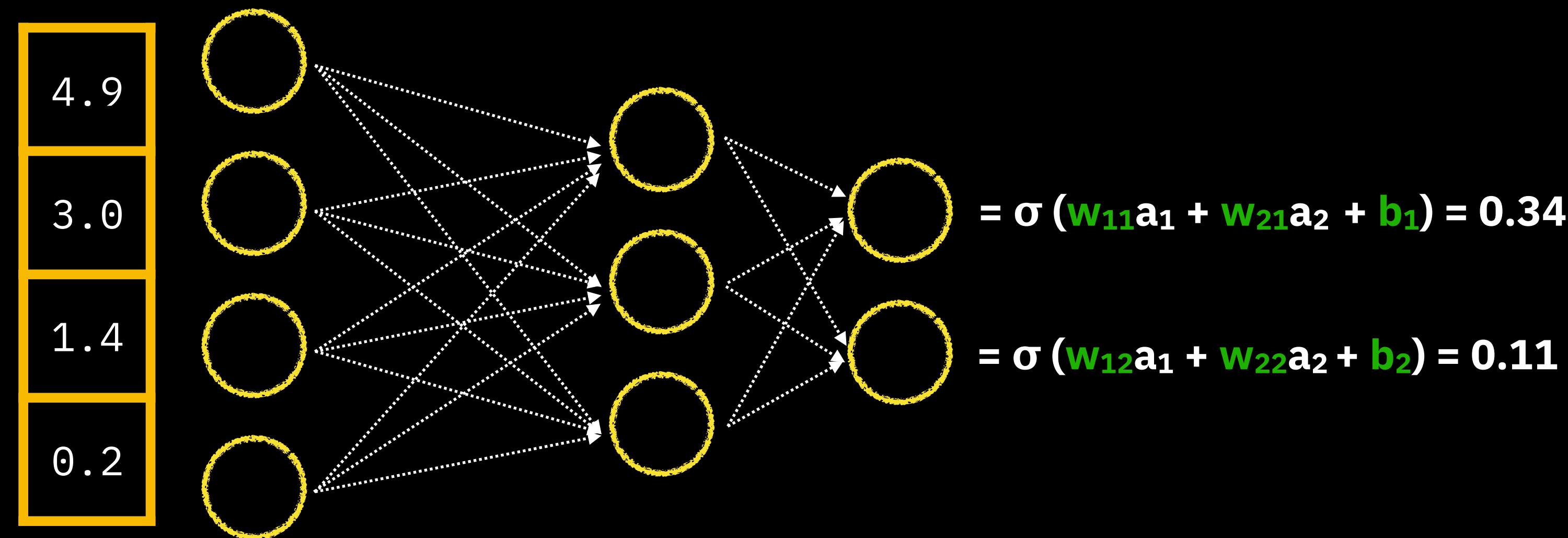
# Fully Connected Layer



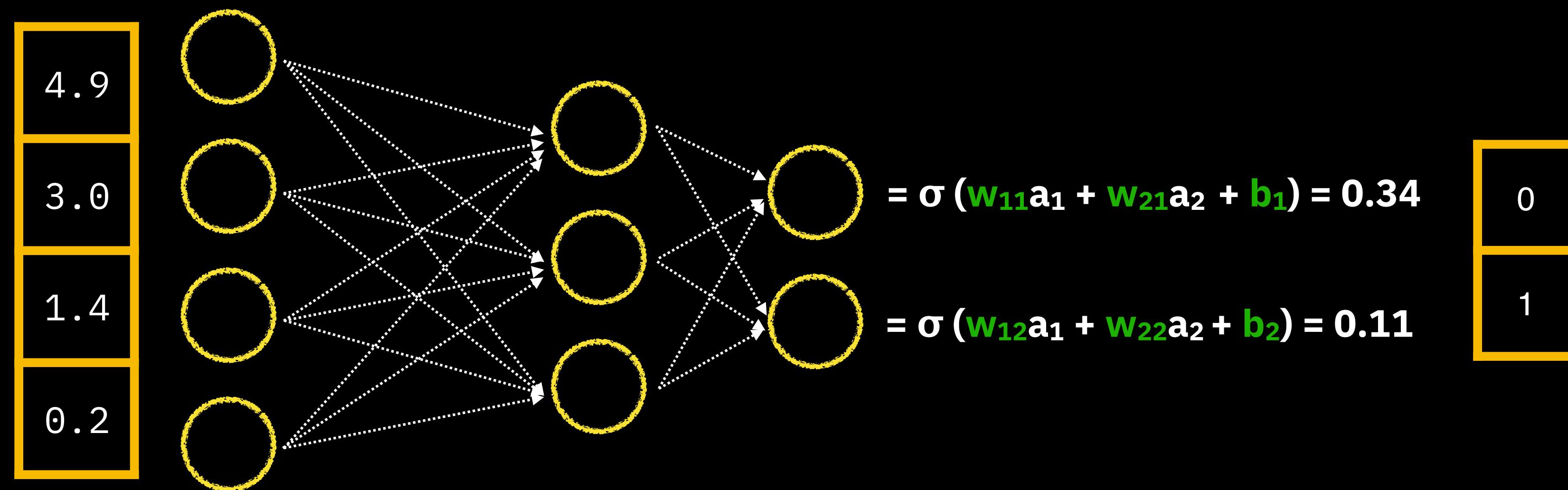
# Fully Connected Layer



# Fully Connected Layer



# Fully Connected Layer



$$(\sigma(w_1a_1 + w_2a_2 + \dots + w_na_n + b) - \text{True})$$

Cost/Loss Function

# Fully Connected Layer

## Neural Function

Input : 4

Output : 3

Parameters :  $((4 \times 3) + 3) + ((3 \times 2) + 2) = 23$

# Fully Connected Layer

## Neural Function

Input : 4

Output : 3

Parameters :  $((4 \times 3) + 3) + ((3 \times 2) + 2) = 23$

## Cost Function

Input : 23

Output : 1 (Cost)

Parameters : Your training examples

## Cost Function

Input : 23

Output : 1 (Cost)

