



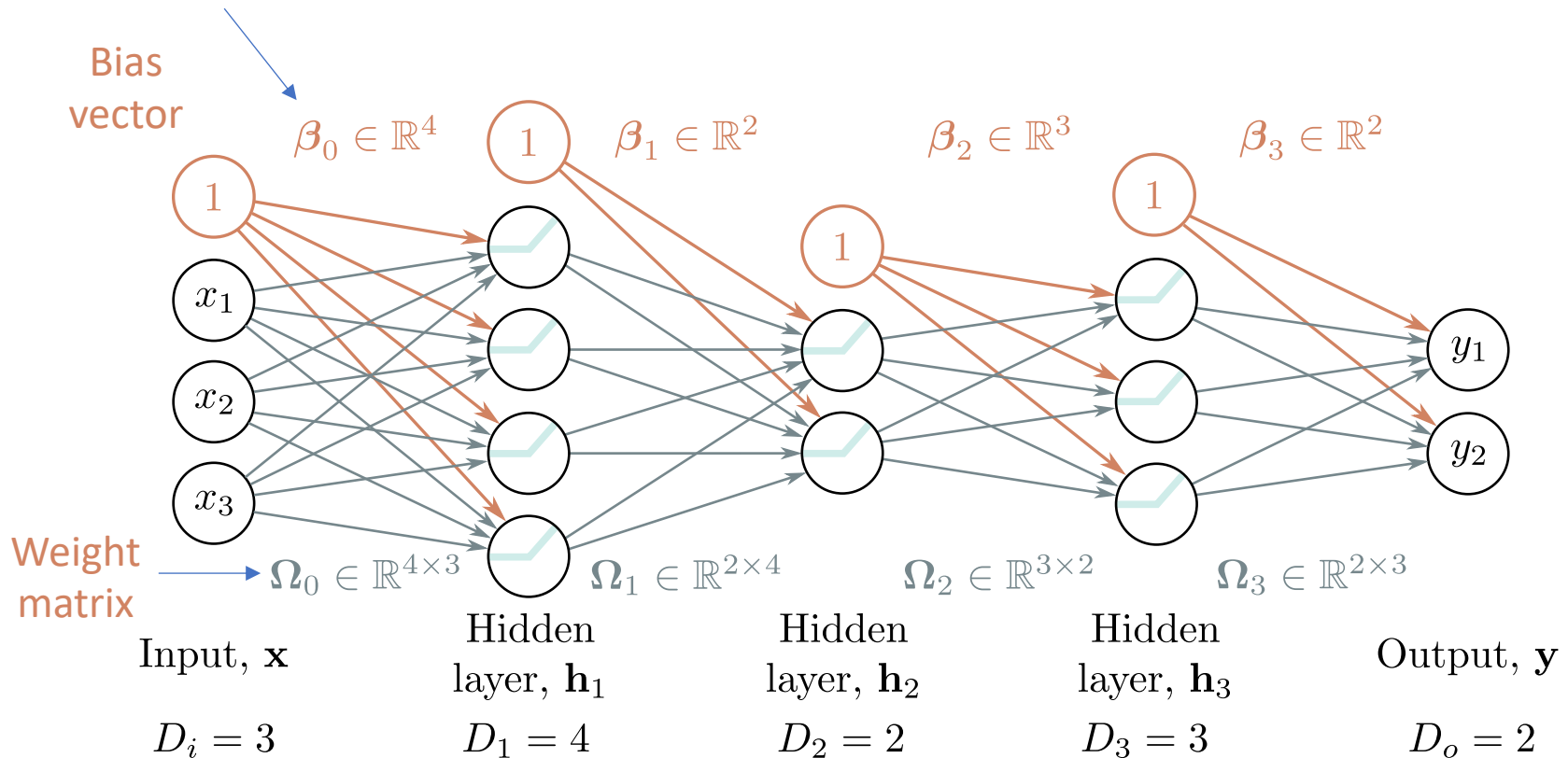
# EECS 230 Deep Learning

## Lecture 4: Back Propagation

Some slides from O. Veksler, Y. Boykov, A. Ng, Y. LeCun, G. Hinton, A. Ranzato, R. Fergus

# How to train neural networks?

□ Training == learning weight and bias



Example of Multi Layer Perceptron (MLP)



# Optimization via Gradient Descent

# Optimization of continuous differentiable functions

- How to minimize a function of a single variable

$$f(x) = (x - 5)^2$$

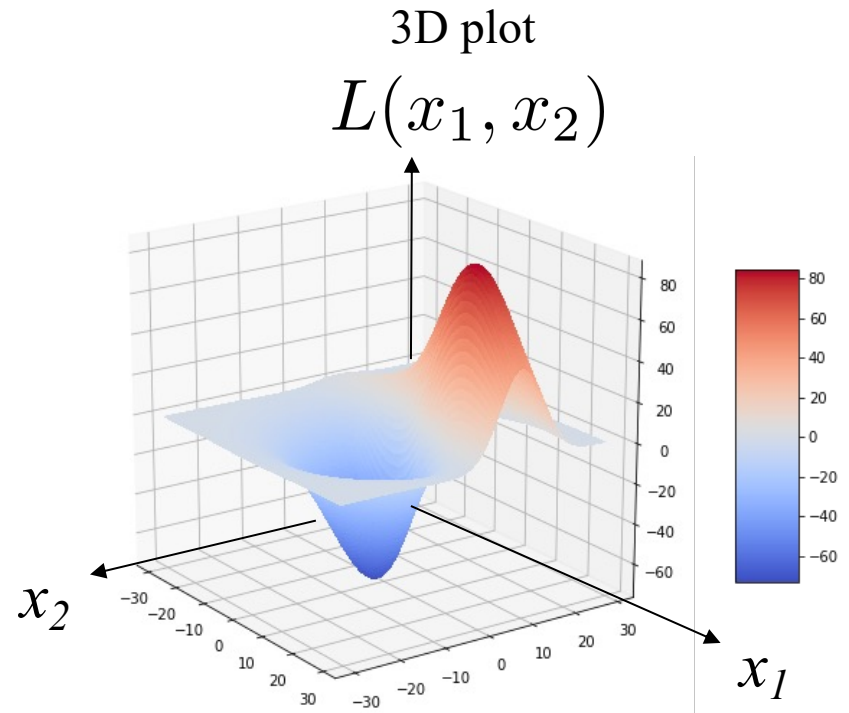
- Take derivative and set it to 0

$$\frac{d}{dx}f(x) = 0$$

- May find a closed form solution

$$\frac{d}{dx}f(x) = 2(x - 5) = 0 \quad \Rightarrow \quad x = 5$$

# Differentiation

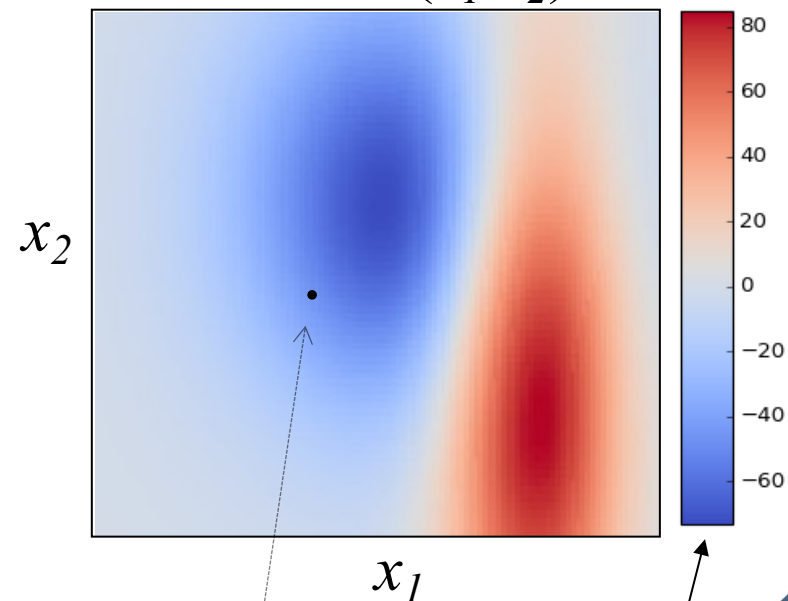


What is “**slope**” of  $L(x_1, x_2)$  at a given point  $\mathbf{x}=(x_1, x_2)$ ?

# Differentiation

“heat-map” visualization of  $L$

domain of  $L(x_1, x_2)$  in  $R^2$



range of  
 $L(x_1, x_2)$

What is “**slope**” of  $L(x_1, x_2)$  at a given point  $\mathbf{x}=(x_1, x_2)$ ?

