

# CS 4803 / 7643: Deep Learning

Topics:

- Analytical Gradients
- Automatic Differentiation
  - Computational Graphs
  - Forward mode vs Reverse mode AD

Dhruv Batra  
Georgia Tech

# Administrativa

- HW1 Reminder
  - Due: 09/09, 11:59pm

# Recap from last time

$$\left[ \min_{\vec{w}} L(\vec{w}) \right]$$

Strategy: Follow the slope

Gradient  
Descent



$$L(\vec{w})$$

# Gradient Descent

$$\min_{\vec{w}} L(\vec{w})$$

```
# Vanilla Gradient Descent
```

```
while True:
```

```
    weights_grad = evaluate_gradient(loss_fun, data, weights)
```

```
    weights += - step_size * weights_grad # perform parameter update
```

backprop

$\vec{w}^{(0)}$  = Initialize

for  $t = 0, 1, 2, \dots$ , exhausted

$$\underline{w}^{(t+1)} = \underline{w}^{(t)} - \boxed{\eta} \nabla_{\vec{w}} L(\vec{w})$$

0.001      Step-size / Learning-rate

# Stochastic Gradient Descent (SGD)

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

Full sum expensive  
when N is large!

Approximate sum  
using a minibatch of  
examples

32 / 64 / 128 common

```
# Vanilla Minibatch Gradient Descent
```

```
while True:
```

```
→ data batch = sample_training_data(data, 256) # sample 256 examples
```

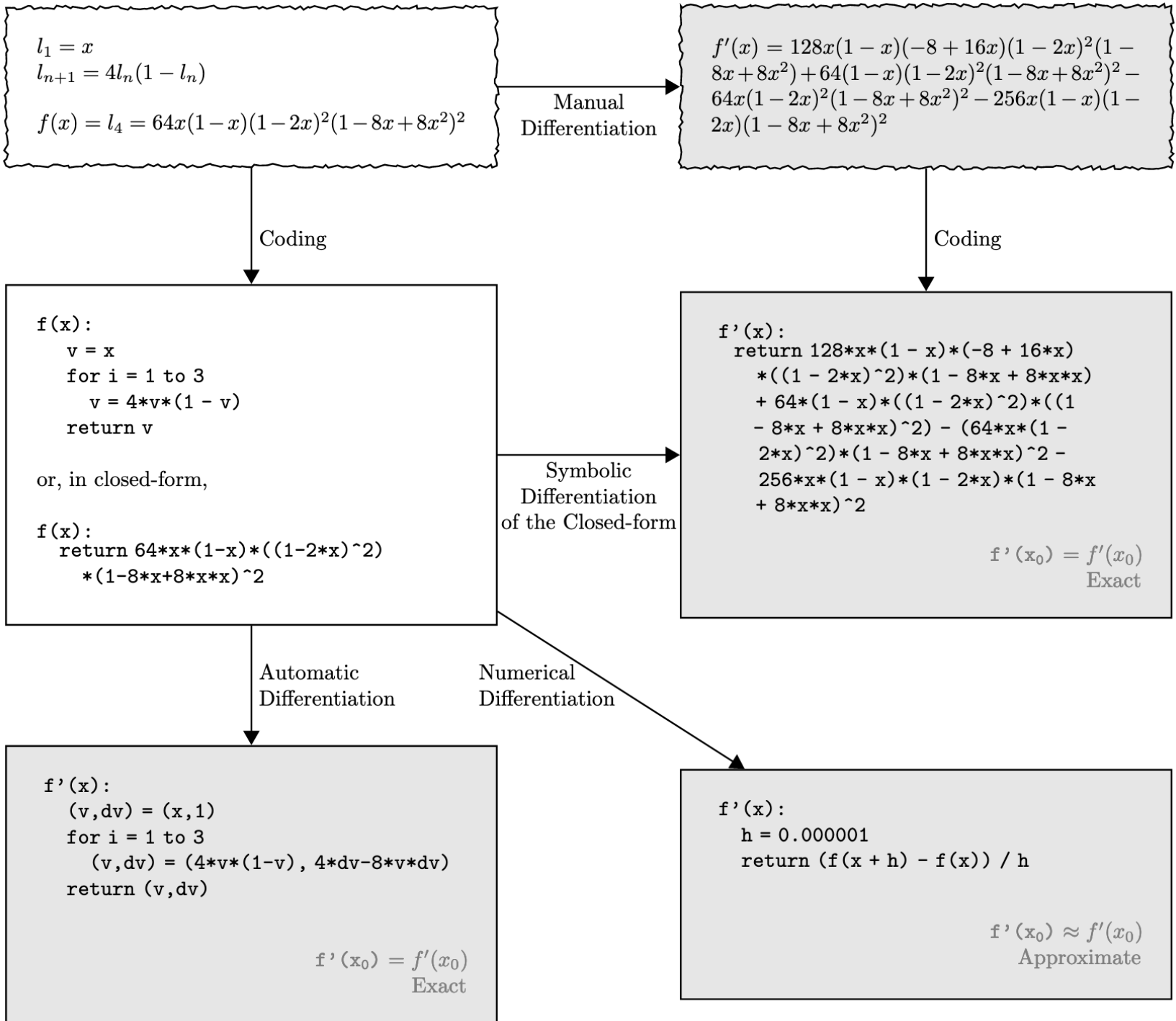
```
weights_grad = evaluate_gradient(loss_fun, data batch, weights)
```

```
weights += - step_size * weights_grad # perform parameter update
```