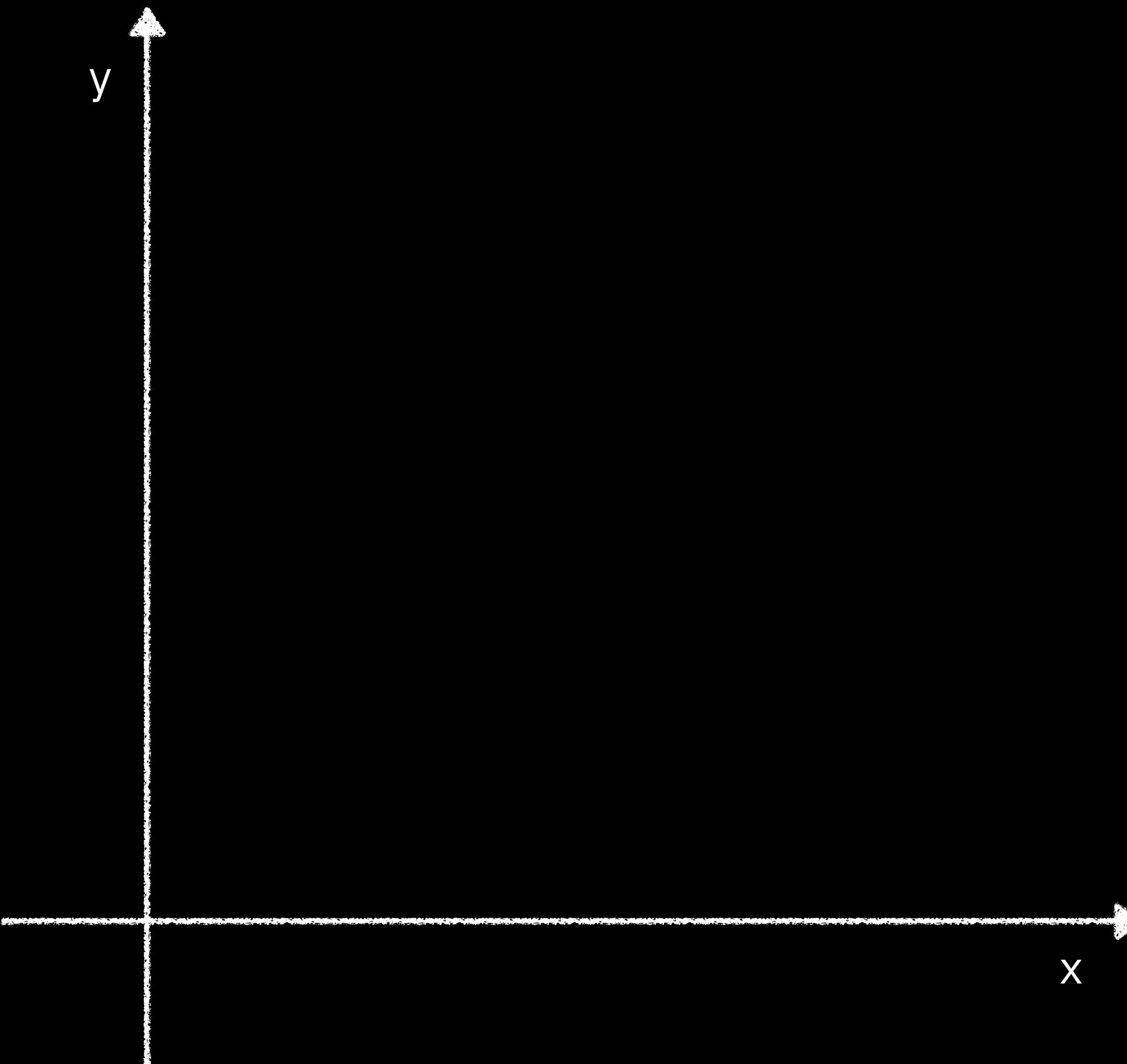
Parameters : Your training examples

$$C(w_1, w_2 \dots w_n)$$

Weights and Biases

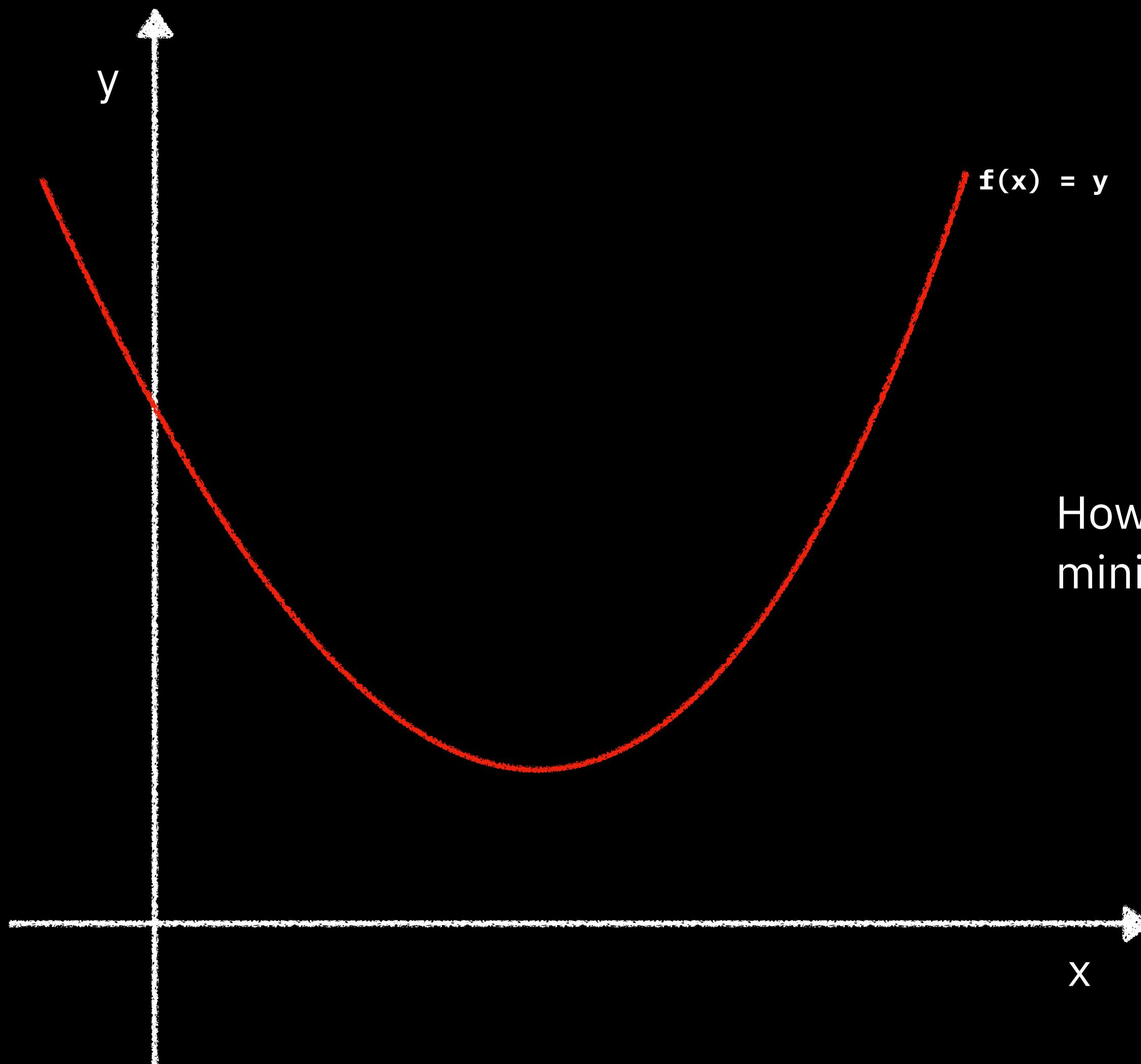
## Cost Function

$C(\mathbf{w})$



## Cost Function

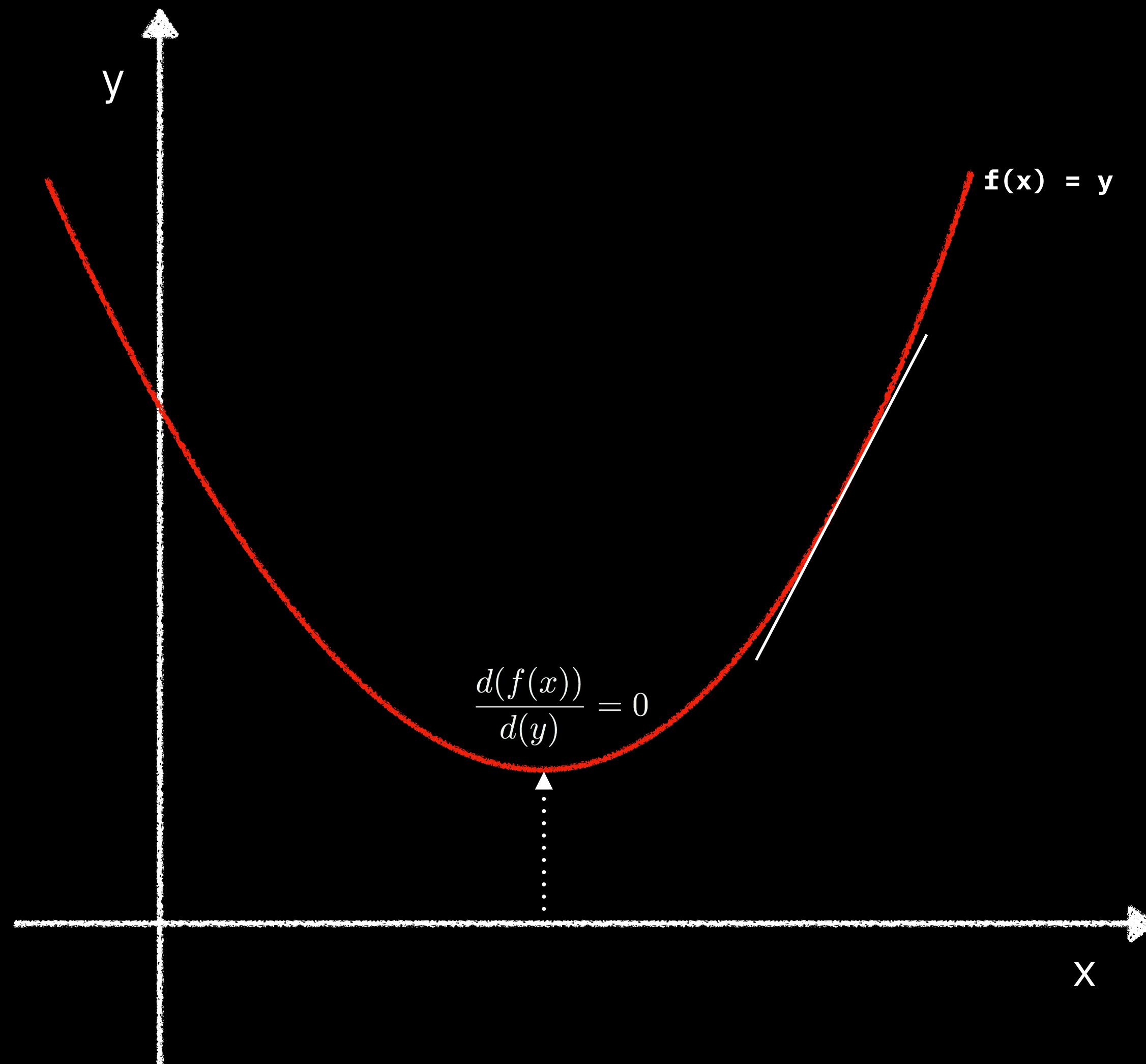
$C(w)$



How do you find the minima in a function?