# Differentiation

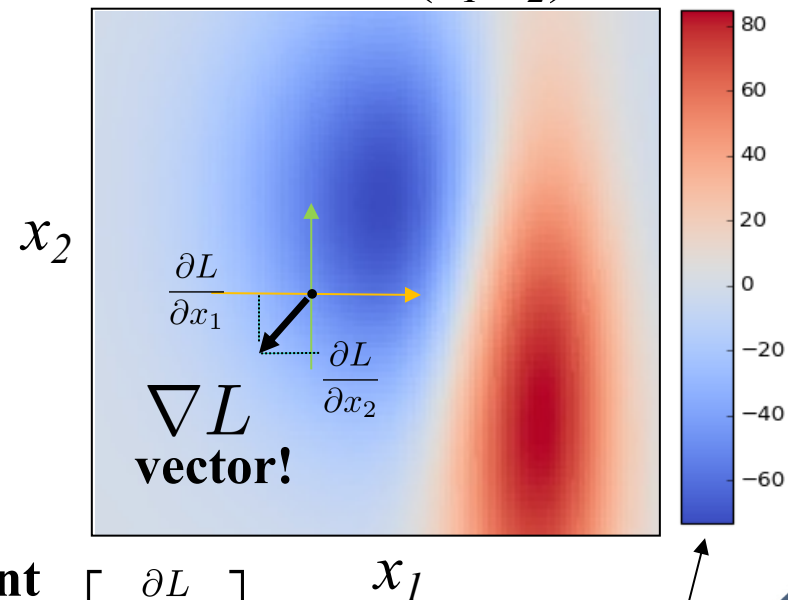
“partial” derivatives

$$\frac{\partial L}{\partial x_1} = \lim_{\epsilon \rightarrow 0} \left( \frac{L(x_1 + \epsilon, x_2) - L(x_1, x_2)}{\epsilon} \right)$$

$$\frac{\partial L}{\partial x_2} = \lim_{\epsilon \rightarrow 0} \left( \frac{L(x_1, x_2 + \epsilon) - L(x_1, x_2)}{\epsilon} \right)$$

“heat-map” visualization of  $L$

domain of  $L(x_1, x_2)$  in  $R^2$



gradient

$$\nabla L = \begin{bmatrix} \frac{\partial L}{\partial x_1} \\ \frac{\partial L}{\partial x_2} \end{bmatrix}$$

**direction of the steepest  
ascent at point  $\mathbf{x}=(x_1, x_2)$**

range of  
 $L(x_1, x_2)$

# Differentiation

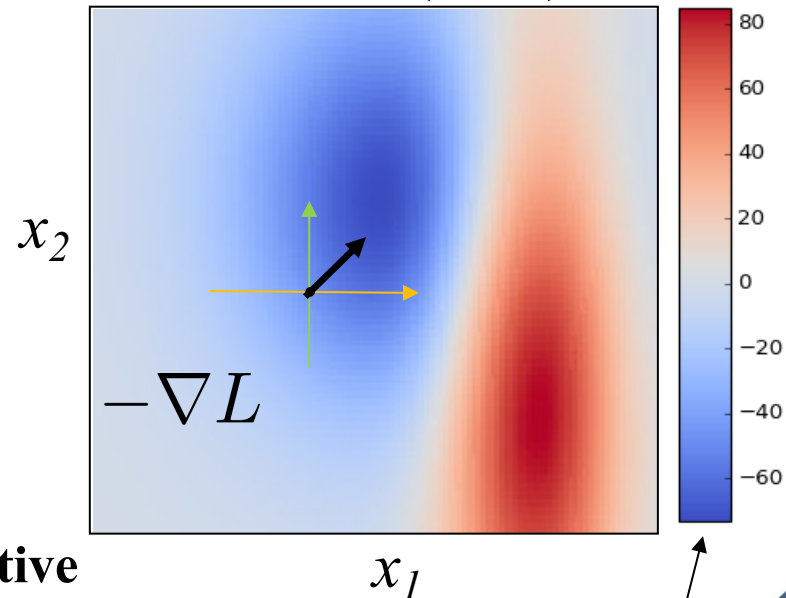
The most common optimization method for continuous differentiable (multi-variate) functions:

## gradient descent

take a step  $\mathbf{x}' = \mathbf{x} - \alpha \nabla L$   
towards lower values  
of the function

“heat-map” visualization of  $L$

domain of  $L(x_1, x_2)$  in  $R^2$



negative  
gradient

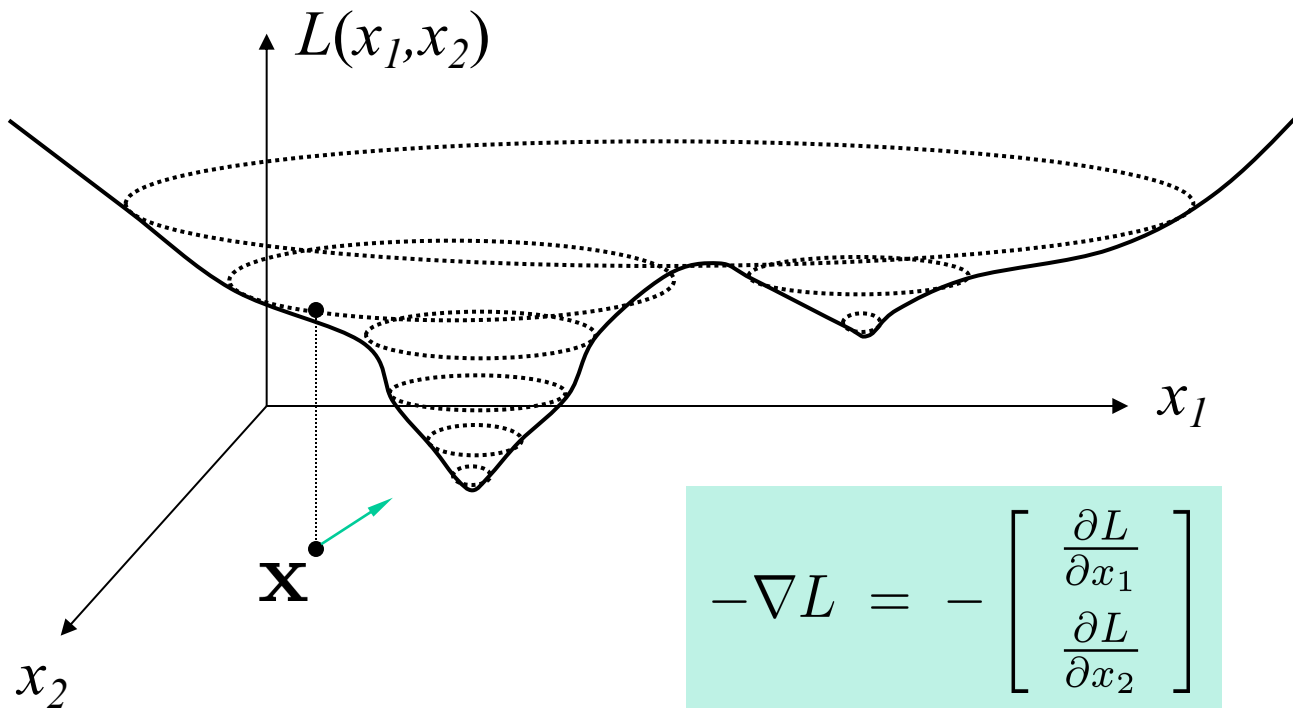
range of  
 $L(x_1, x_2)$

direction of the steepest  
descent at point  $\mathbf{x}=(x_1, x_2)$



# Gradient Descent

Example: for a function of two variables

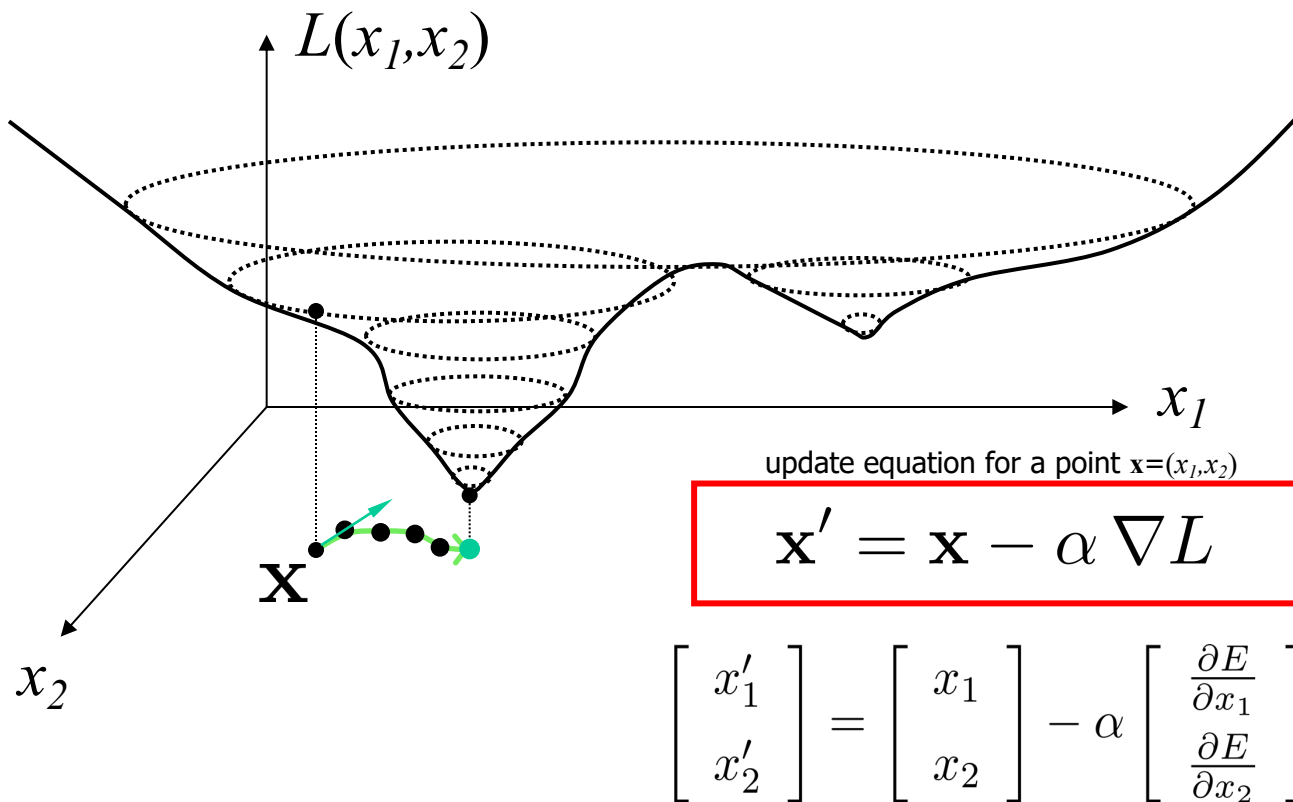


$$-\nabla L = - \begin{bmatrix} \frac{\partial L}{\partial x_1} \\ \frac{\partial L}{\partial x_2} \end{bmatrix}$$