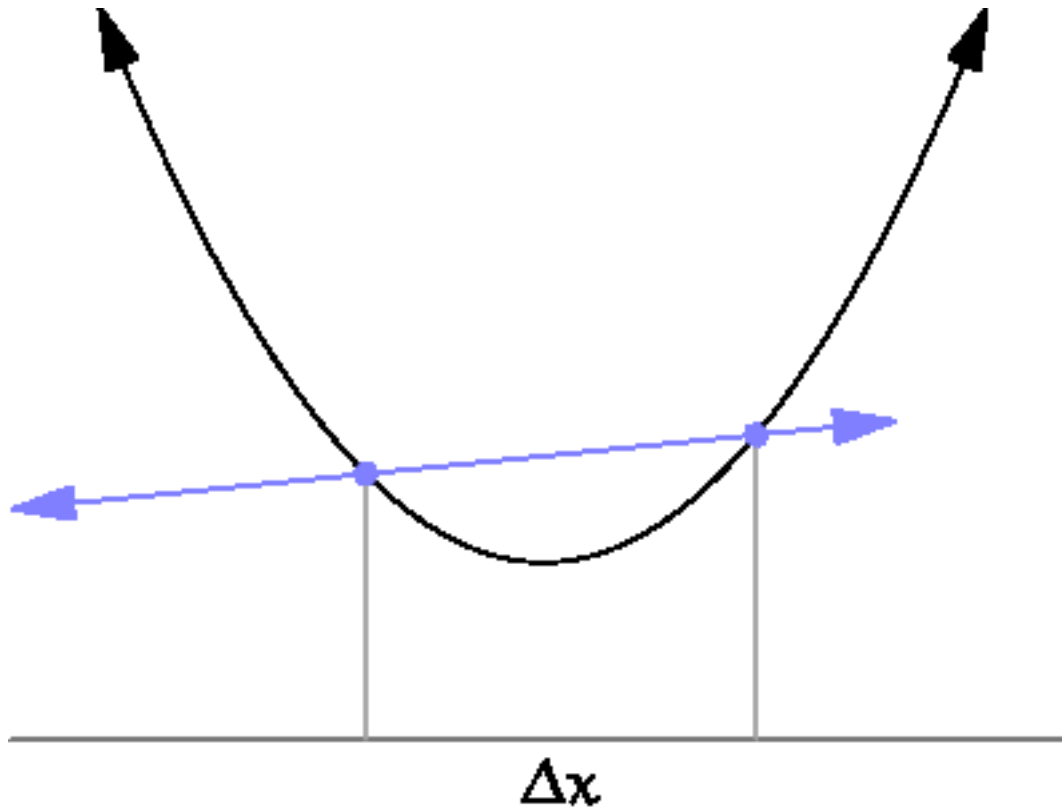


# How do we compute gradients?

- Analytic or “Manual” Differentiation
- Symbolic Differentiation X
- Numerical Differentiation
- Automatic Differentiation
  - Forward mode AD
  - Reverse mode AD
    - aka “backprop”







current W:

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

loss 1.25347

W + h (first dim):

[0.34 + 0.0001,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

loss 1.25322

gradient dW:

[-2.5,  
?,  
?,

$$(1.25322 - 1.25347)/0.0001 = -2.5$$

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

?,  
?,...]