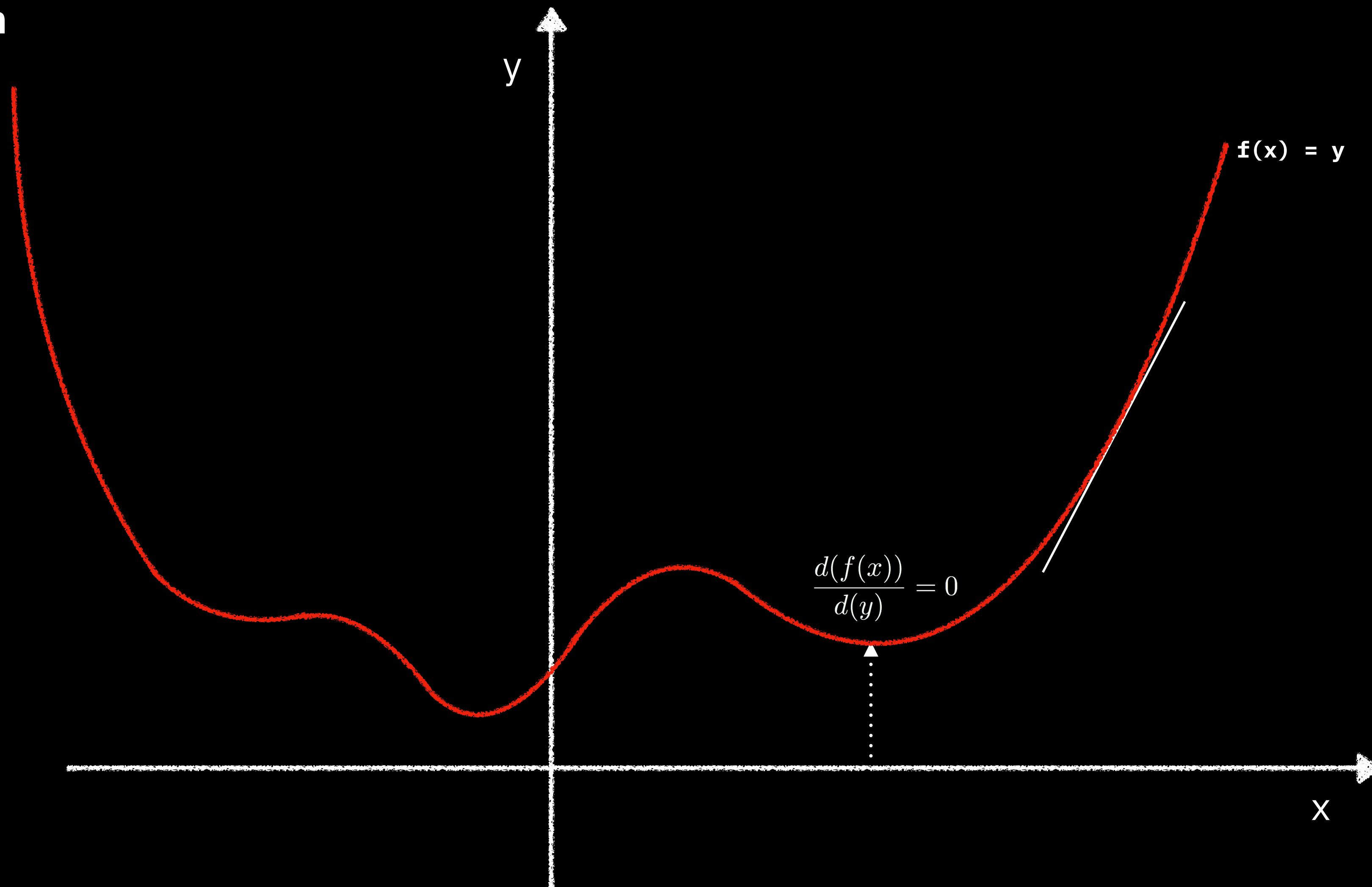
## Cost Function

$C(w)$



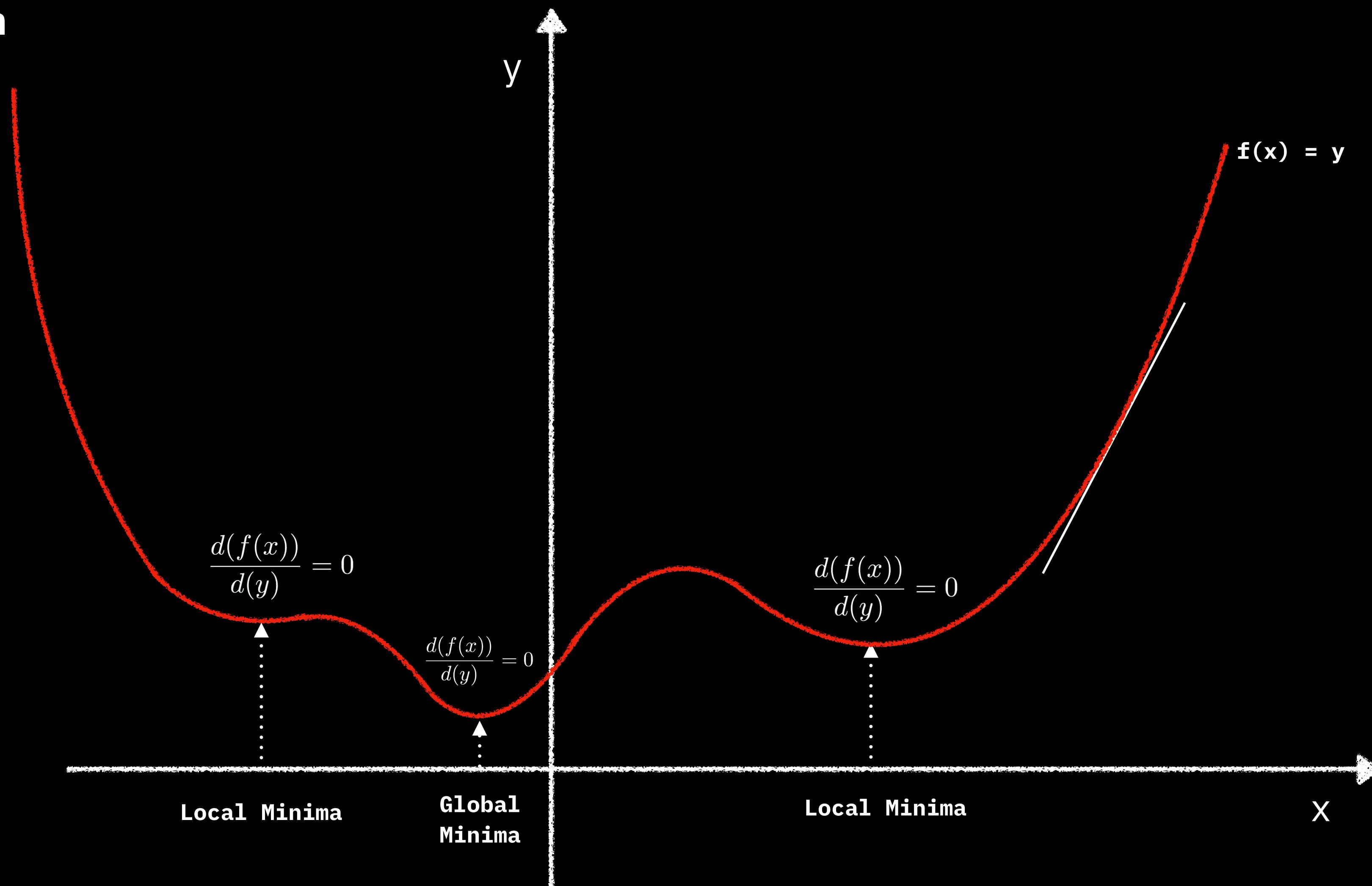
## Cost Function

$C(w)$



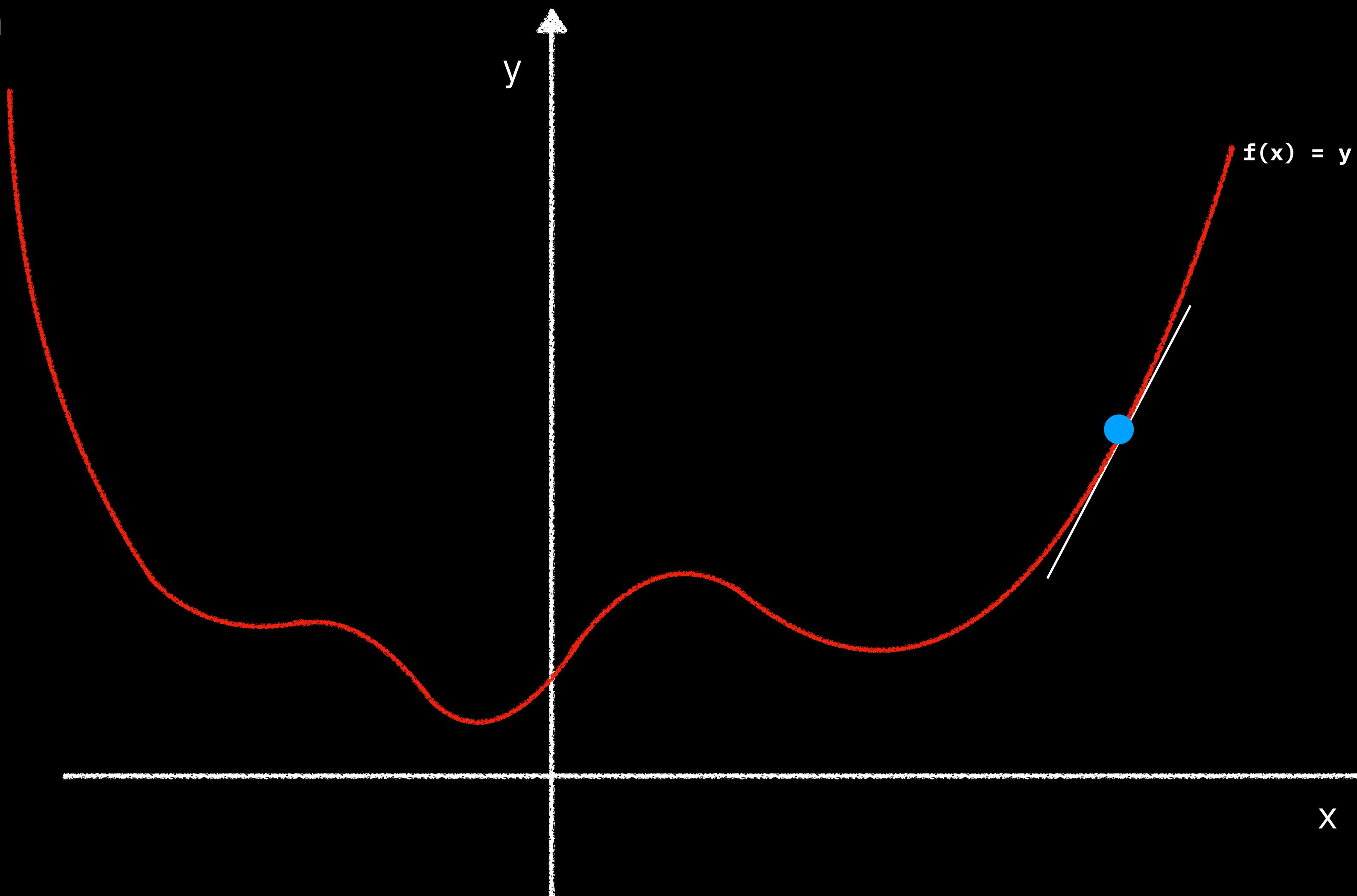
## Cost Function

$C(\mathbf{w})$



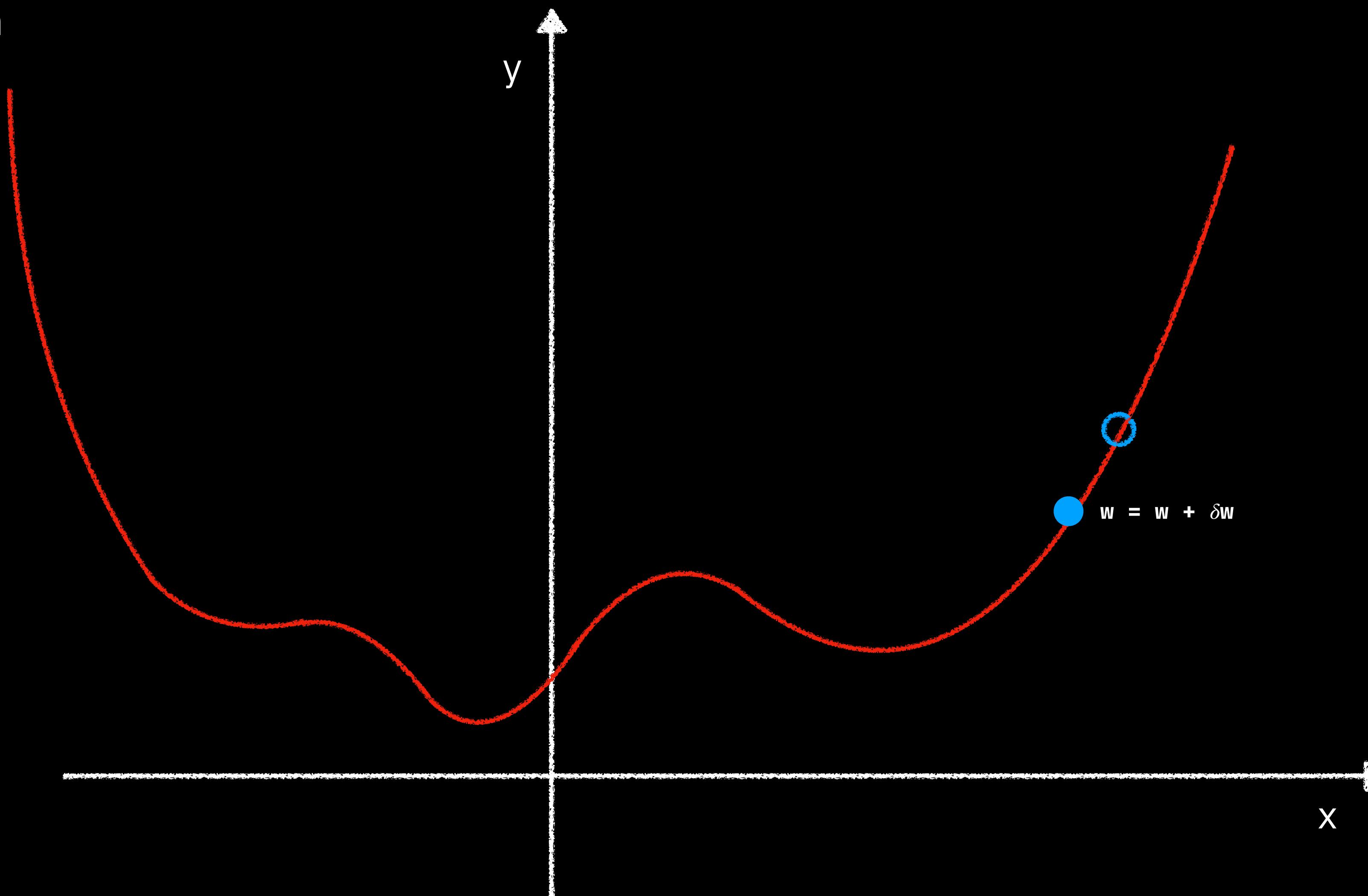
## Cost Function

$C(w)$



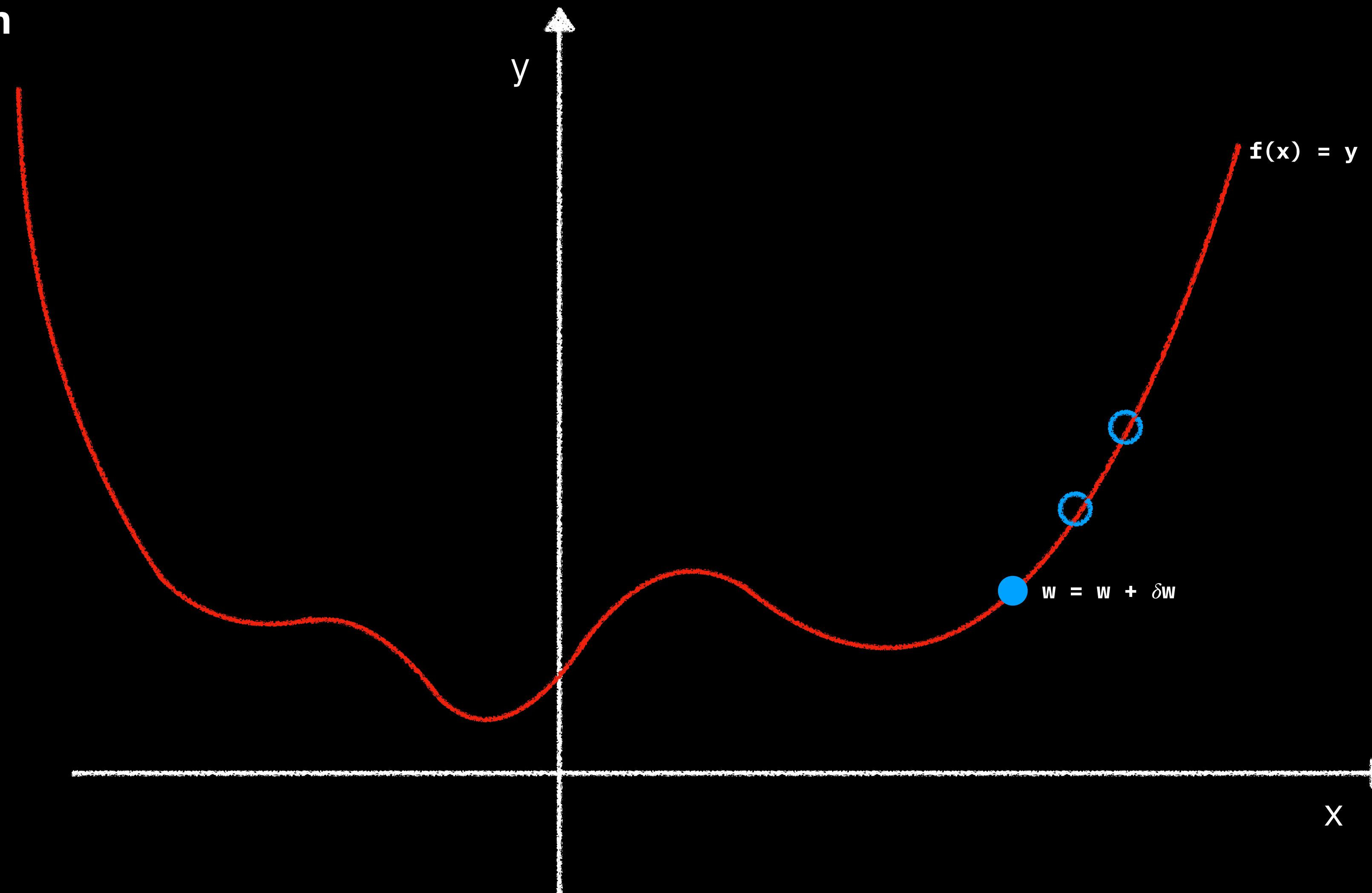
## Cost Function

$C(w)$



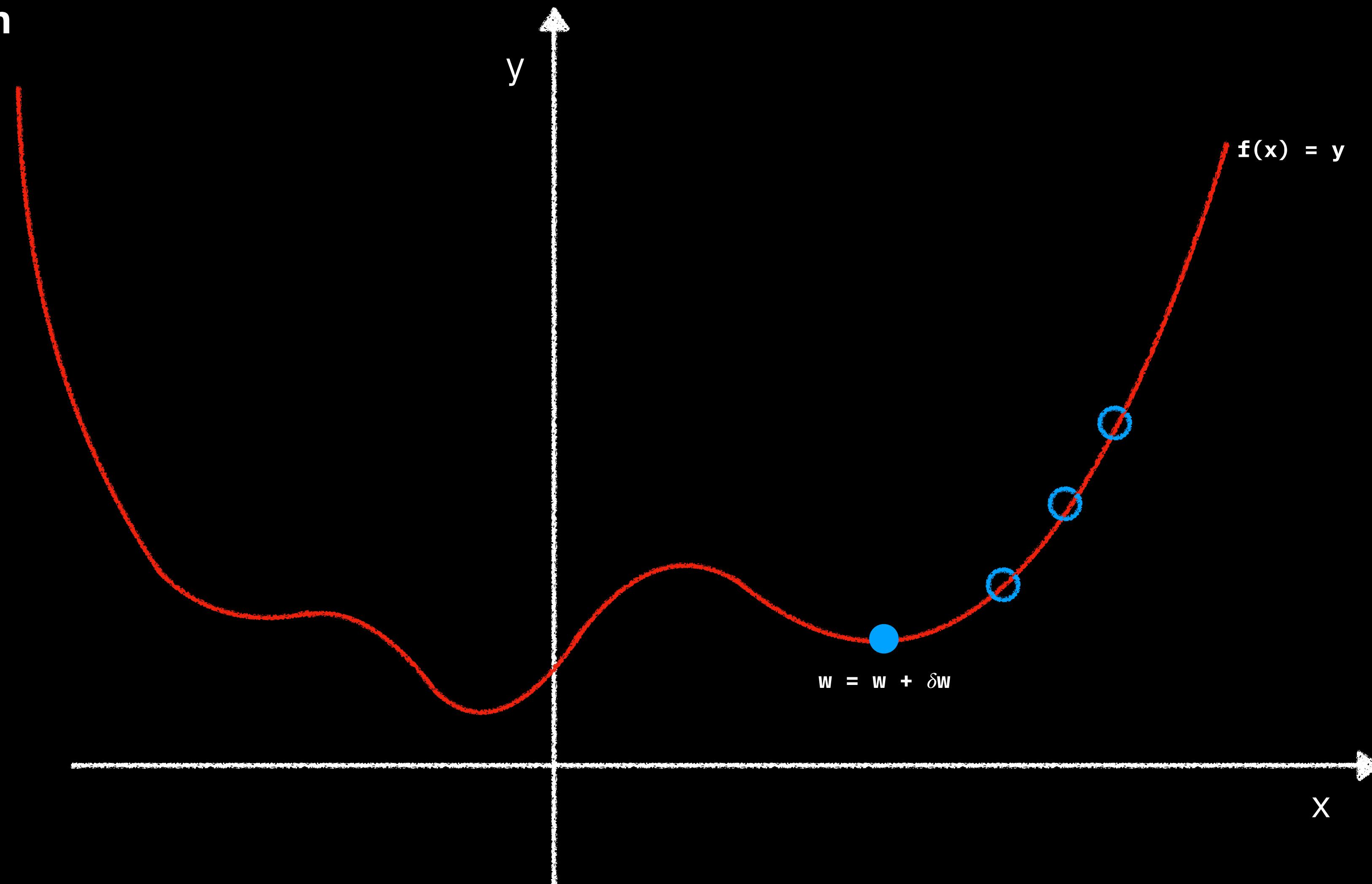
## Cost Function

$C(w)$



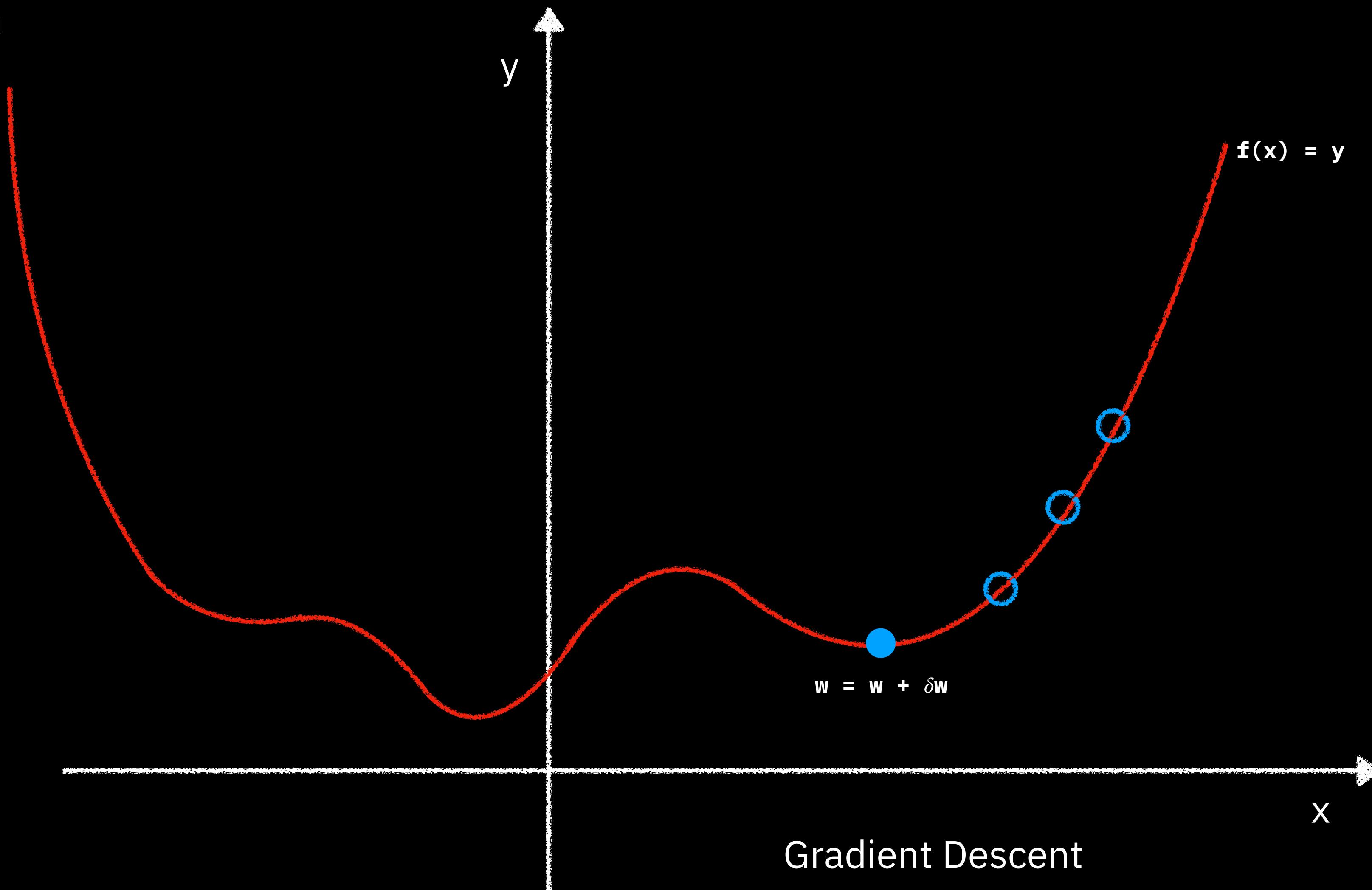
## Cost Function

$C(w)$