- direction of (negative) **gradient** at point  $\mathbf{x}=(x_1, x_2)$  is direction of the steepest descent towards lower values of function  $L$
- magnitude of gradient at  $\mathbf{x}=(x_1, x_2)$  gives the value of the slope

# Gradient Descent

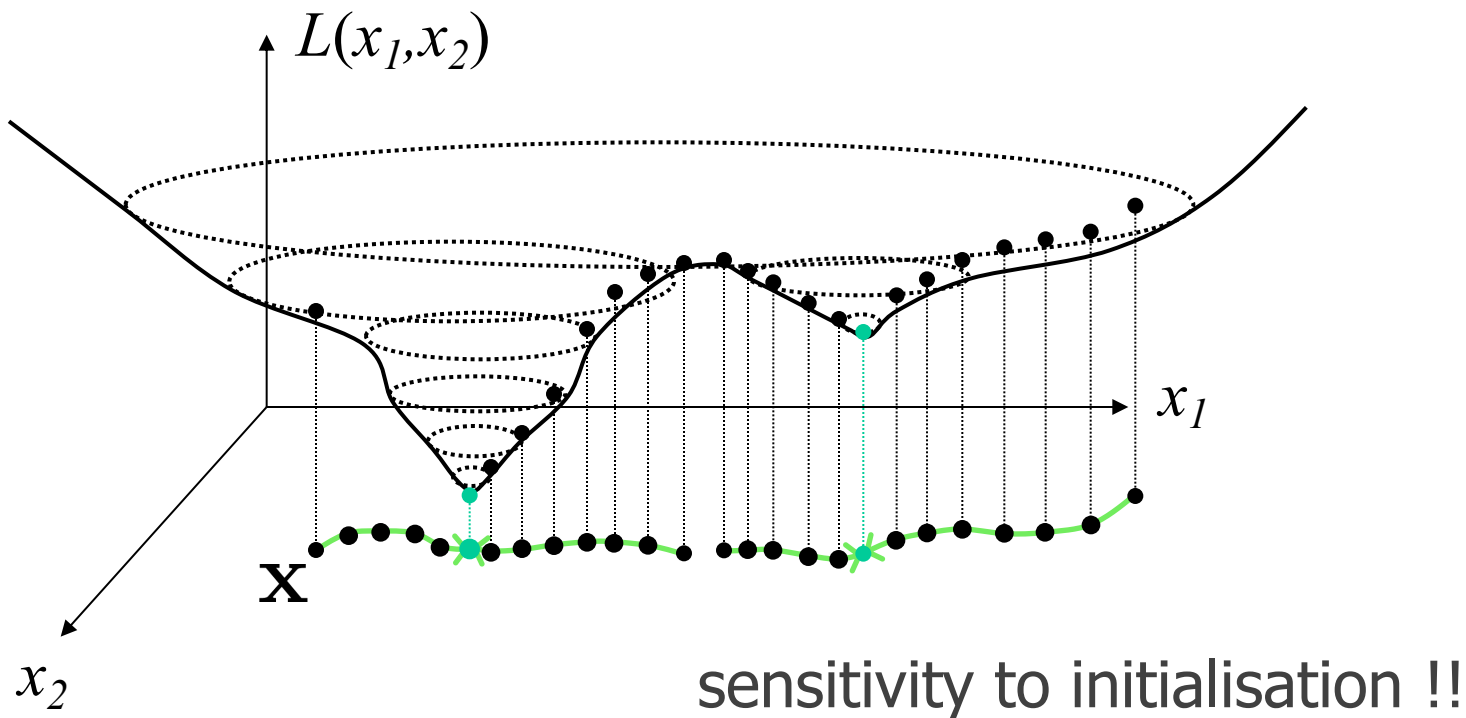
Example: for a function of two variables



Stop at a **local minima** where  $\nabla L = \vec{0}$

# Gradient Descent

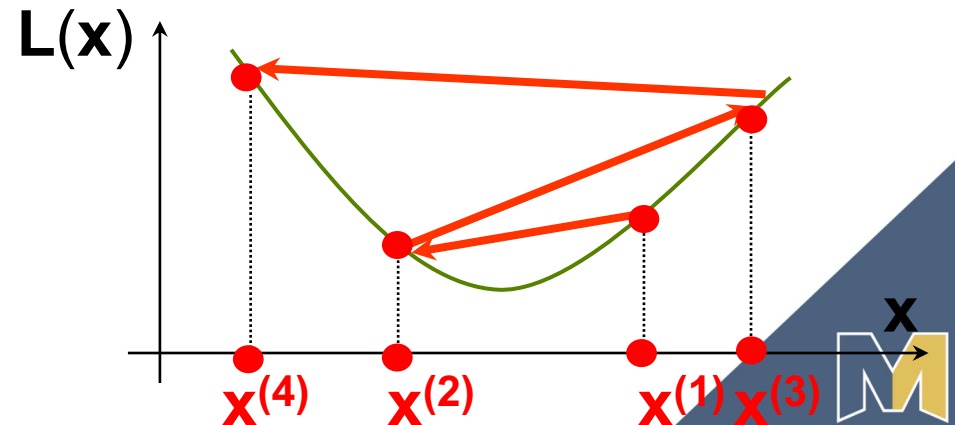
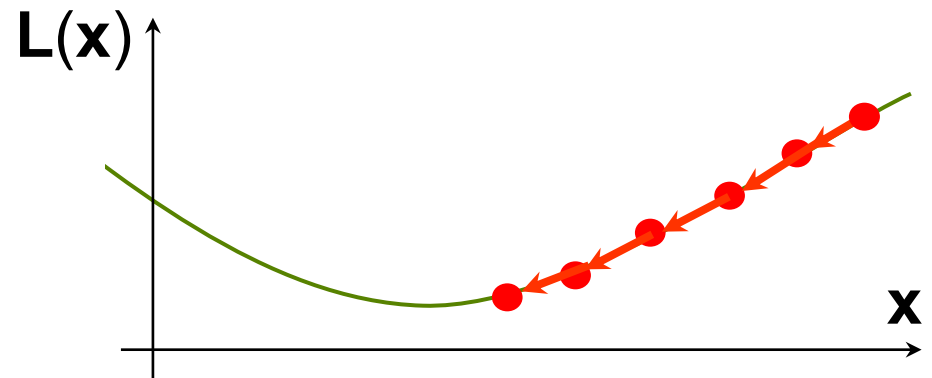
Example: for a function of two variables



# How to Set Learning Rate $\alpha$ ?

$$\mathbf{x}' = \mathbf{x} - \alpha \nabla L$$

- If  $\alpha$  too small, too many iterations to converge
- If  $\alpha$  too large, may overshoot the local minimum and possibly never even converge



# Variable Learning Rate

If desired, can change learning rate  $\alpha$  at each iteration

$k = 1$   
 $\mathbf{x}^{(1)} = \text{any initial guess}$   
choose  $\alpha, \epsilon$   
**while**  $\alpha \|\nabla L(\mathbf{x}^{(k)})\| > \epsilon$   
     $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha \nabla L(\mathbf{x}^{(k)})$   
     $k = k + 1$

fixed  $\alpha$   
gradient descent

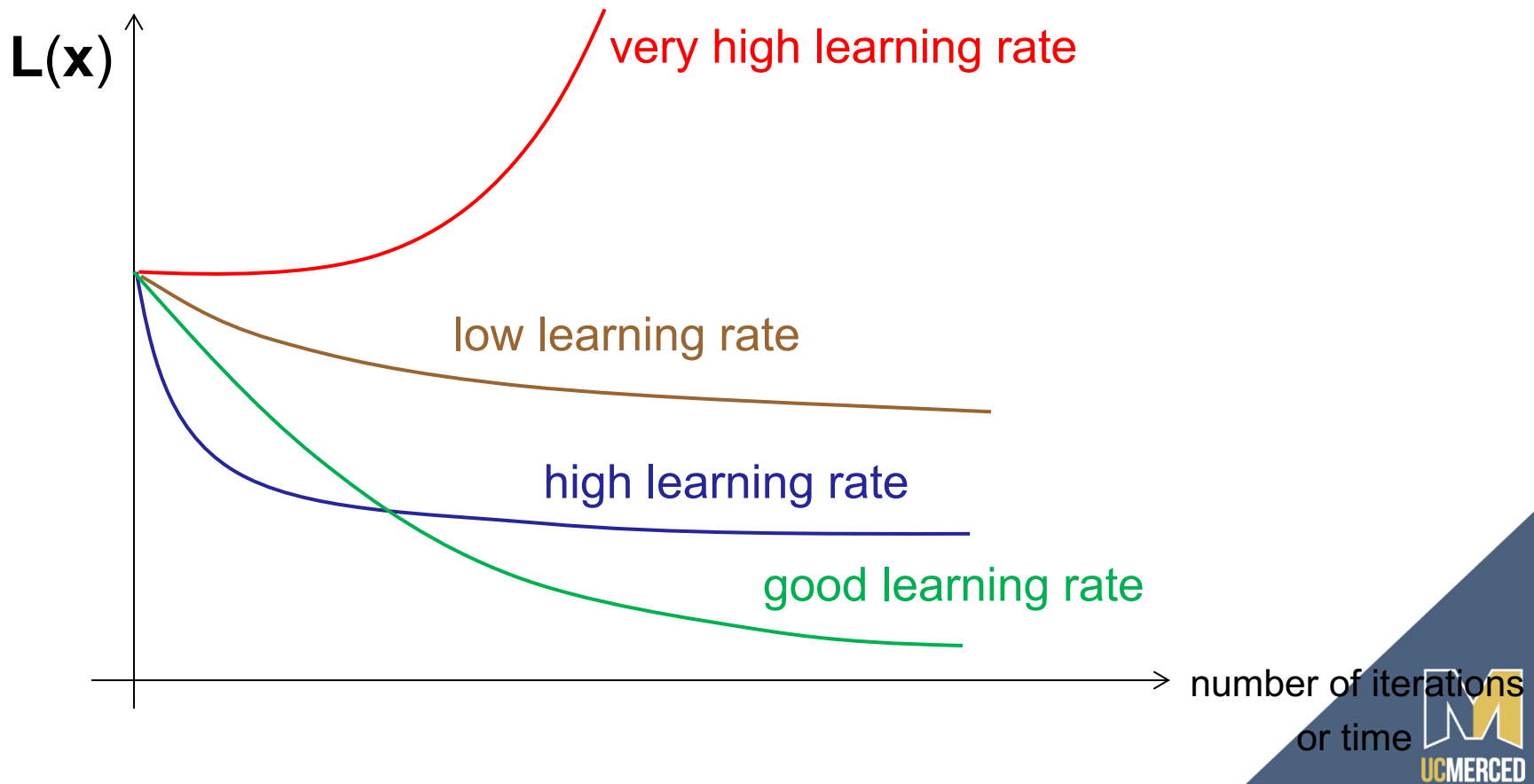


$k = 1$   
 $\mathbf{x}^{(1)} = \text{any initial guess}$   
choose  $\epsilon$   
**while**  $\alpha \|\nabla L(\mathbf{x}^{(k)})\| > \epsilon$   
    **choose**  $\alpha^{(k)}$   
     $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha^{(k)} \nabla L(\mathbf{x}^{(k)})$   
     $k = k + 1$