# Numerical vs Analytic Gradients

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

**Numerical gradient:** slow :(), approximate :(), easy to write :)  
**Analytic gradient:** fast :), exact :), error-prone :(

In practice: Derive analytic gradient, check your implementation with numerical gradient.

This is called a **gradient check**.

# Plan for Today

- Analytical Gradients
- Automatic Differentiation
  - Computational Graphs
  - Forward mode vs Reverse mode AD

# Example: Logistic Regression

Input:  $\mathbf{x} \in \mathbb{R}^D$

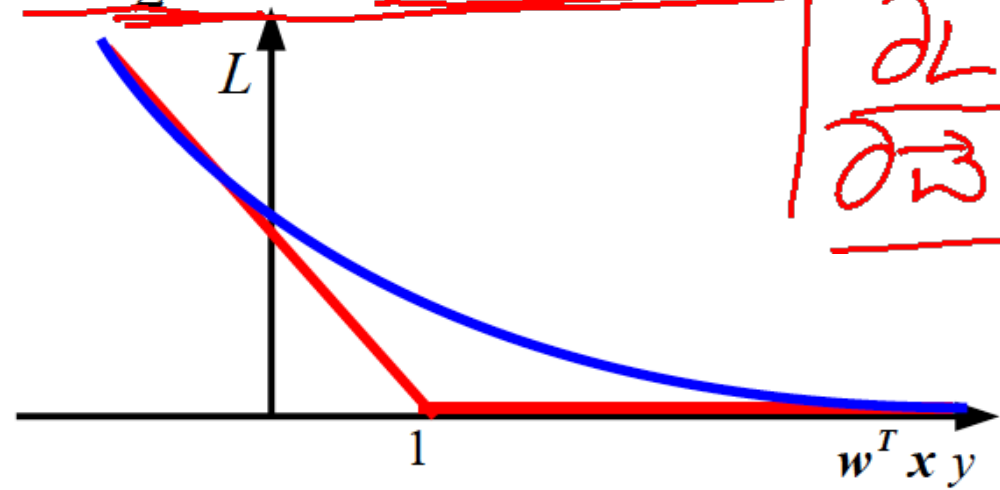
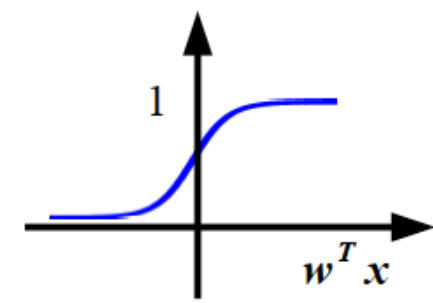
Binary label:  $y \in \{-1, +1\}$

Parameters:  $\mathbf{w} \in \mathbb{R}^D$   $\vec{w}_1, \vec{w}_2$

Output prediction:  $p(y=1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$

Loss:  $L = \frac{1}{2} \|\mathbf{w}\|^2 - \lambda \log(p(y|\mathbf{x}))$

$$\frac{\partial L}{\partial \vec{w}}$$



Log Loss

# Vector/Matrix Derivatives [Notation]

	$S$	$V$	$M$
$S$	$\frac{\partial y}{\partial x}$	$\frac{\partial \vec{y}}{\partial x}$	$\frac{\partial y}{\partial x}$
$V$	$\frac{\partial \vec{y}}{\partial x}$	$\frac{\partial \vec{y}}{\partial x}$	tensors
$M$	$\frac{\partial Y}{\partial x}$		

$$x, y \in \mathbb{R}$$

$$\vec{x} \in \mathbb{R}^d \quad y \in \mathbb{R}^c$$

$$X, Y \in \mathbb{R}^{m \times n}$$

$$\frac{\partial \vec{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \vdots \\ \frac{\partial y_c}{\partial x} \end{bmatrix}_{c \times 1}$$

$$\text{num} = \frac{\text{dim } 1}{\text{rows}}$$

$$\frac{\partial L}{\partial \vec{w}} = \begin{bmatrix} \frac{\partial L}{\partial w_1} & \dots & \frac{\partial L}{\partial w_n} \end{bmatrix}$$