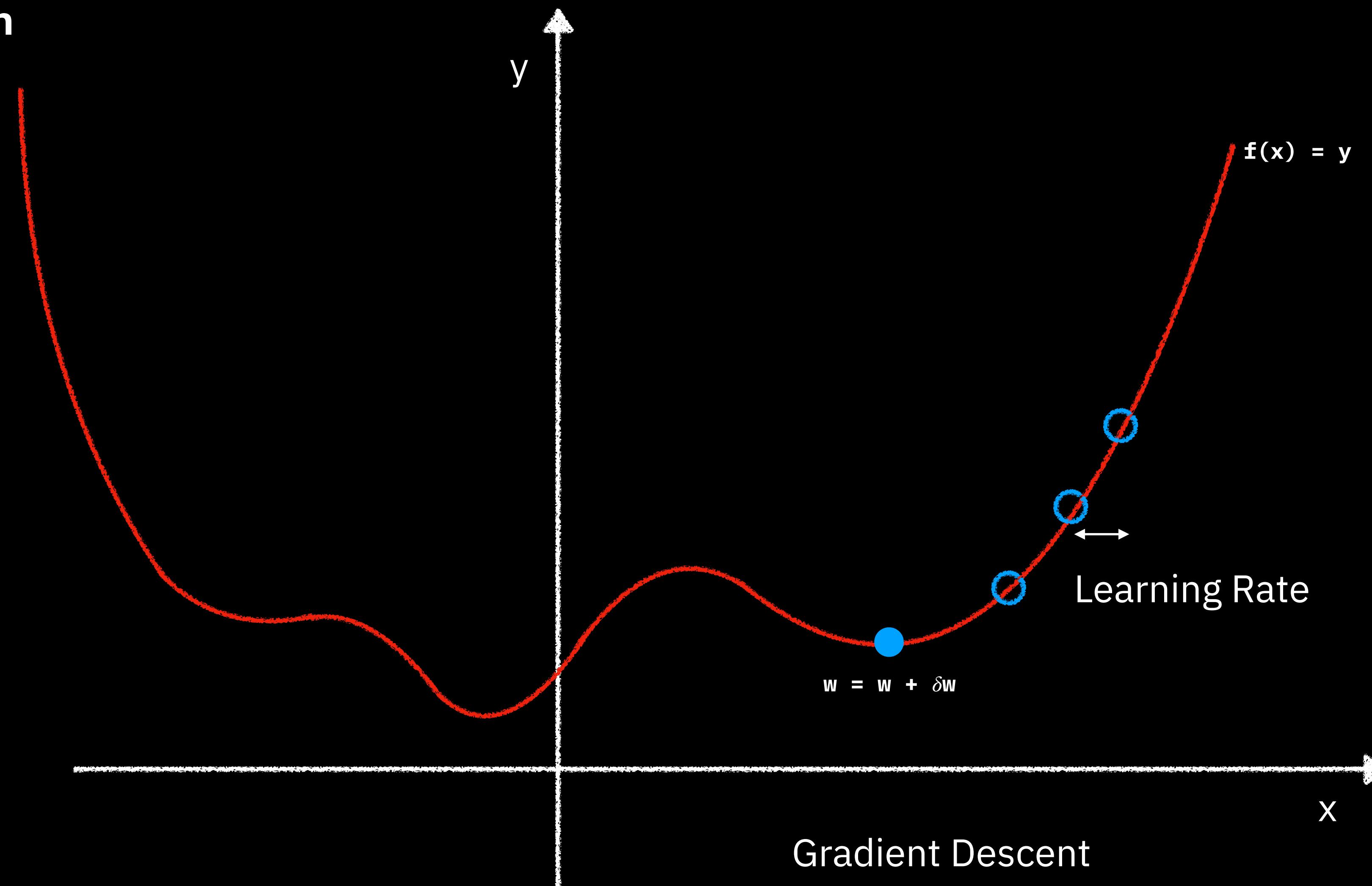
## Cost Function

$C(w)$



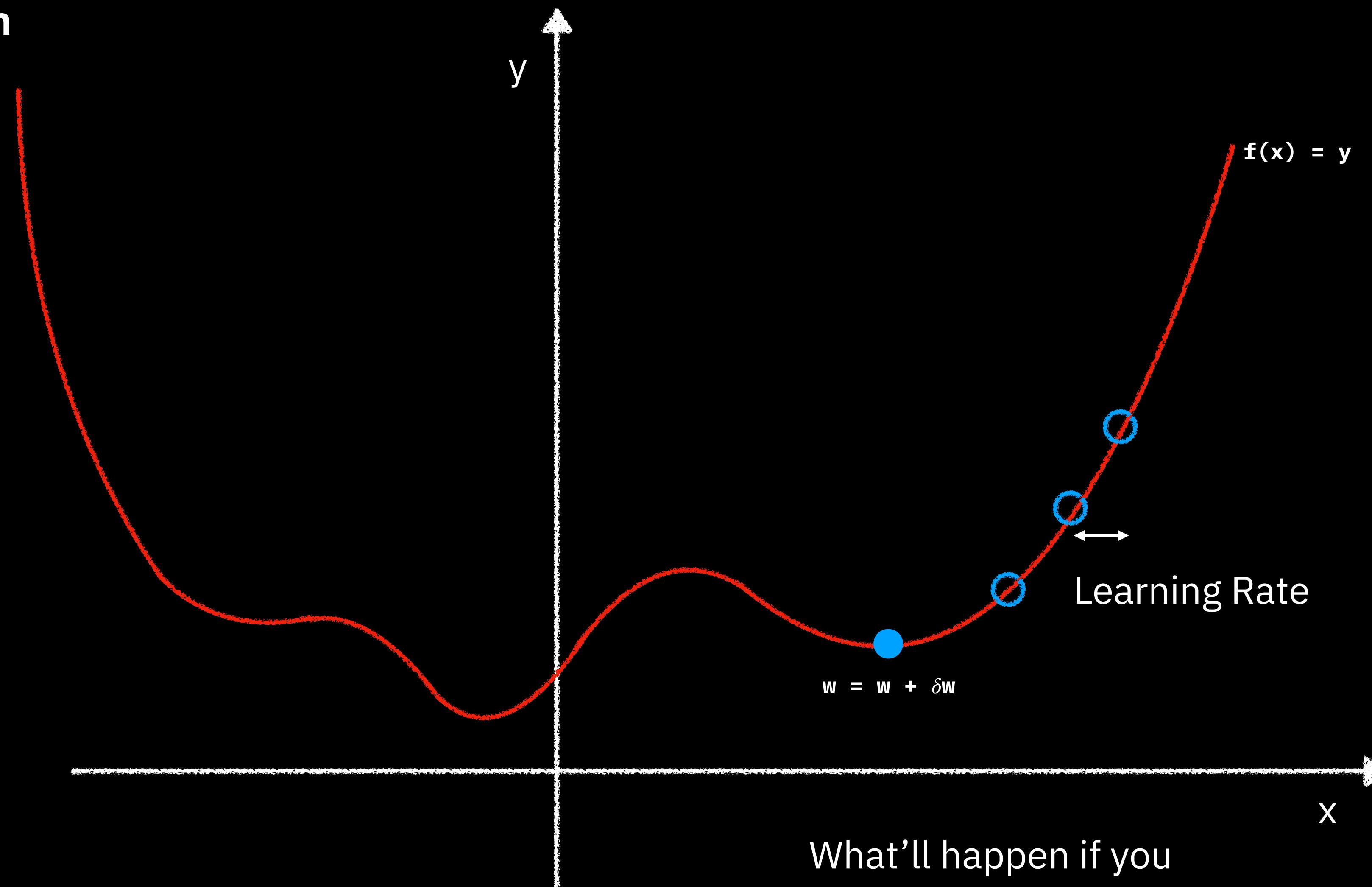
## Cost Function

$C(w)$



## Cost Function

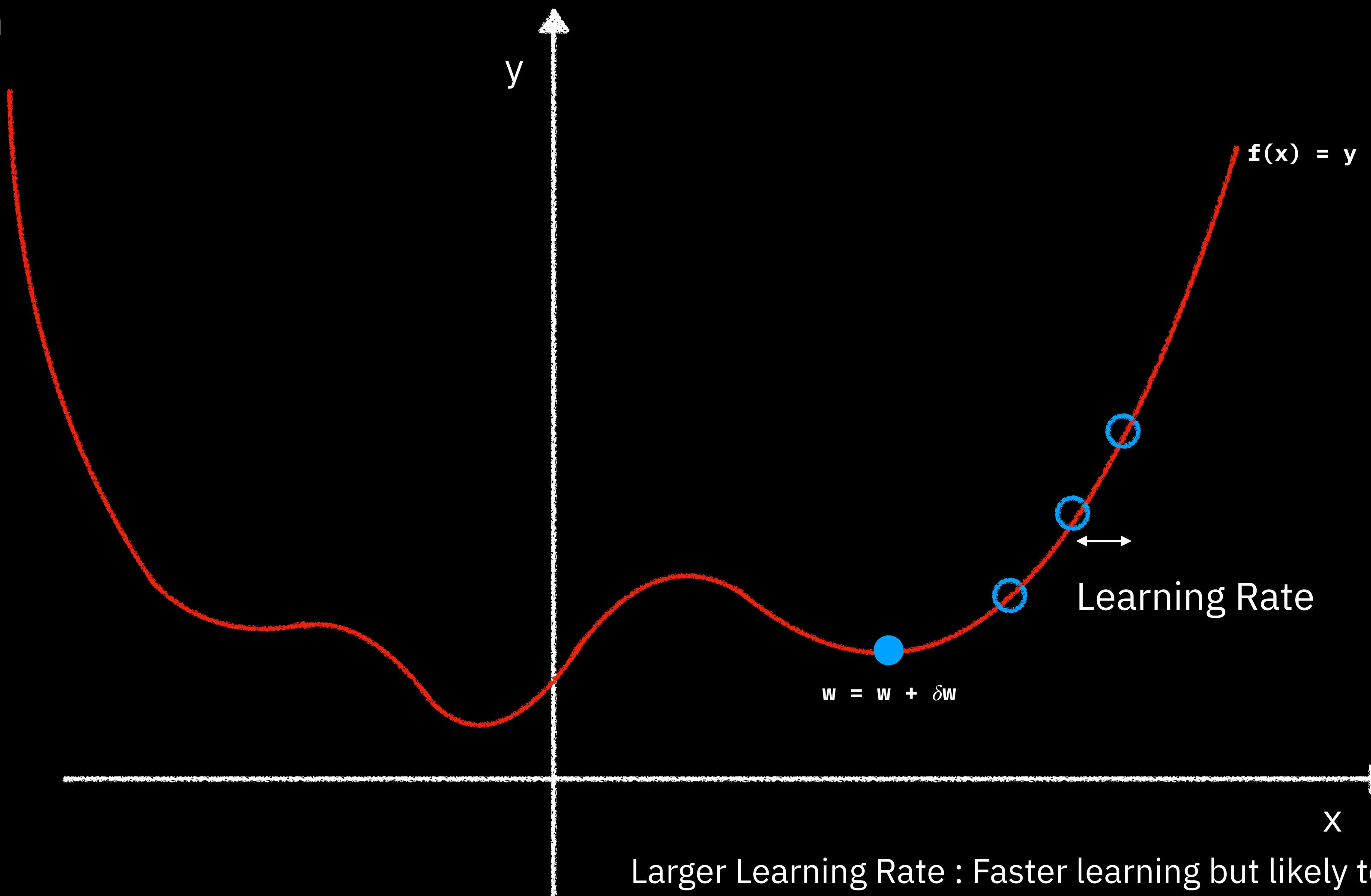
$C(w)$



What'll happen if you  
have a large learning  
rate value?

## Cost Function

$C(w)$

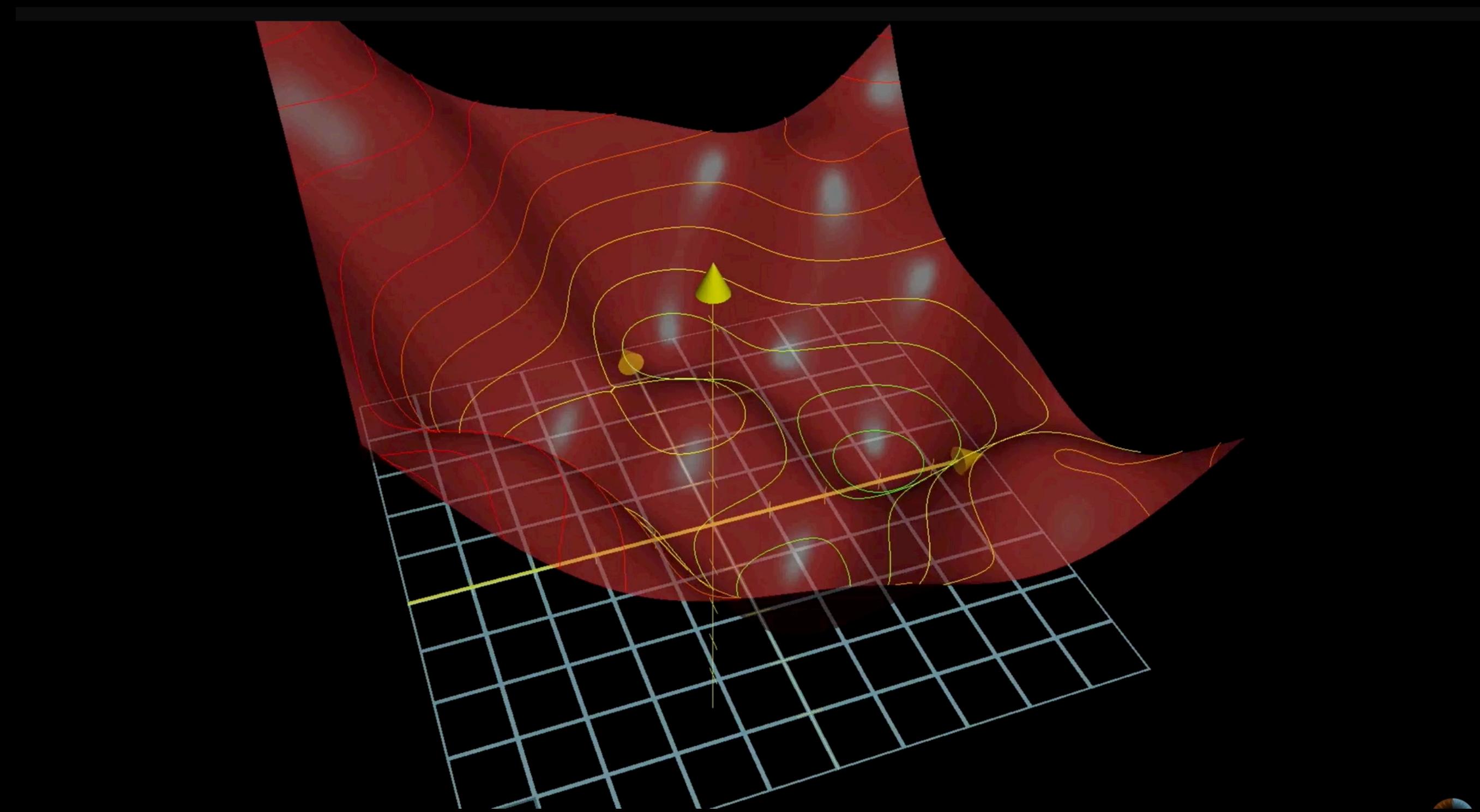


Larger Learning Rate : Faster learning but likely to overshoot the minima.  
(Possibly might never reach minima)  
Small learning rate : Maybe not even learn before your lifetime.