variable  $\alpha$   
gradient descent

# Learning Rate

- Monitor learning rate by looking at how fast the objective function decreases

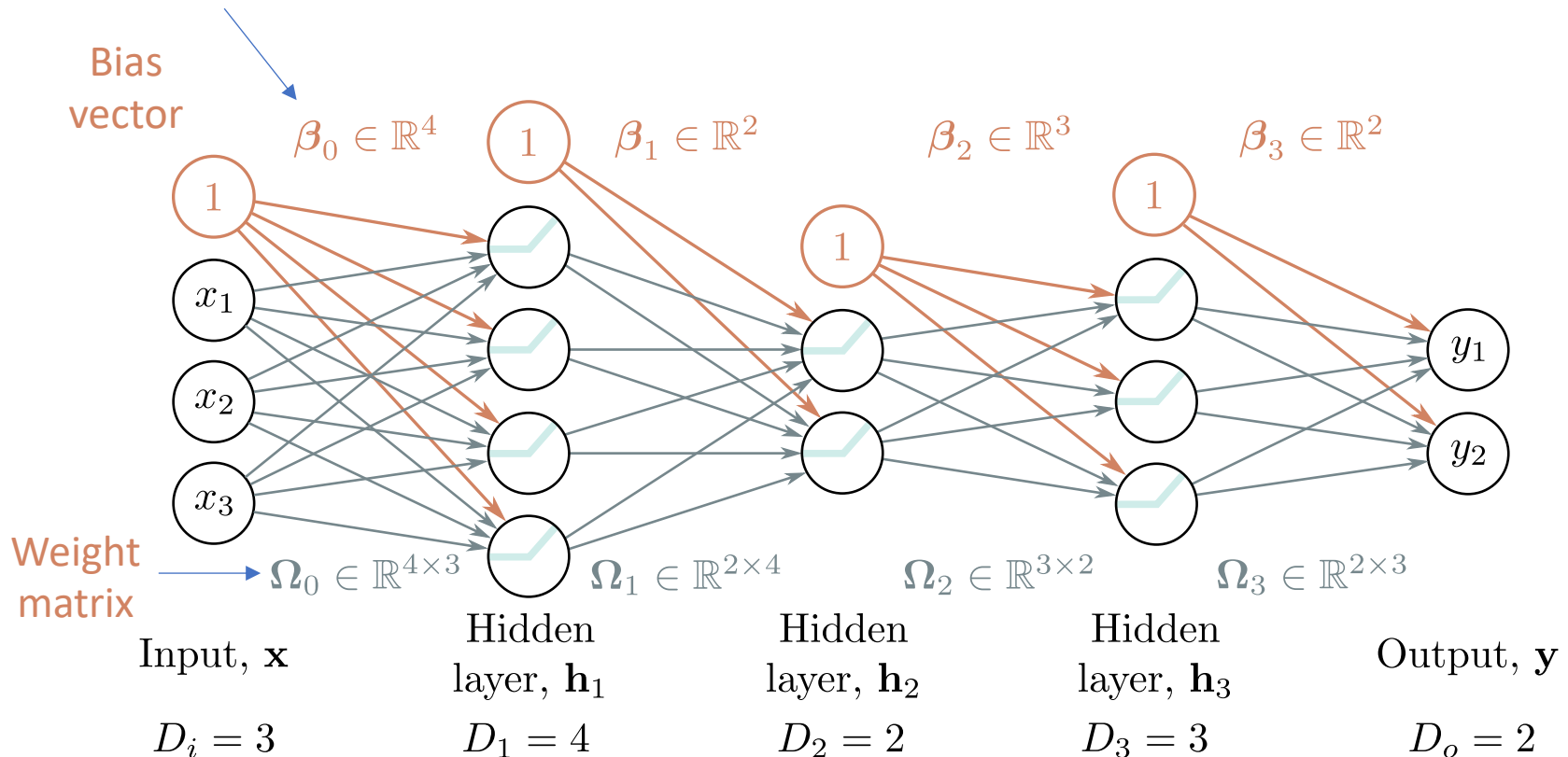




# Derivative and Back Propagation

# How to take derivatives w.r.t. weights?

□ Training == learning weight and bias

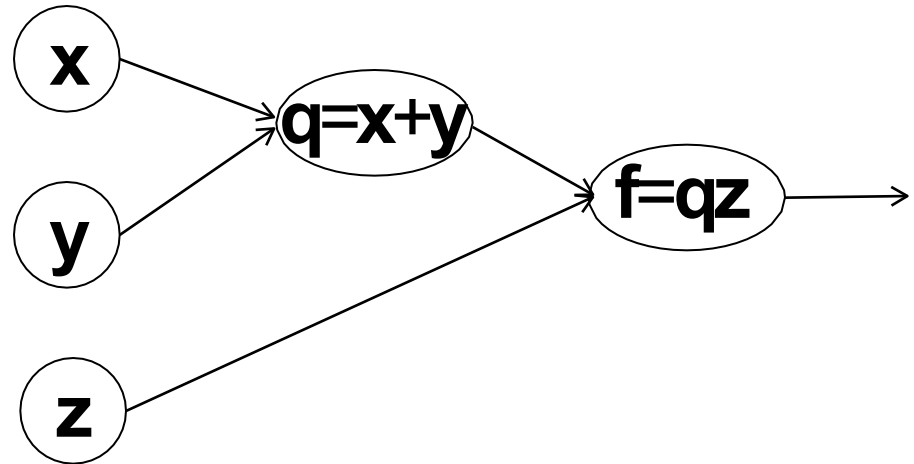


Example of Multi Layer Perceptron (MLP)



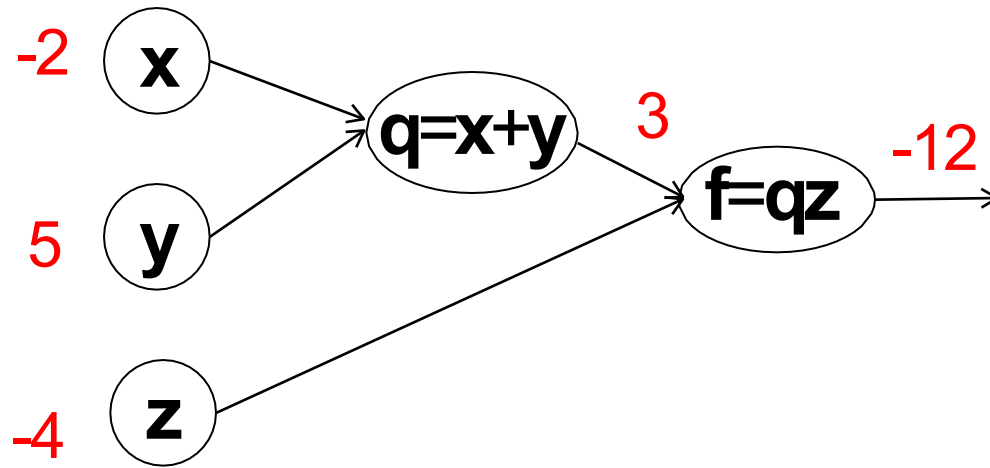
# Computing Derivatives: Small Example

- Small network  $f(x,y,z) = (x+y)z$
- Rewrite using
  - $q = x + y$
- $f(x,y,z) = qz$
- each node does one operation



# Computing Derivatives: Small Example

- Small network  $f(\mathbf{x}, \mathbf{y}, \mathbf{z}) = (\mathbf{x} + \mathbf{y})\mathbf{z}$
- Rewrite using
  - $\mathbf{q} = \mathbf{x} + \mathbf{y}$
  - $f(\mathbf{x}, \mathbf{y}, \mathbf{z}) = \mathbf{q}\mathbf{z}$
- Example of computing  $f(-2, 5, -4)$