# Cost Function

Neural Network Learning = Minimising a Cost Function

Method to do that : Gradient Descent



// Youtube Video : 3Blue1Brown Gradient Descent

# Improving Neural Networks

Consider a function : Input 1 ; Output 0

Let's play with different initial weights.

Learning Rate constant.

// <http://neuralnetworksanddeeplearning.com/chap3.html>

# Loss Function

Quadratic Loss function

$$\text{Loss} = \sum \frac{1}{2}(\text{target} - \text{output})^2$$

# Loss Function

Quadratic Loss function

$$\text{Loss} = \sum \frac{1}{2}(\text{target} - \text{output})^2$$

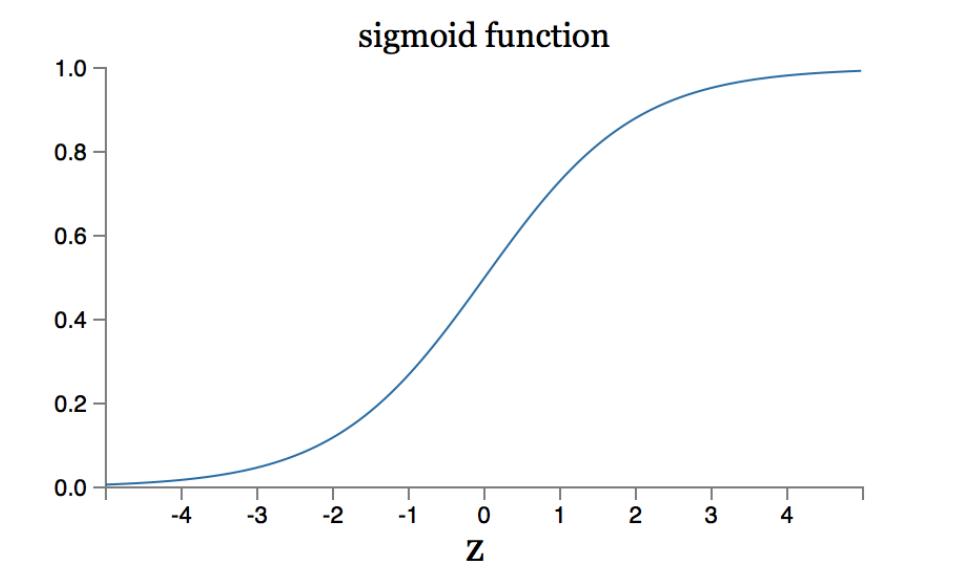
$$z = wx + b$$

$$a = \sigma(z)$$

$$C = \frac{(y - a)^2}{2}$$

$$\frac{\partial C}{\partial w} = (a - y)\sigma'(z)x = a\sigma'(z)$$

$$\frac{\partial C}{\partial b} = (a - y)\sigma'(z) = a\sigma'(z),$$



# Loss Function

## Quadratic Loss function

$$\text{Loss} = \sum \frac{1}{2}(\text{target} - \text{output})^2$$

$$z = wx + b$$

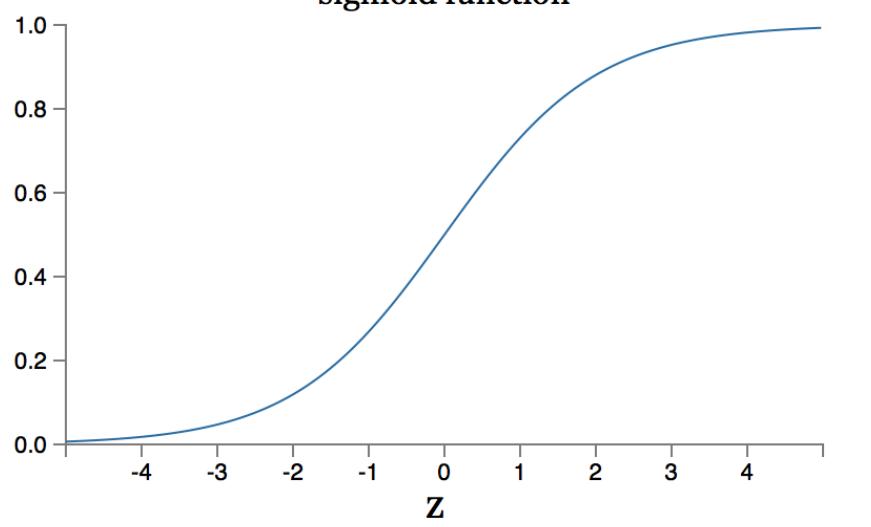
$$a = \sigma(z)$$

$$C = \frac{(y - a)^2}{2}$$

$$\frac{\partial C}{\partial w} = (a - y)\sigma'(z)x = a\sigma'(z)$$

$$\frac{\partial C}{\partial b} = (a - y)\sigma'(z) = a\sigma'(z),$$

sigmoid function



## Cross-entropy Loss function

$$C = -\frac{1}{n} \sum_x [y \ln a + (1 - y) \ln(1 - a)]$$

- Non-Negative ( $C > 0$ )
- Tends towards 0, if the desired output is close to actual output

$$\begin{aligned}\frac{\partial C}{\partial w_j} &= -\frac{1}{n} \sum_x \left( \frac{y}{\sigma(z)} - \frac{(1-y)}{1-\sigma(z)} \right) \frac{\partial \sigma}{\partial w_j} \\ &= -\frac{1}{n} \sum_x \left( \frac{y}{\sigma(z)} - \frac{(1-y)}{1-\sigma(z)} \right) \sigma'(z) x_j.\end{aligned}$$

$$\frac{\partial C}{\partial w_j} = \frac{1}{n} \sum_x x_j (\sigma(z) - y).$$

# Loss Function

## Quadratic Loss function

$$\text{Loss} = \sum \frac{1}{2}(\text{target} - \text{output})^2$$

$$z = wx + b$$

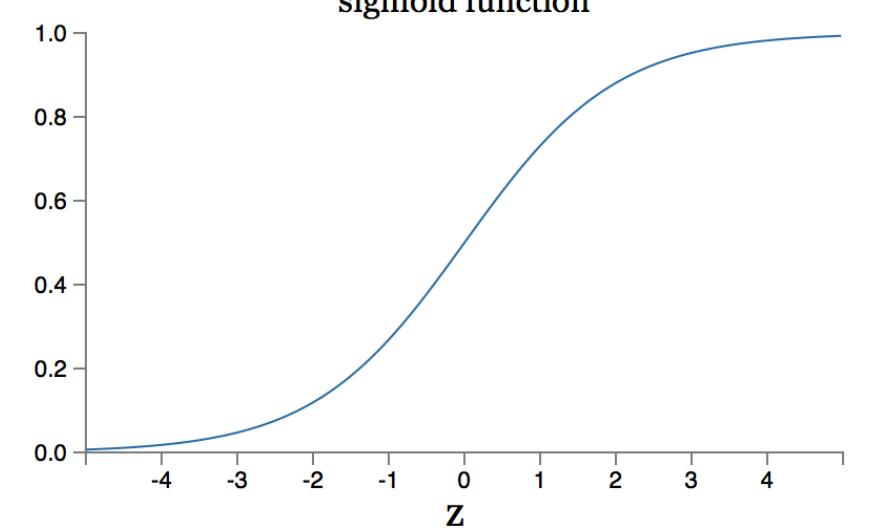
$$a = \sigma(z)$$

$$C = \frac{(y - a)^2}{2}$$

$$\frac{\partial C}{\partial w} = (a - y)\sigma'(z)x = a\sigma'(z)$$

$$\frac{\partial C}{\partial b} = (a - y)\sigma'(z) = a\sigma'(z),$$

sigmoid function



## Cross-entropy Loss function

$$C = -\frac{1}{n} \sum_x [y \ln a + (1 - y) \ln(1 - a)]$$

- Non-Negative ( $C > 0$ )
- Tends towards 0, if the desired output is close to actual output

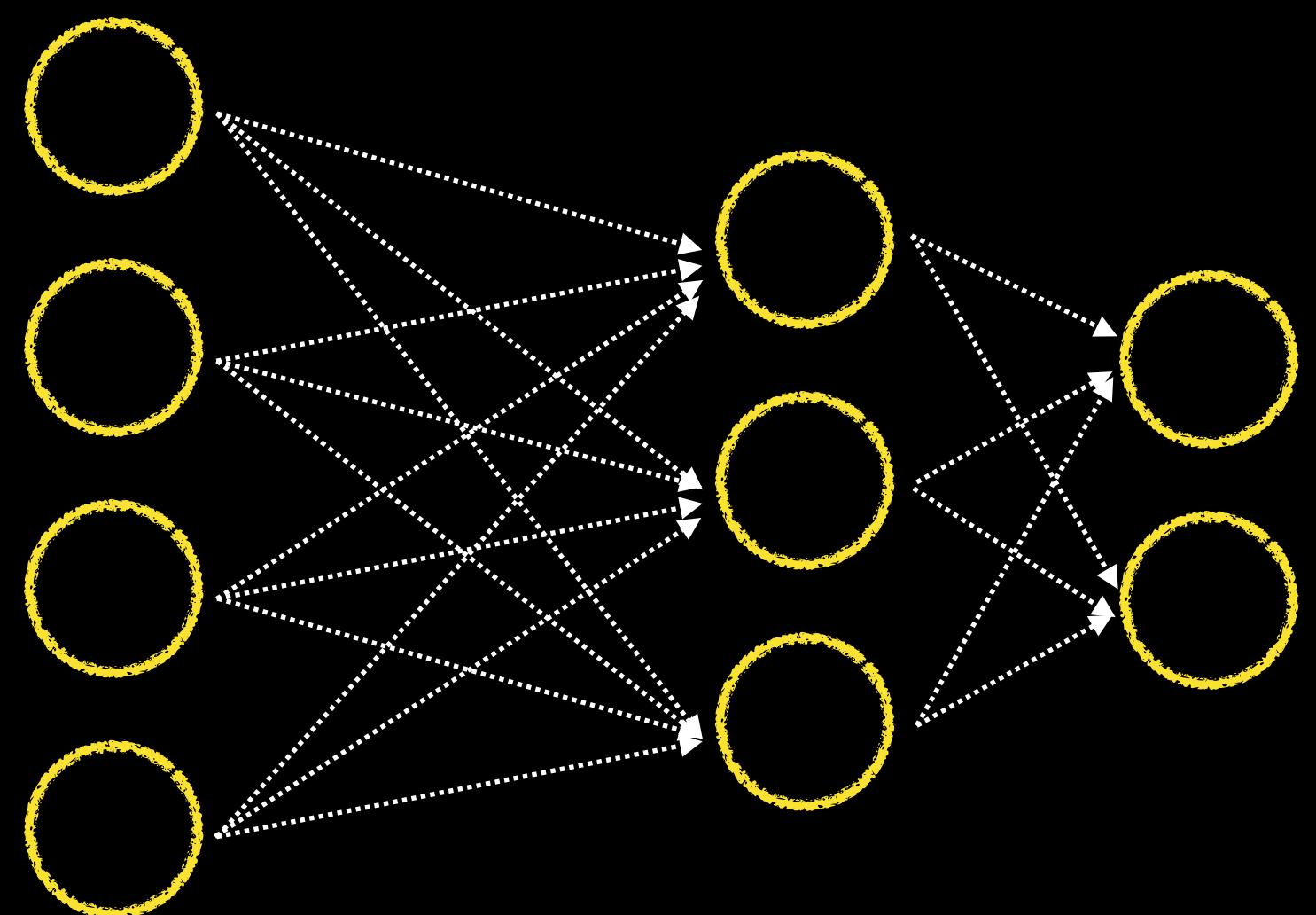
$$\begin{aligned}\frac{\partial C}{\partial w_j} &= -\frac{1}{n} \sum_x \left( \frac{y}{\sigma(z)} - \frac{(1-y)}{1-\sigma(z)} \right) \frac{\partial \sigma}{\partial w_j} \\ &= -\frac{1}{n} \sum_x \left( \frac{y}{\sigma(z)} - \frac{(1-y)}{1-\sigma(z)} \right) \sigma'(z) x_j.\end{aligned}$$

$$\frac{\partial C}{\partial w_j} = \frac{1}{n} \sum_x x_j (\sigma(z) - y).$$

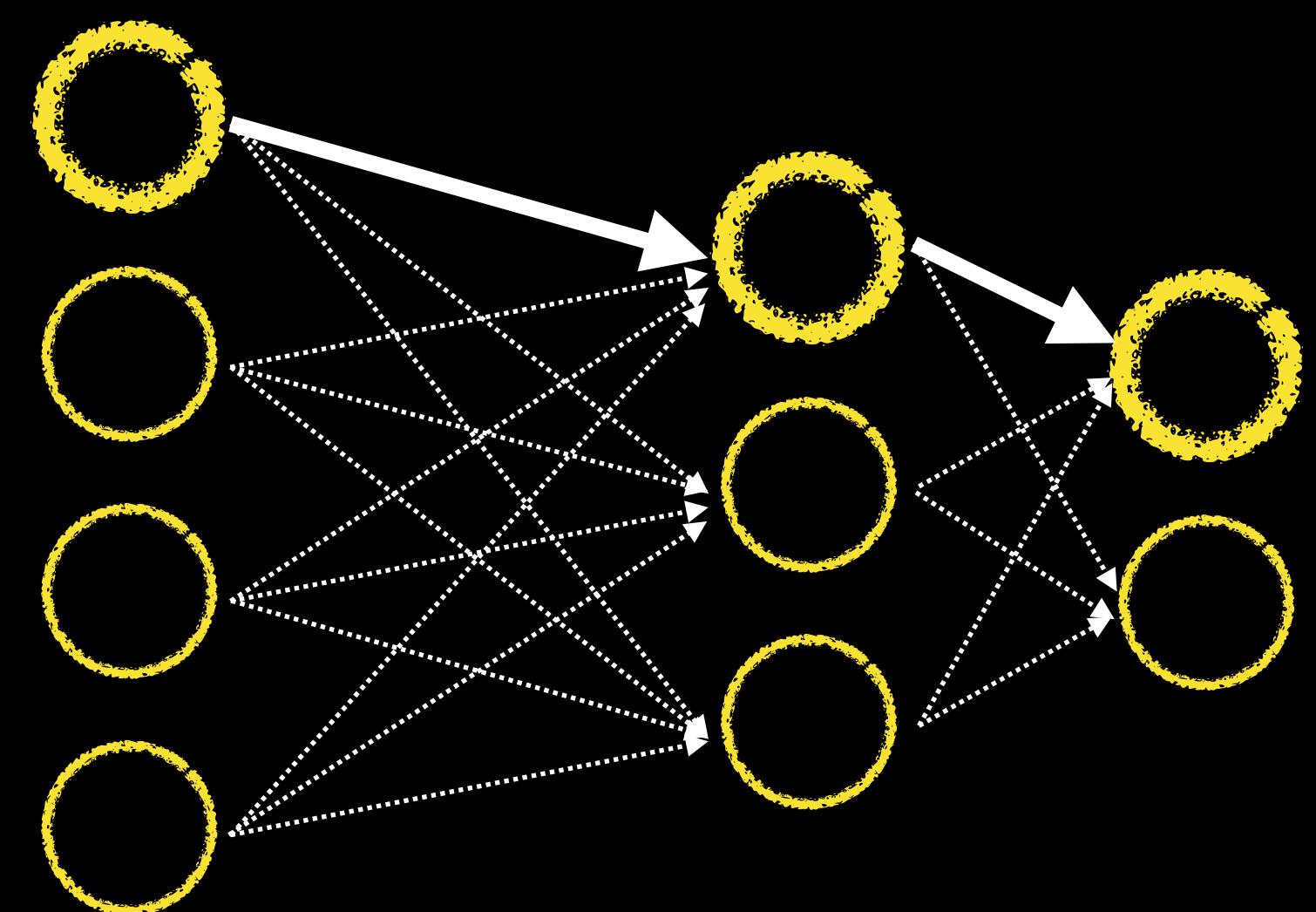
### Homework :

- Think a little more about the origin of this problem.
- Come up with how cross-entropy term can derived from that problem.
- Find the derivative of the sigmoid function.

# Computational Graphs (Backpropagation)



# Computational Graphs (Backpropagation)

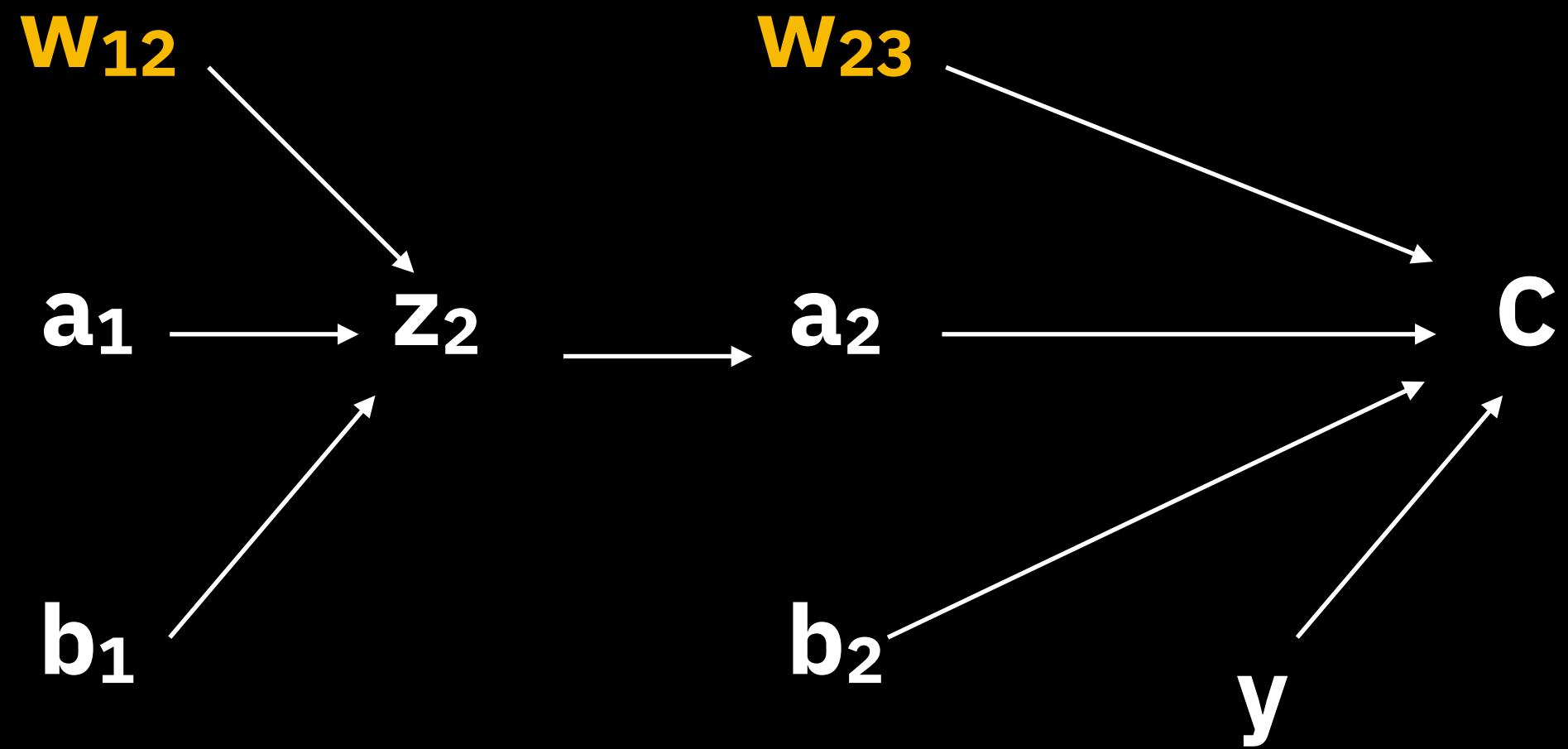
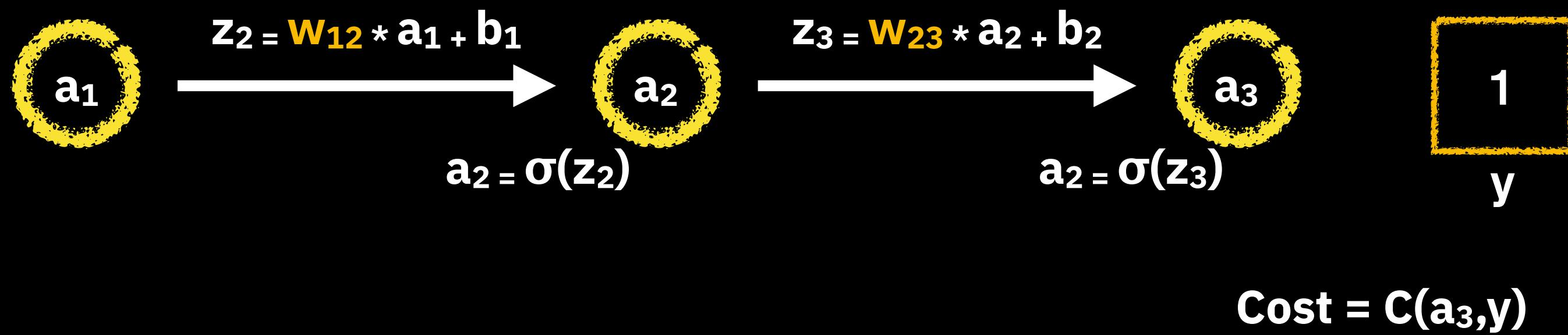


# Computational Graphs (Backpropagation)

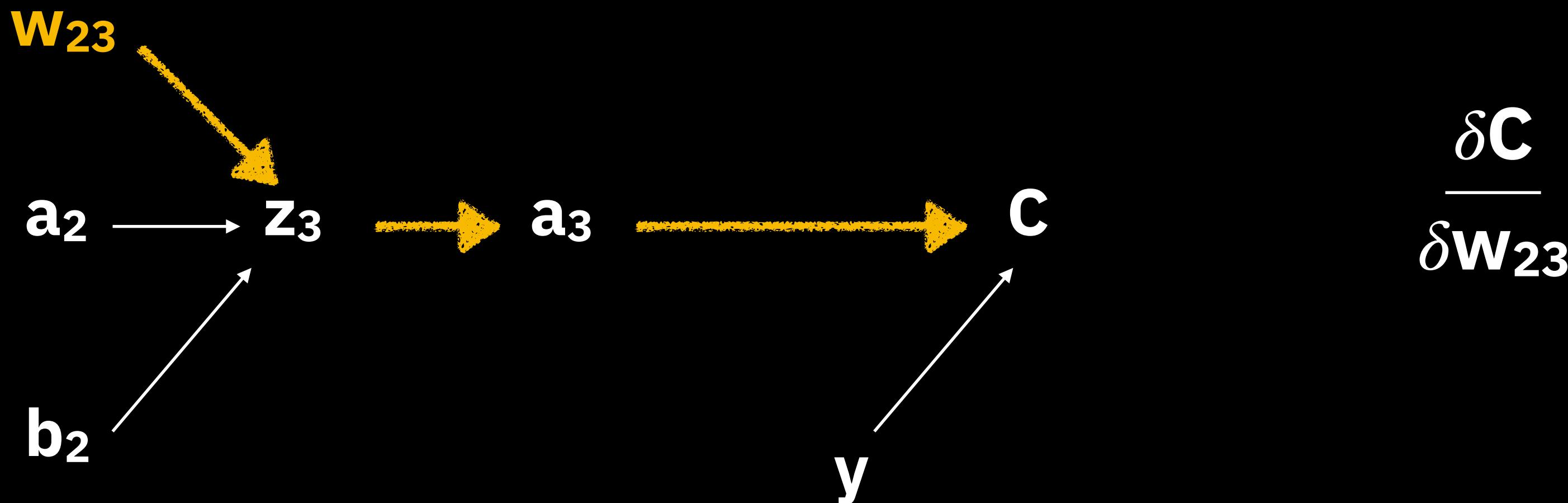
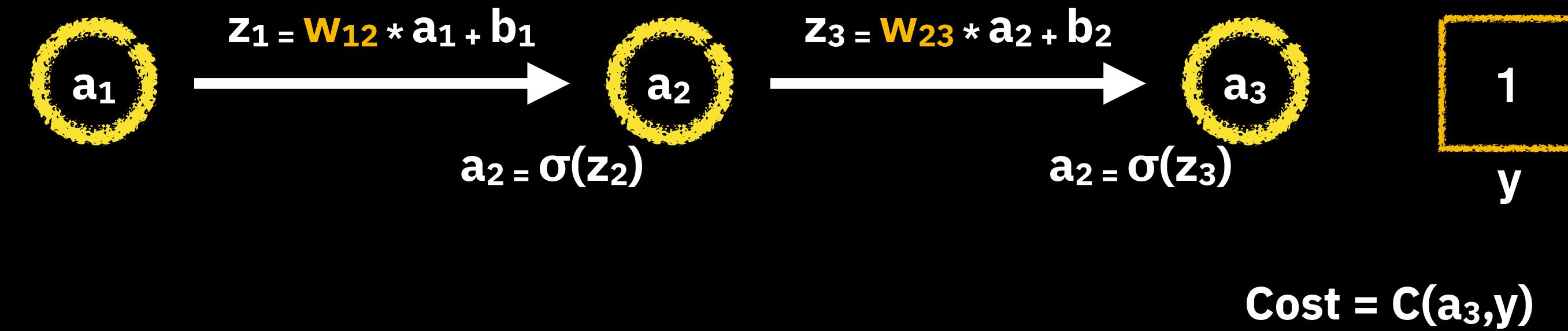