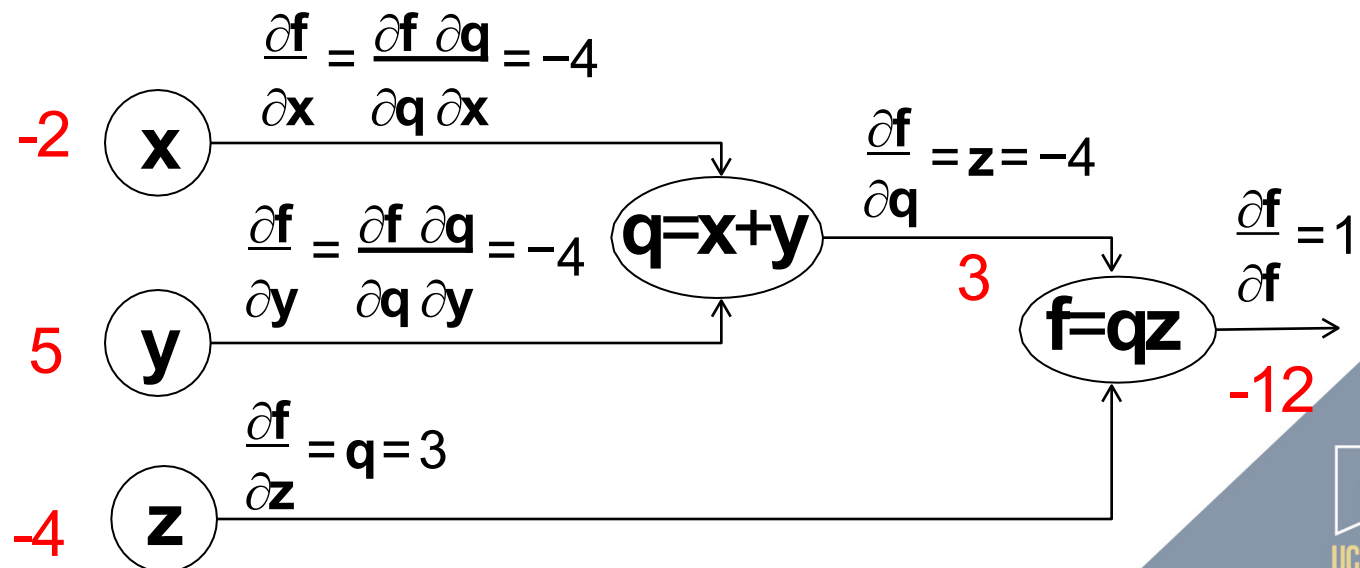


# Computing Derivatives: Small Example

- Small network  $f(\mathbf{x}, \mathbf{y}, \mathbf{z}) = (\mathbf{x} + \mathbf{y})\mathbf{z}$
- Rewrite using  $\mathbf{q} = \mathbf{x} + \mathbf{y} \Rightarrow f(\mathbf{x}, \mathbf{y}, \mathbf{z}) = \mathbf{q}\mathbf{z}$
- Want  $\frac{\partial f}{\partial \mathbf{x}}, \frac{\partial f}{\partial \mathbf{y}}, \frac{\partial f}{\partial \mathbf{z}}$
- Compute  $\frac{\partial f}{\partial}$  from the end backwards
  - for each edge, with respect to the main variable at edge origin
  - using chain rule with respect to the variable at edge end, if needed

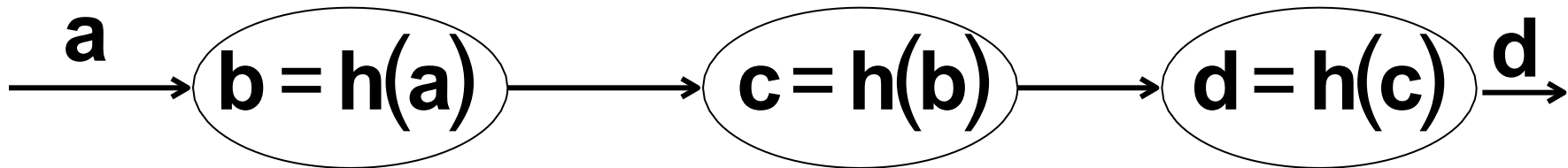
chain rule for  $f(\mathbf{y}(\mathbf{x}))$

$$\frac{\partial f}{\partial \mathbf{x}} = \frac{\partial f}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}}$$



# Computing Derivatives: Chain of Chain Rule

- Compute  $\frac{\partial d}{\partial a}$  from the end backwards
  - for each edge, with respect to the main variable at edge origin
  - using chain rule with respect to the variable at edge end, if needed



$$\frac{\partial d}{\partial a} = \frac{\partial d}{\partial b} \frac{\partial b}{\partial a}$$

prev local

$$\frac{\partial d}{\partial b} = \frac{\partial d}{\partial c} \frac{\partial c}{\partial b}$$

prev local

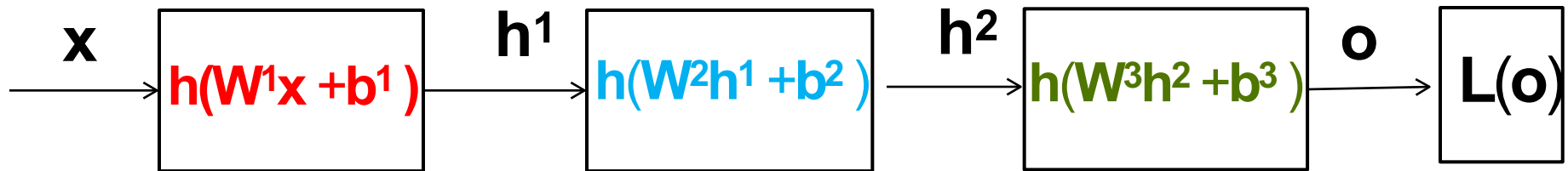
$$\frac{\partial d}{\partial c}$$

local



example: if  $h(c) = c^2$ , then  $\frac{\partial d}{\partial c} = \frac{\partial h}{\partial c} = 2c$

# Computing Derivatives Backwards

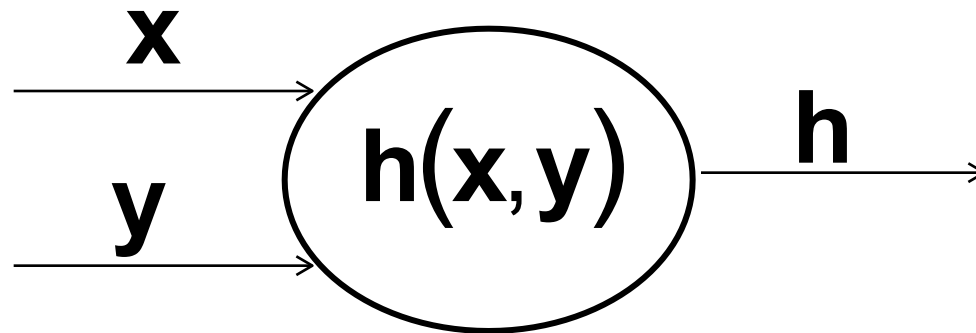


← direction of computation

- Have loss function  $L(o)$
- Need derivatives for all  $\frac{\partial L}{\partial w}$ ,  $\frac{\partial L}{\partial b}$
- Will compute derivatives from end to front, backwards
- On the way will also compute intermediate derivatives  $\frac{\partial L}{\partial h}$

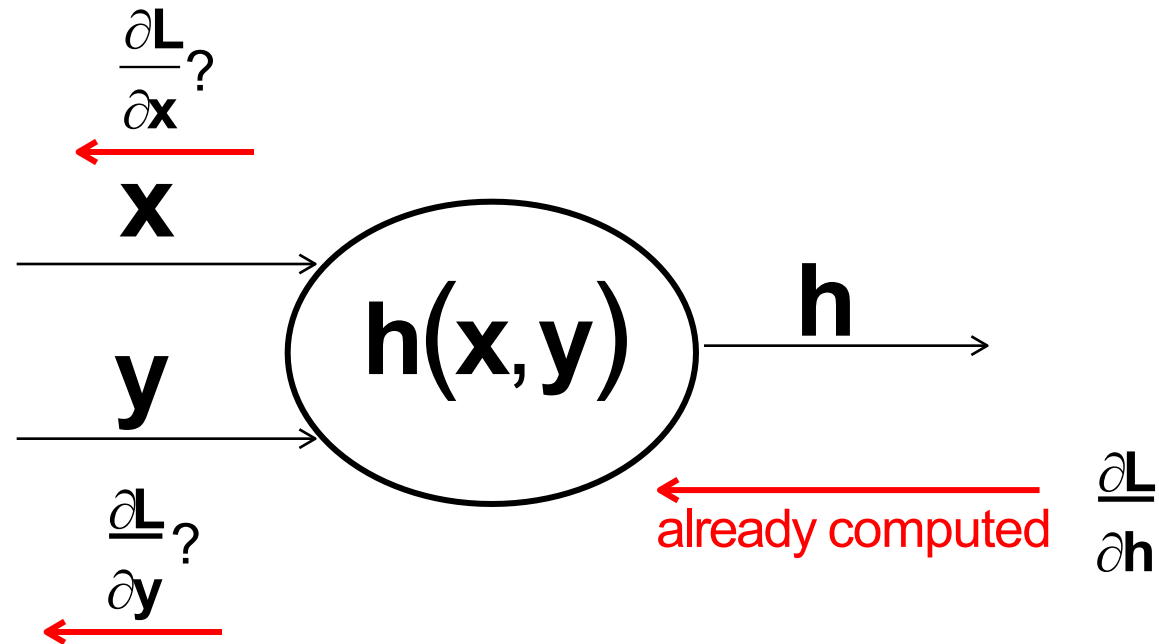
# Computing Derivatives: Look at One Node

- Simplified view at a network node
  - inputs  $\mathbf{x}, \mathbf{y}$  come in
  - node computes some function  $\mathbf{h}(\mathbf{x}, \mathbf{y})$



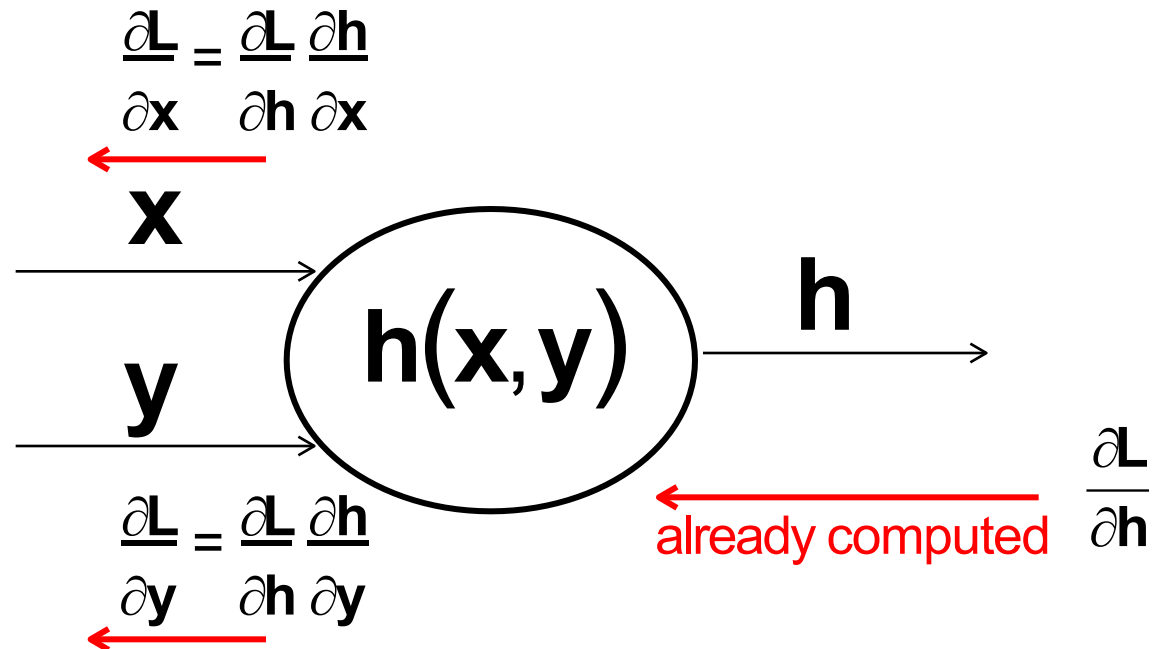
# Computing Derivatives: Look at One Node

- At each network node
  - inputs  $\mathbf{x}, \mathbf{y}$  come in
  - nodes computes activation function  $\mathbf{h}(\mathbf{x}, \mathbf{y})$
- Have loss function  $\mathbf{L}(\cdot)$



# Computing Derivatives: Look at One Node

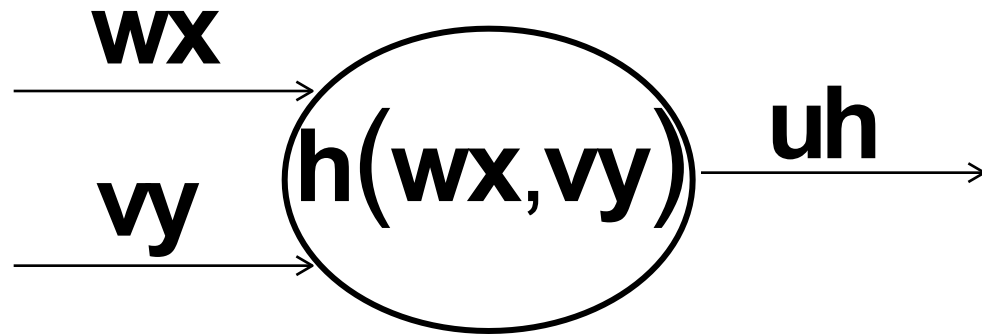
- Need  $\frac{\partial \mathbf{L}}{\partial \mathbf{x}}, \frac{\partial \mathbf{L}}{\partial \mathbf{y}}$
- Easy to compute local node derivatives  $\frac{\partial \mathbf{h}}{\partial \mathbf{x}}, \frac{\partial \mathbf{h}}{\partial \mathbf{y}}$



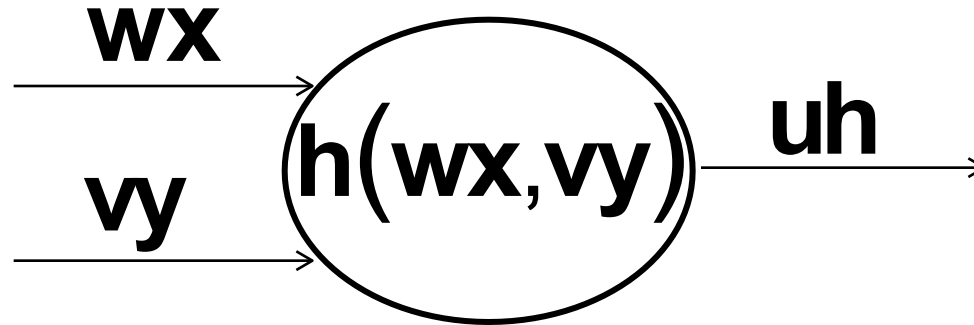


# Computing Derivatives: Look at One Node

- More complete view at a network node
  - inputs  $\mathbf{x}, \mathbf{y}$  come in, get multiplied by weight  $\mathbf{w}$  and  $\mathbf{v}$
  - node computes function  $\mathbf{h}(\mathbf{wx}, \mathbf{vy})$
  - node output  $\mathbf{h}$  gets multiplied by  $\mathbf{u}$