$$\text{Cost} = C(a_3, y)$$

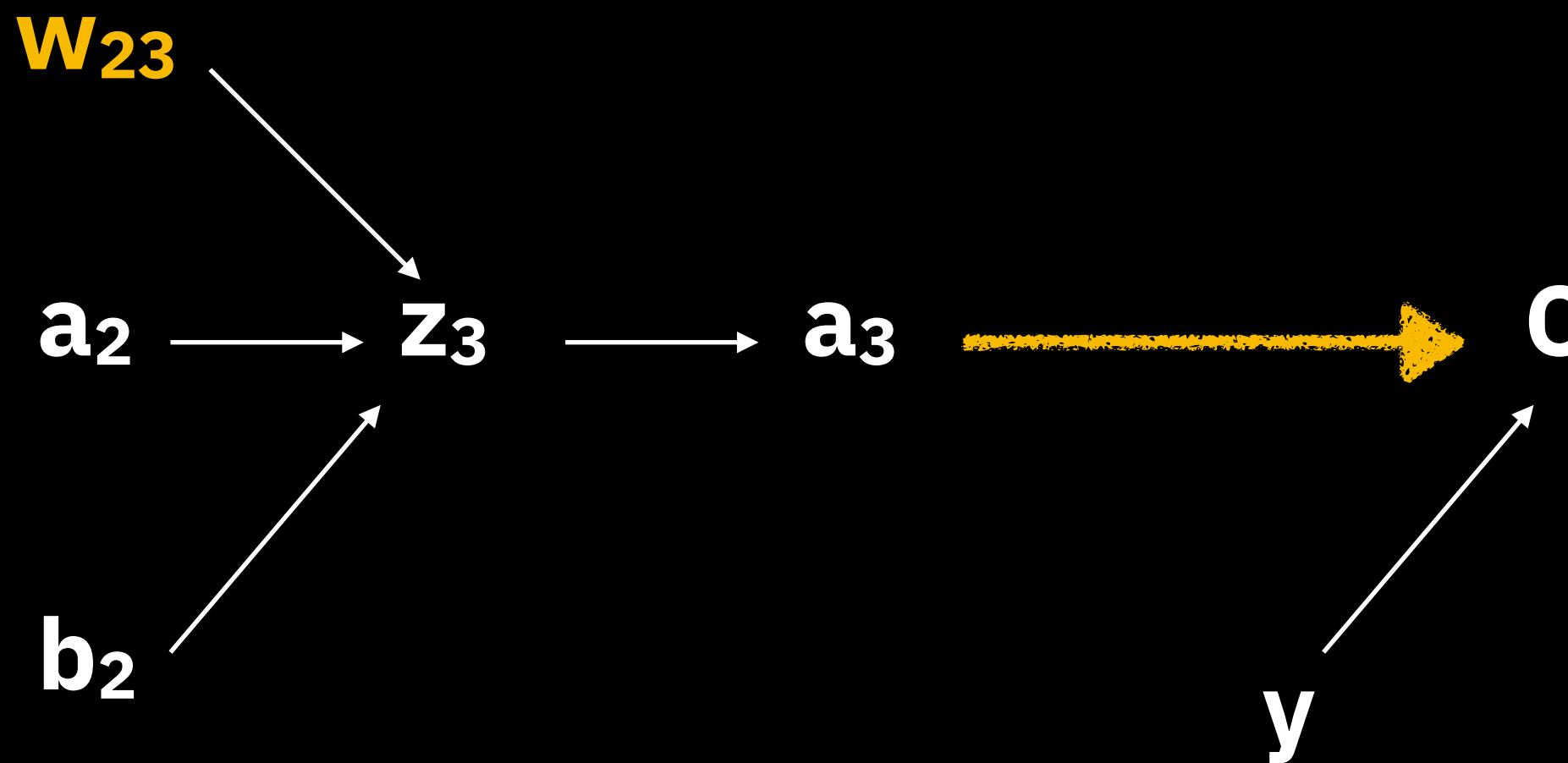
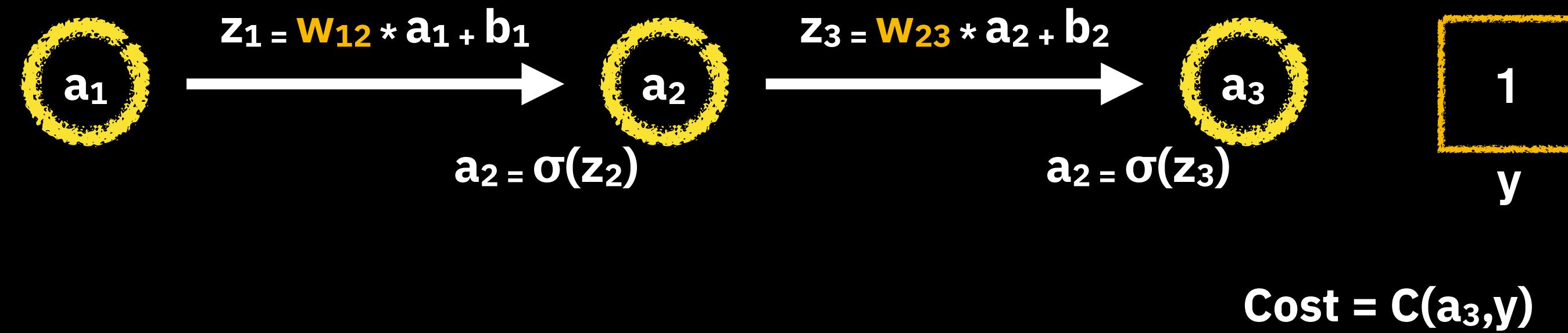
# Computational Graphs (Backpropagation)



# Computational Graphs (Backpropagation)

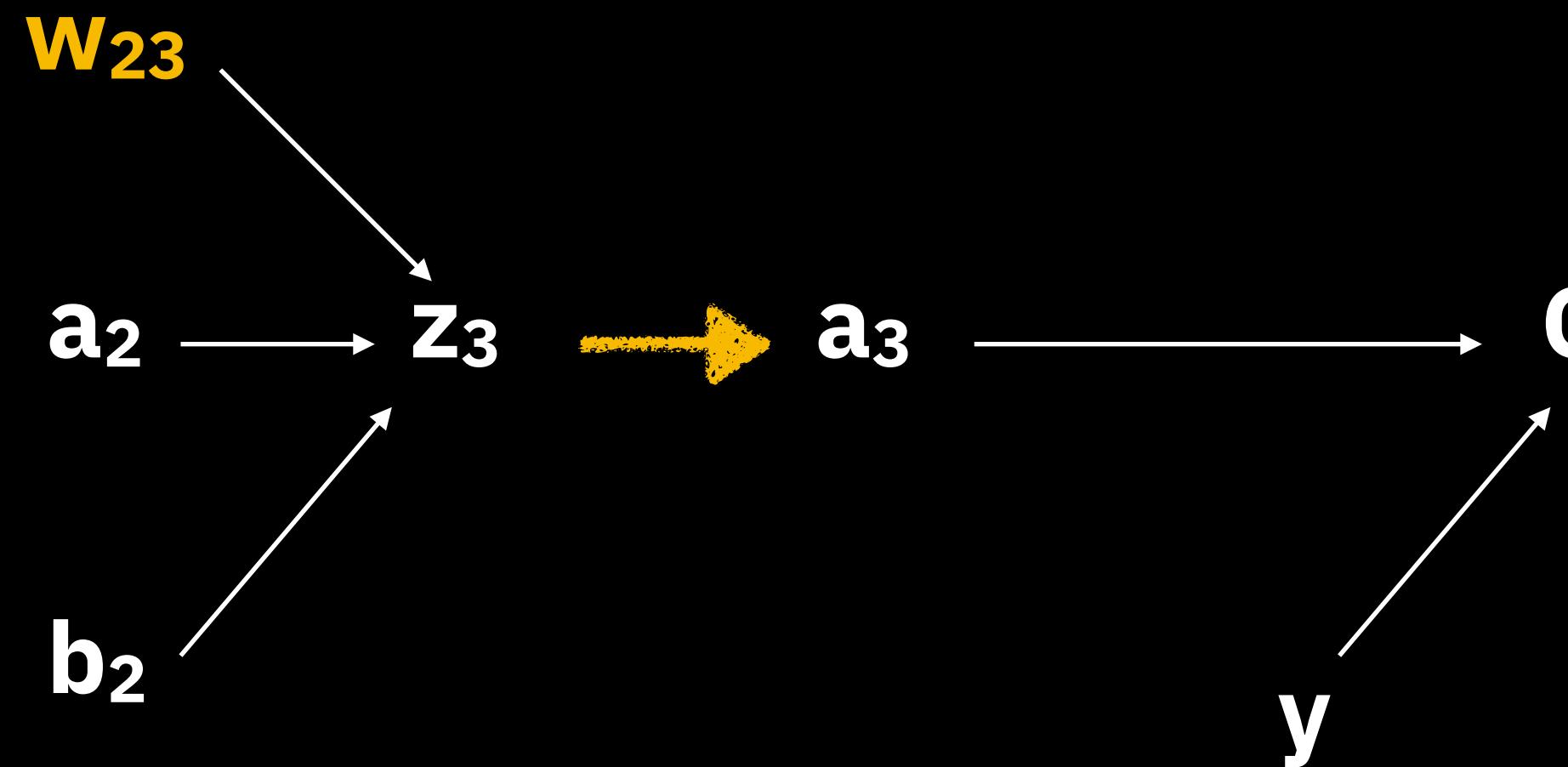
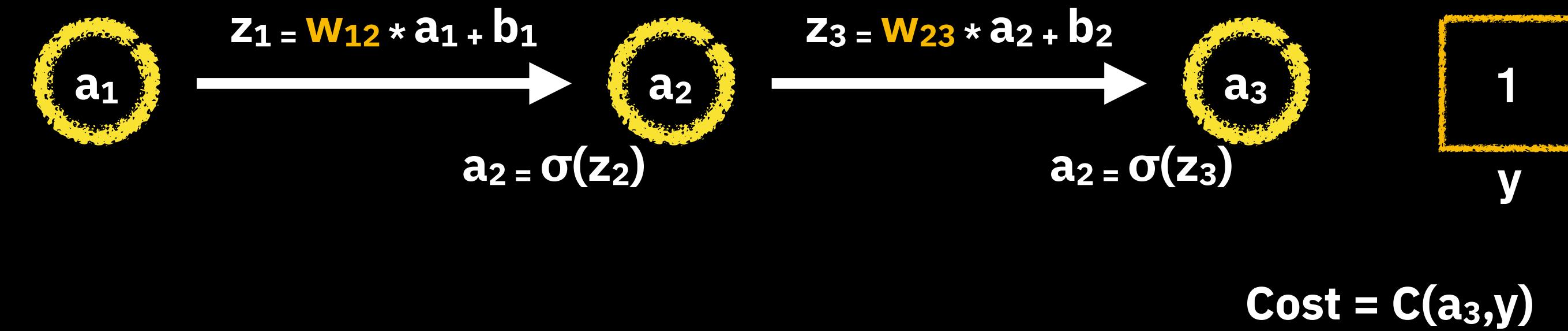


# Computational Graphs (Backpropagation)



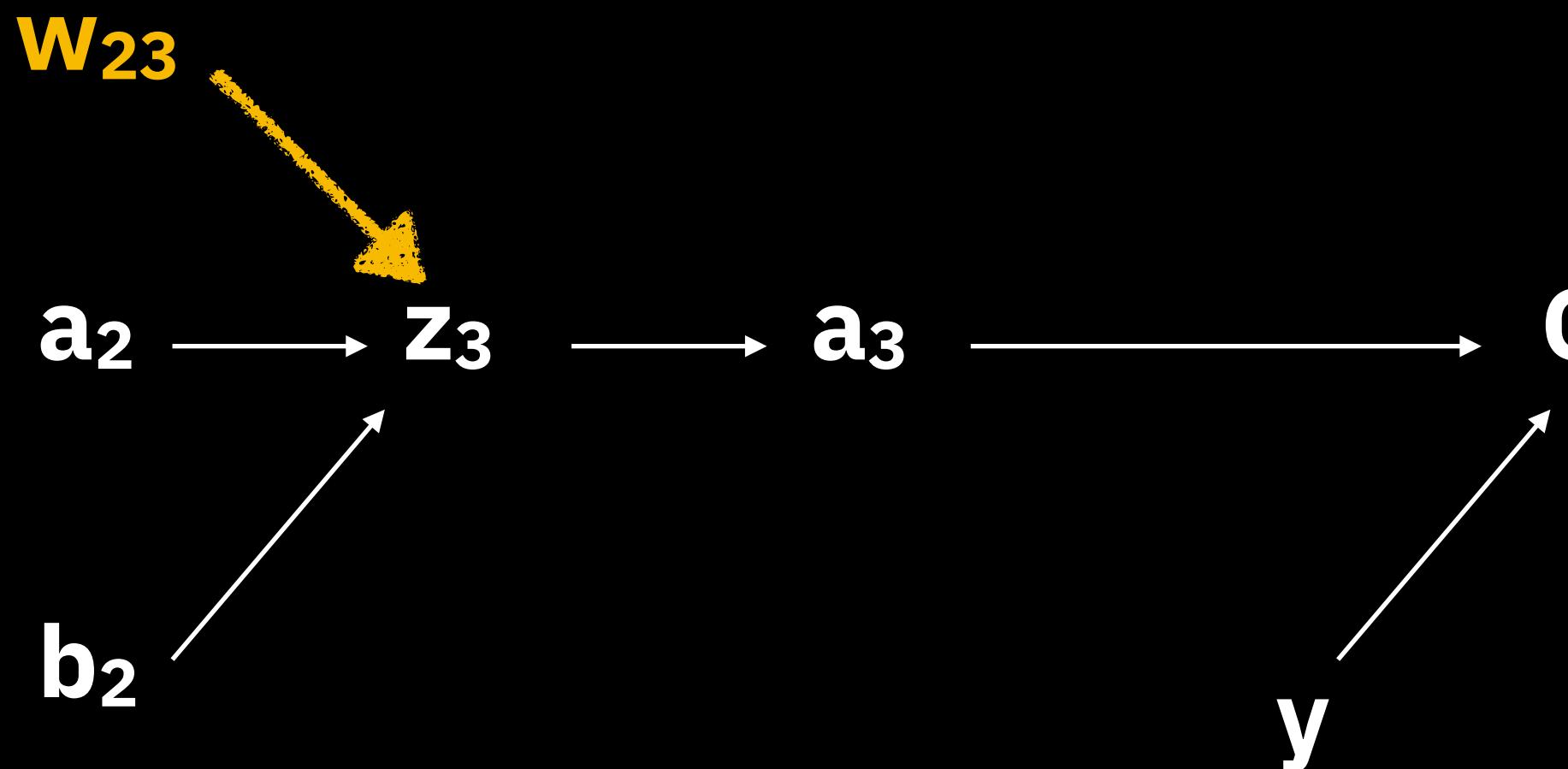
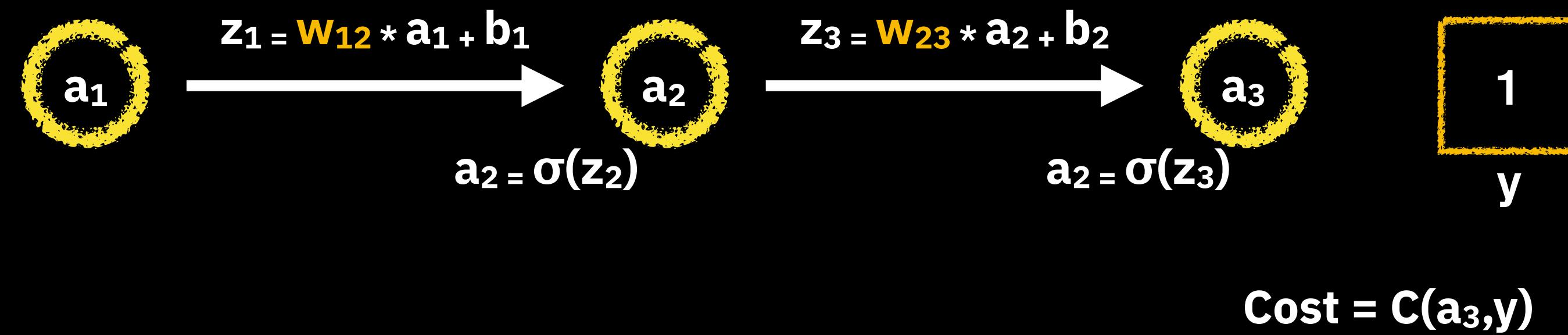
$$\frac{\delta C}{\delta W_{23}} = \frac{\delta C}{\delta a_3}$$

# Computational Graphs (Backpropagation)



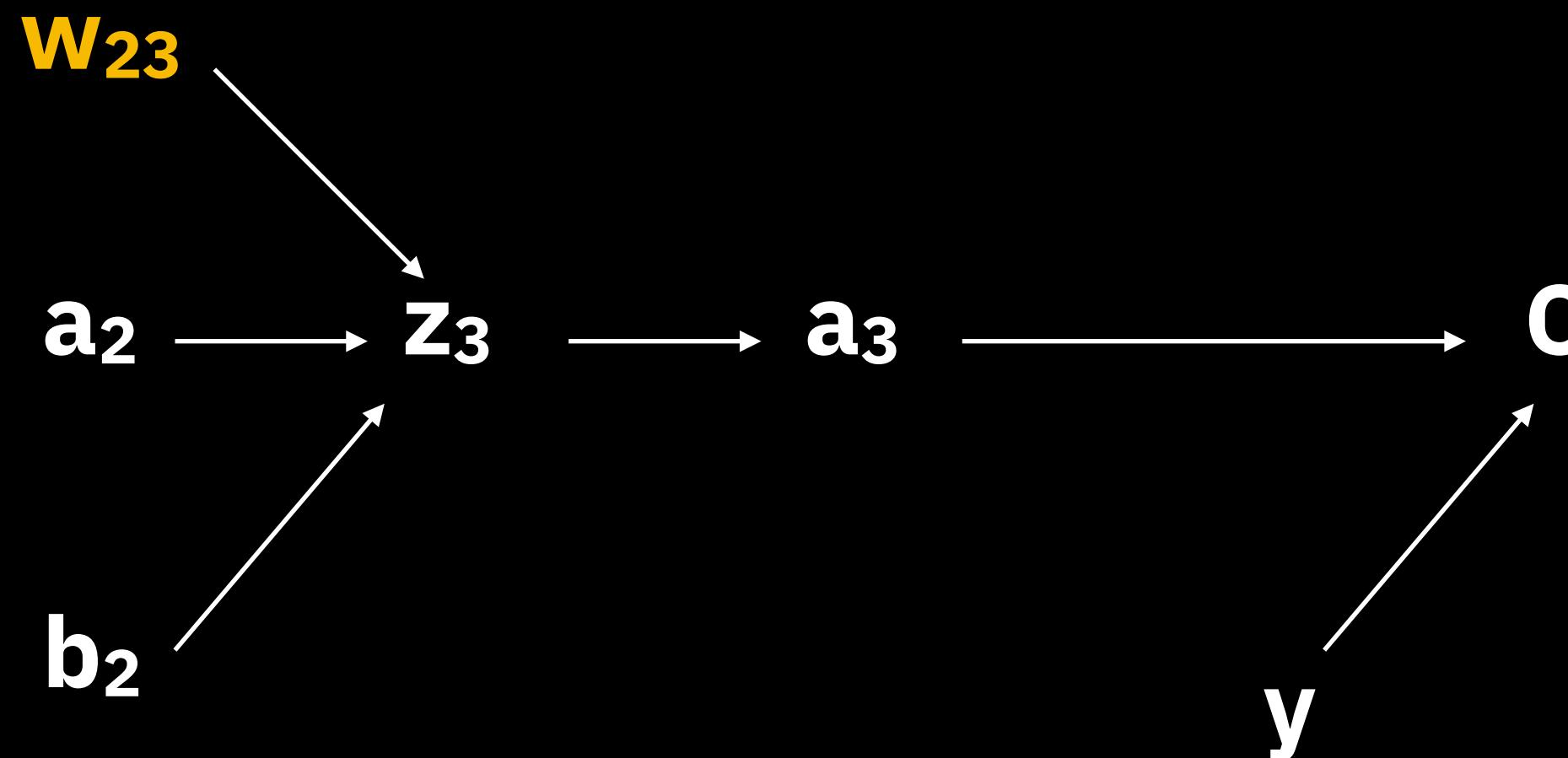
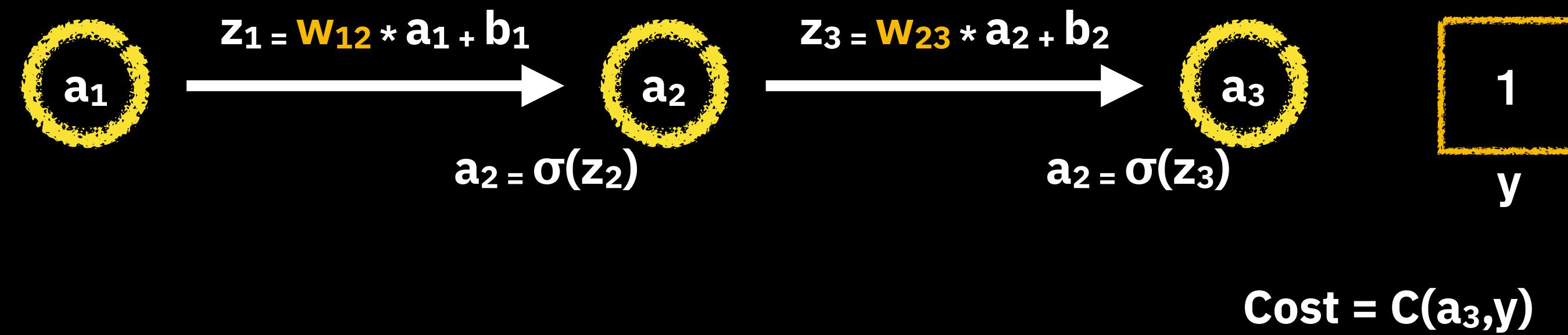
$$\frac{\delta C}{\delta W_{23}} = \frac{\delta C}{\delta a_3} \frac{\delta a_3}{\delta z_3}$$

# Computational Graphs (Backpropagation)



$$\frac{\delta C}{\delta W_{23}} = \frac{\delta C}{\delta a_3} \frac{\delta a_3}{\delta z_3} \frac{\delta z_3}{\delta W_{23}}$$

# Computational Graphs (Backpropagation)



$$\frac{\delta C}{\delta W_{23}} = \frac{\delta C}{\delta a_3} \frac{\delta a_3}{\delta z_3} \frac{\delta z_3}{\delta W_{23}}$$

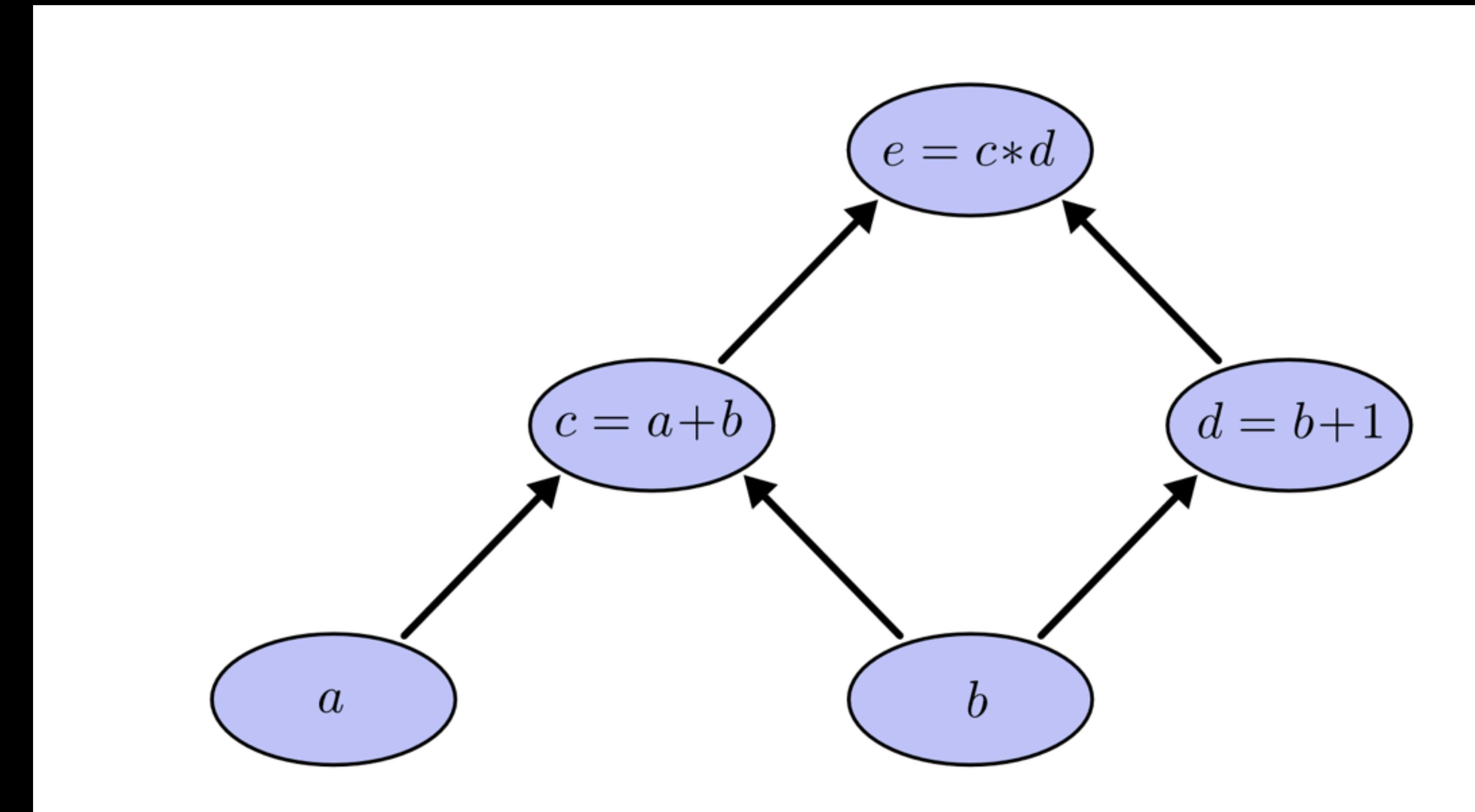
Chain Rule

# Computational Graphs

$$c = a + b$$

$$d = b + 1$$

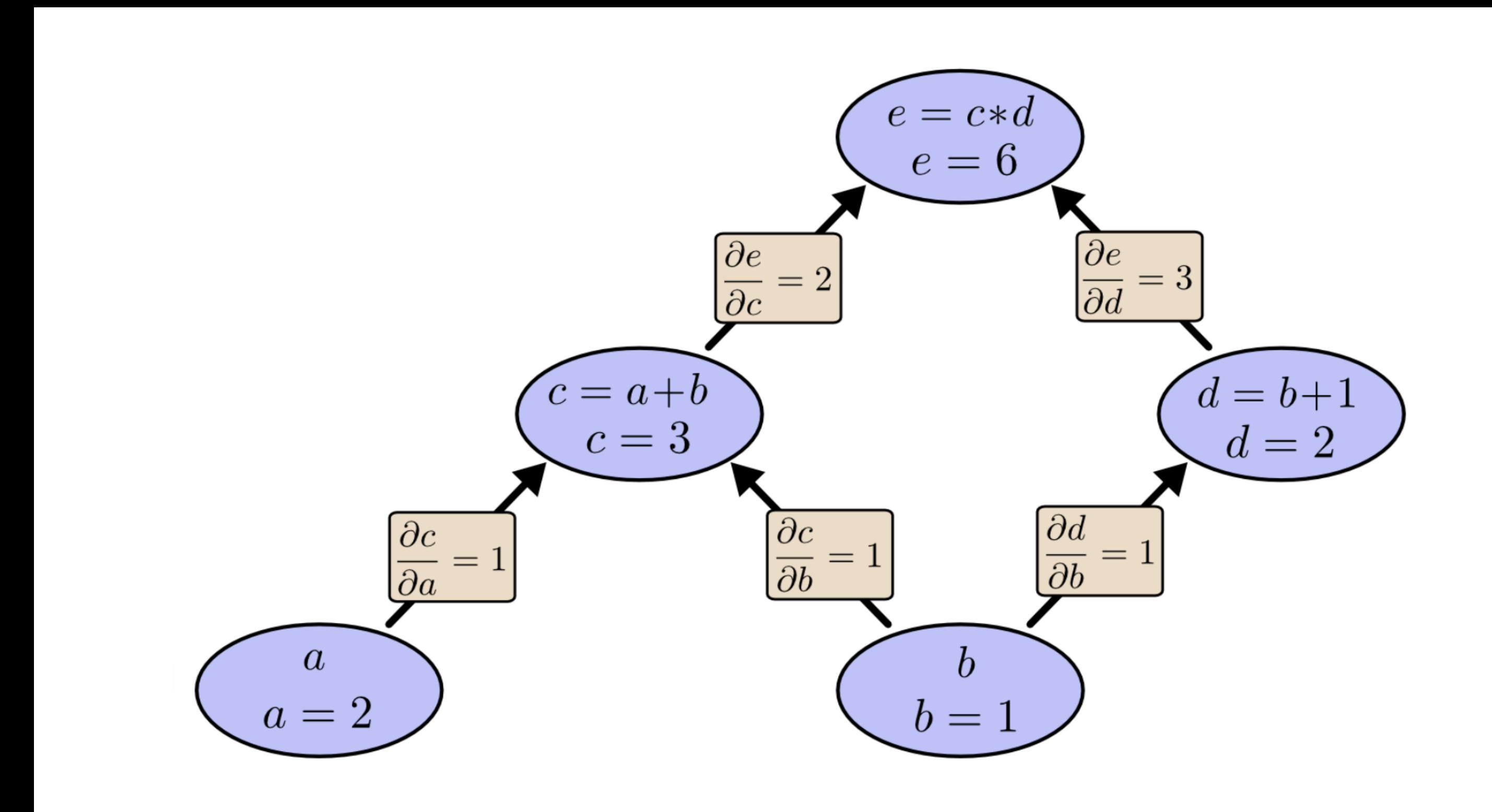
$$e = c * d$$



# Computational Graphs

## Derivatives in Computational Graphs

$$\frac{\partial e}{\partial b} = 1 * 2 + 1 * 3$$





# Any Questions?

email id : harisris@gmail.com