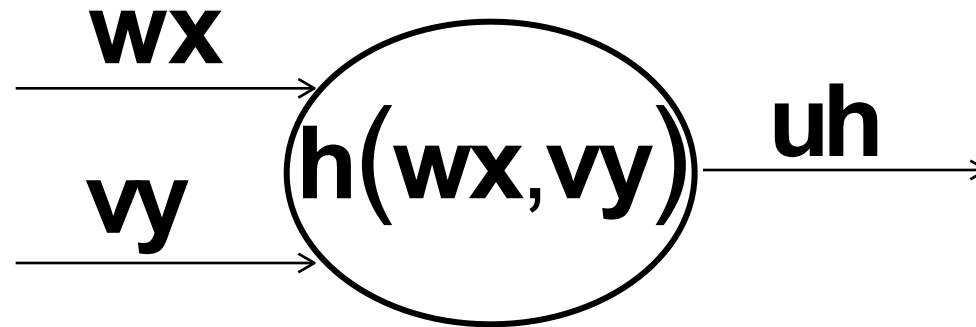


# Computing Derivatives: Look at One Node

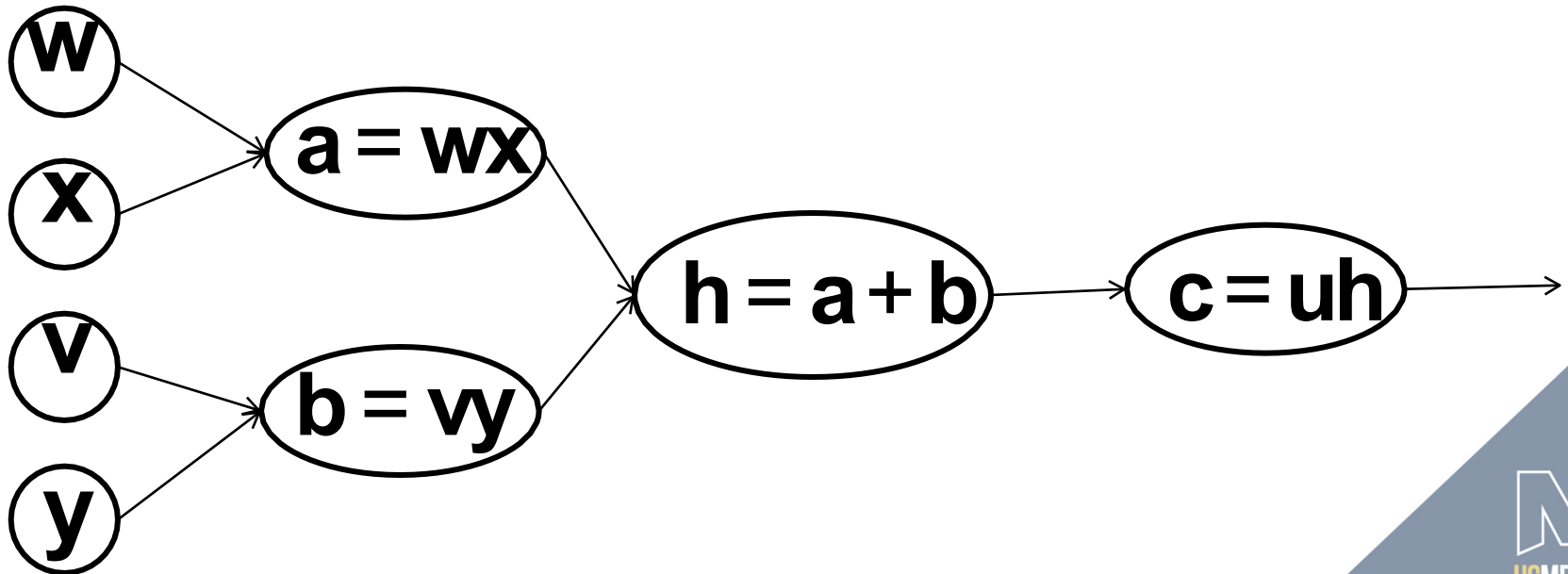


- To be concrete, let  $h(i,j) = i + j$

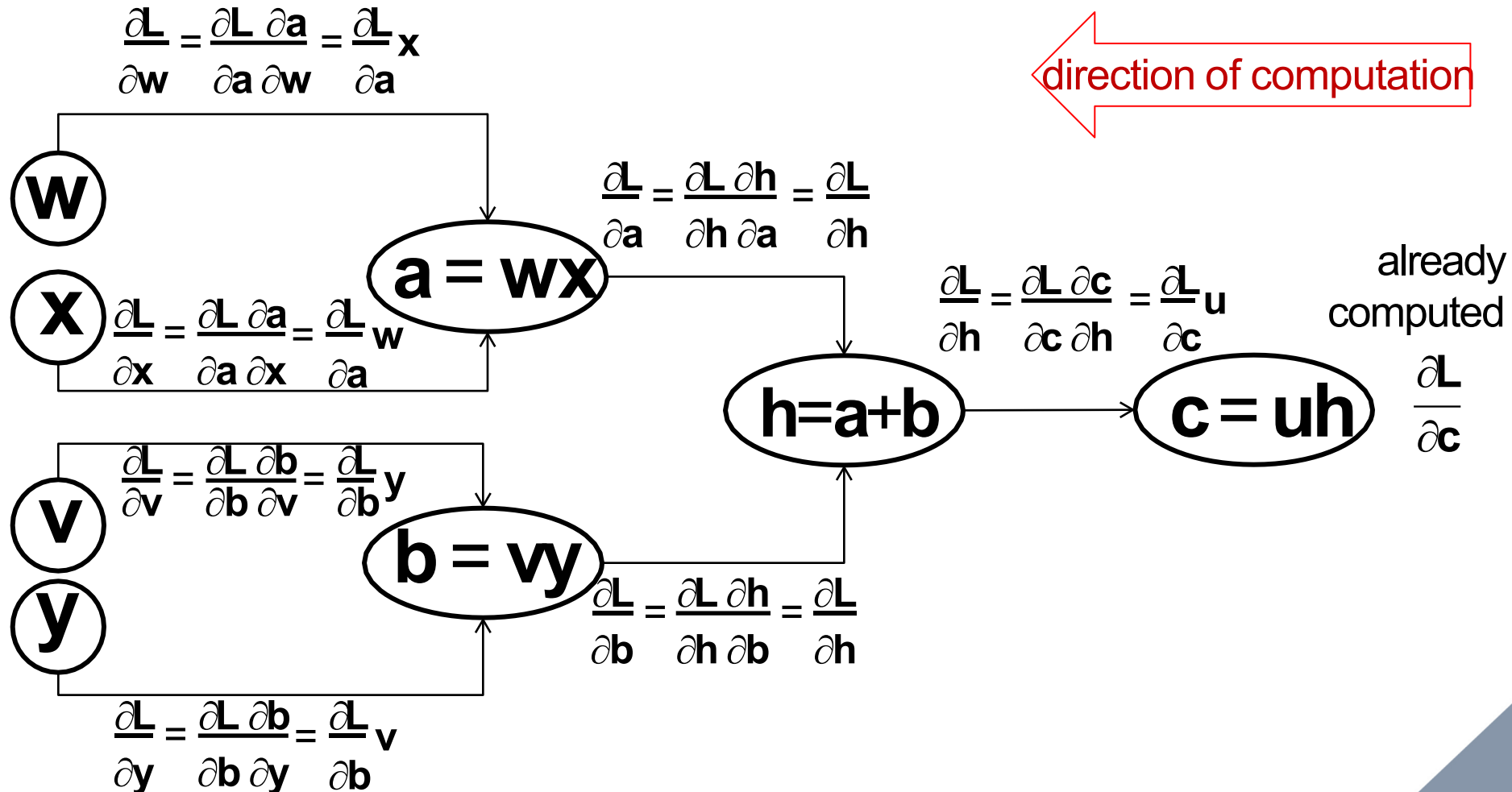
# Computing Derivatives: Look at One Node



- $h(i, j) = i + j$
- Break into more computational nodes
  - all computation happens inside nodes, not on edges

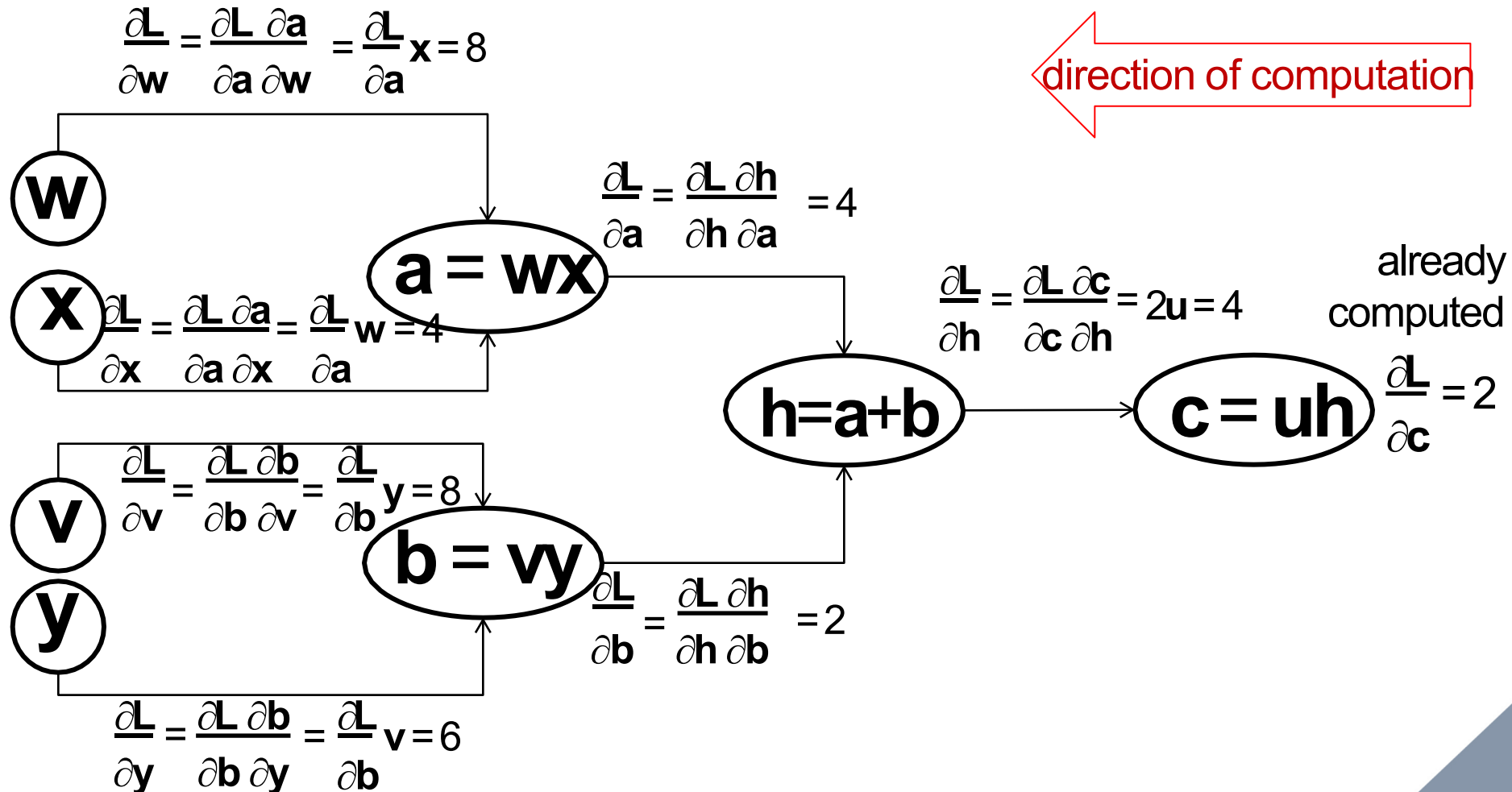


# Computing Derivatives: Look at One Node



- Some of these partial derivatives are intermediate
  - their values will not be used for gradient descent

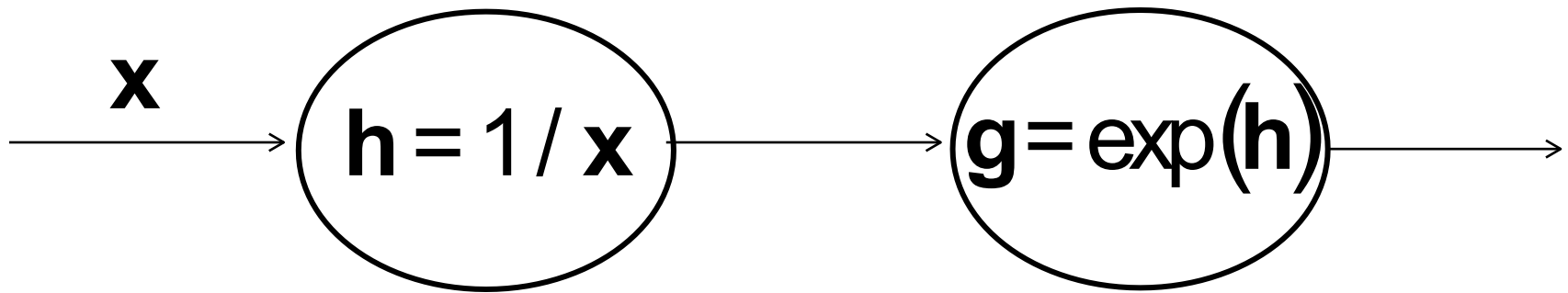
# Computing Derivatives: Look at One Node



- Example when  $\mathbf{w} = 1$ ,  $\mathbf{x} = 2$ ,  $\mathbf{v} = 3$ ,  $\mathbf{y} = 4$ ,  $\mathbf{u} = 2$ ,  $\frac{\partial \mathbf{L}}{\partial \mathbf{c}} = 2$

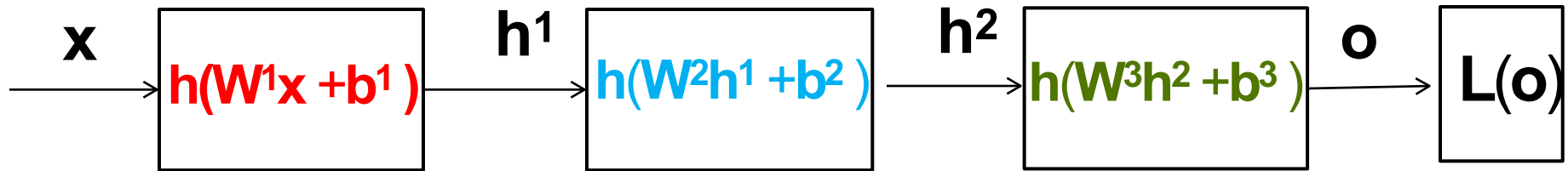
# Computing Derivatives: Staging Computation

- Each node is responsible for one function
- To compute  $\exp(1/\mathbf{x})$



# Computing Derivatives: Vector Notation

- Inputs outputs are often vectors



- $h(a)$  is a function from  $\mathbf{R}^n$  to  $\mathbf{R}^m$
- Chain rule generalizes to vector functions

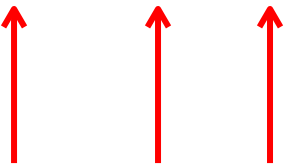
# Computing Derivatives: Vector Notation

- Let  $\mathbf{f}(\mathbf{x}): \mathbf{R}^n \rightarrow \mathbf{R}^m$ ,
  - $\mathbf{x}$  is  $n$ -dimensional vector and output  $\mathbf{f}(\mathbf{x})$  is  $m$ -dimensional vector
- Jacobian matrix
  - has  $m$  rows and  $n$  columns
  - has  $\frac{\partial \mathbf{f}_i}{\partial \mathbf{x}_j}$  in row  $i$ , column  $j$



# Computing Derivatives: Vector Notation

- $\mathbf{f}(\mathbf{x}): \mathbf{R}^n \rightarrow \mathbf{R}^m$  and  $\mathbf{g}(\mathbf{x}): \mathbf{R}^k \rightarrow \mathbf{R}^n$
- $\mathbf{f}(\mathbf{g}(\mathbf{x})): \mathbf{R}^k \rightarrow \mathbf{R}^m$
- Chain rule for vector functions

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \frac{\partial \mathbf{f}}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial \mathbf{x}}$$


Jacobian matrices

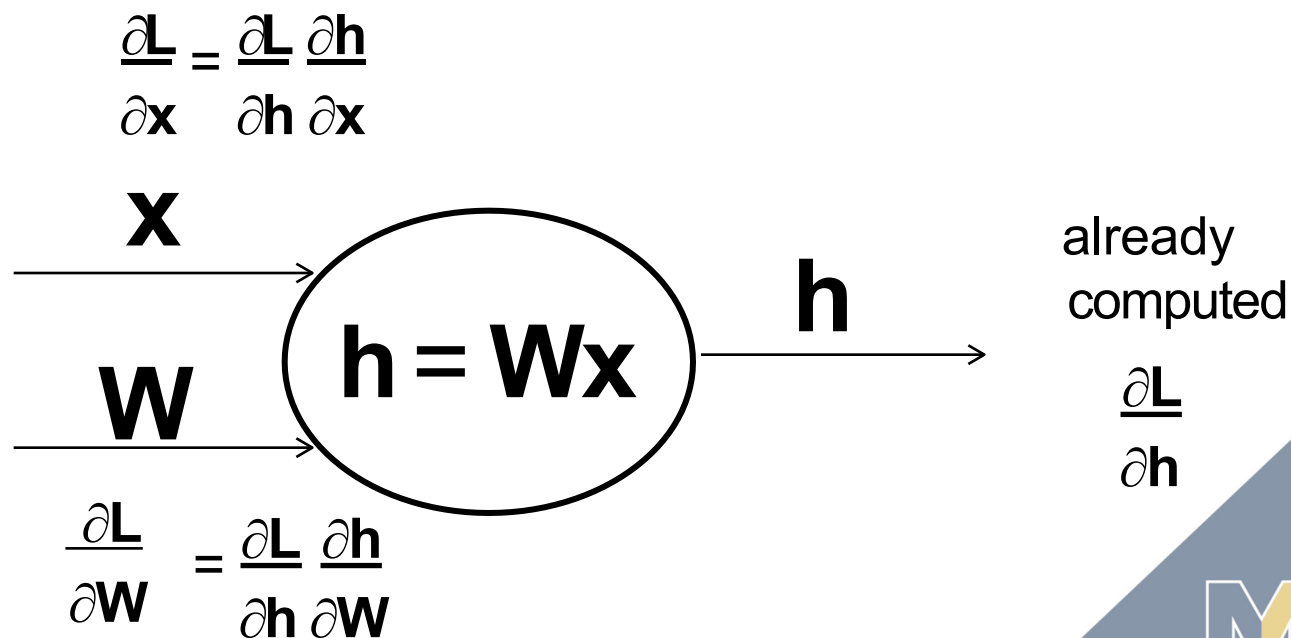
# Vector Notation: Look at One Node

- $\mathbf{h}$ ,  $\mathbf{x}$ ,  $\mathbf{y}$  are vectors
- already computed Jacobian  $\frac{\partial \mathbf{L}}{\partial \mathbf{h}}$
- Need Jacobians  $\frac{\partial \mathbf{L}}{\partial \mathbf{x}}$ ,  $\frac{\partial \mathbf{L}}{\partial \mathbf{y}}$

- Easy to compute local node Jacobians  $\frac{\partial \mathbf{h}}{\partial \mathbf{x}}$ ,  $\frac{\partial \mathbf{h}}{\partial \mathbf{y}}$
- $\frac{\partial \mathbf{L}}{\partial \mathbf{x}} = \frac{\partial \mathbf{L}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{x}}$  ← Jacobian matrices
- $\frac{\partial \mathbf{L}}{\partial \mathbf{y}} = \frac{\partial \mathbf{L}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{y}}$
- 
- The diagram illustrates a node  $h(\mathbf{x}, \mathbf{y})$  represented by an oval. Two input vectors,  $\mathbf{x}$  and  $\mathbf{y}$ , are shown as arrows pointing into the oval from the left. An output vector  $\mathbf{h}$  is shown as an arrow pointing out of the oval to the right. A red triangle is drawn with its vertices at the input/output Jacobian expressions and the node itself. Red arrows point from the text 'Jacobian matrices' to the two vertices of this triangle. The expression  $\frac{\partial \mathbf{L}}{\partial \mathbf{h}}$  is shown to the right of the node, with a red arrow pointing from the text 'already computed' to it.

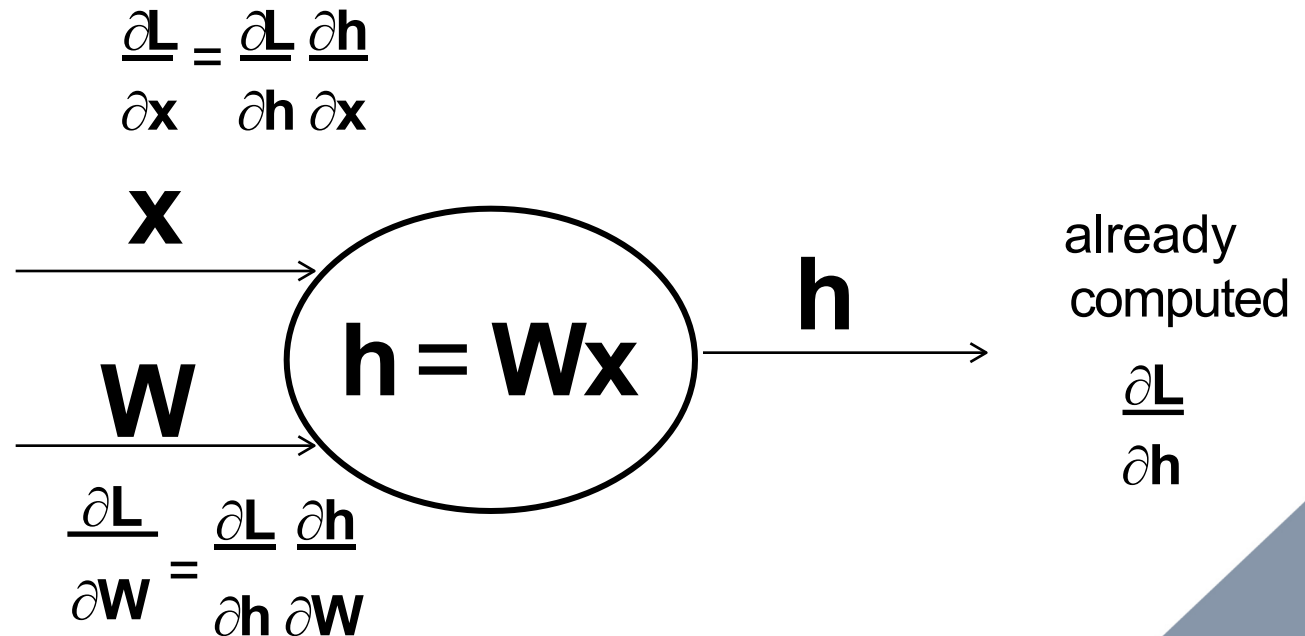
# Vector Notation: Look at One Node

- Can apply to matrices (and tensors) as well
- But first vectorize matrix (or tensor)
- Say  $\mathbf{W}$  is 10 x 5, stretch into 50x1 vector
- Still denote Jacobian by  $\frac{\partial \mathbf{h}}{\partial \mathbf{W}}$



# Vector Notation: Look at One Node

- Easy to compute local node Jacobians  $\frac{\partial \mathbf{h}}{\partial \mathbf{x}}$ ,  $\frac{\partial \mathbf{h}}{\partial \mathbf{W}}$
- But they can get very large (although sparse)
- Say  $\mathbf{h}$  is 1000 x 1,  $\mathbf{W}$  is 1000 x 500, then  $\frac{\partial \mathbf{h}}{\partial \mathbf{W}}$  is 1000 x 500,000



# Summary

- ❑ Gradient Descent Optimization
- ❑ Chain rule of derivatives
- ❑ Back propagation

## Next

- ❑ Advanced optimization methods
- ❑ Network regularization
- ❑ Practical tricks for training neural networks