$$\text{den} = \frac{\text{dim } 2}{\text{col}}$$

# Vector/Matrix Derivatives Notation

$$\frac{\partial \vec{y}}{\partial \vec{x}} = \begin{matrix} & \xrightarrow{\quad} & i & \\ \downarrow & \left[ \begin{array}{ccc} \frac{\partial y_i}{\partial x_1} & \dots & \frac{\partial y_i}{\partial x_3} & \dots & \frac{\partial y_i}{\partial x_d} \\ & & \frac{\partial y_c}{\partial x_j} & & \end{array} \right] & \downarrow & c \times d \\ & \underbrace{\hspace{10em}} & & \end{matrix}$$

Jacobian

# Vector Derivative Example

$$\vec{y} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} x \\ x^2 \end{bmatrix} \quad \frac{\partial \vec{y}}{\partial x} = \begin{bmatrix} 1 \\ 2x \end{bmatrix}$$

$$y = \vec{w}^T \vec{x} = \sum_i w_i x_i$$

$$\frac{\partial y}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & \dots & \frac{\partial y}{\partial x_n} \end{bmatrix} \quad \frac{\partial (\sum_i w_i x_i)}{\partial x_1}$$

$$= [w_1 \quad w_2 \quad \dots \quad w_n]$$

$$= \vec{w}^T$$

$$\frac{\partial (\vec{w}^T \vec{x})}{\partial \vec{x}} = \vec{w}^T$$



# Vector Derivative Example

$$\frac{\partial (\vec{w}^T A \vec{x})}{\partial \vec{w}} = \underline{2 \vec{w}^T [A]}$$

$$\vec{y} = A \vec{x}$$

$$\frac{\partial \vec{y}}{\partial \vec{x}} = A$$

$$y_i = \sum_j a_{ij} x_j$$

$$i \cdot \left[ \frac{\partial y_i}{\partial x_j} \right] = \left[ a_{ij} \right]$$

$$\frac{\partial y_i}{\partial x_j} = a_{ij}$$

# Extension to Tensors

$$\underline{X} \in \mathbb{R}^{d_1 \times \dots \times d_n}$$

$$\underline{Y} \in \mathbb{R}^{c_1 \times \dots \times c_n}$$



$$y\text{-vec} = Y(:, :)$$

$$x\text{-vec} = X(:, :)$$

$$\frac{\partial Y[i_1, \dots, i_m]}{\partial X[j_1, \dots, j_m]}$$

$$\frac{\partial \overrightarrow{y\text{-vec}}}{\partial \overrightarrow{x\text{-vec}}} = \begin{bmatrix} ( ) \\ ( ) \\ ( ) \end{bmatrix}$$

# Chain Rule: Composite Functions

$$f(g(x)) = (f \circ g)(x)$$

$$\begin{aligned} \overbrace{f(x)}^{L(\vec{w})} &= g_n(g_{n-1} \dots g_1(x)) \\ &= (g_n \circ g_{n-1} \circ \dots \circ g_1)(x) \end{aligned}$$

# Chain Rule: Scalar Case

$$\underline{x} \in \mathbb{R}^1 \xrightarrow{g_1(\cdot)} z \in \mathbb{R}^1 \xrightarrow{g_2(\cdot)} \underline{y} \in \mathbb{R}^1$$
$$y = g_2(\underbrace{g_1(x)}_z)$$

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial x}$$

Scalar mult.

# Chain Rule: Vector Case

$$\vec{x} \in \mathbb{R}^d \xrightarrow[\substack{g_1(\cdot) \\ : \mathbb{R}^d \rightarrow \mathbb{R}^m}]{\quad} \vec{z} \in \mathbb{R}^m \xrightarrow[\substack{g_2(\cdot) \\ : \mathbb{R}^m \rightarrow \mathbb{R}^c}]{\quad} \vec{y} \in \mathbb{R}^c$$

$$\boxed{\left[ \frac{\partial \vec{y}}{\partial \vec{x}} \right]} = \boxed{\left[ \frac{\partial \vec{y}}{\partial \vec{z}} \right]} \cdot \boxed{\left[ \frac{\partial \vec{z}}{\partial \vec{x}} \right]}$$

Matrix Mult

$J_{g_2}$        $J_{g_1}$

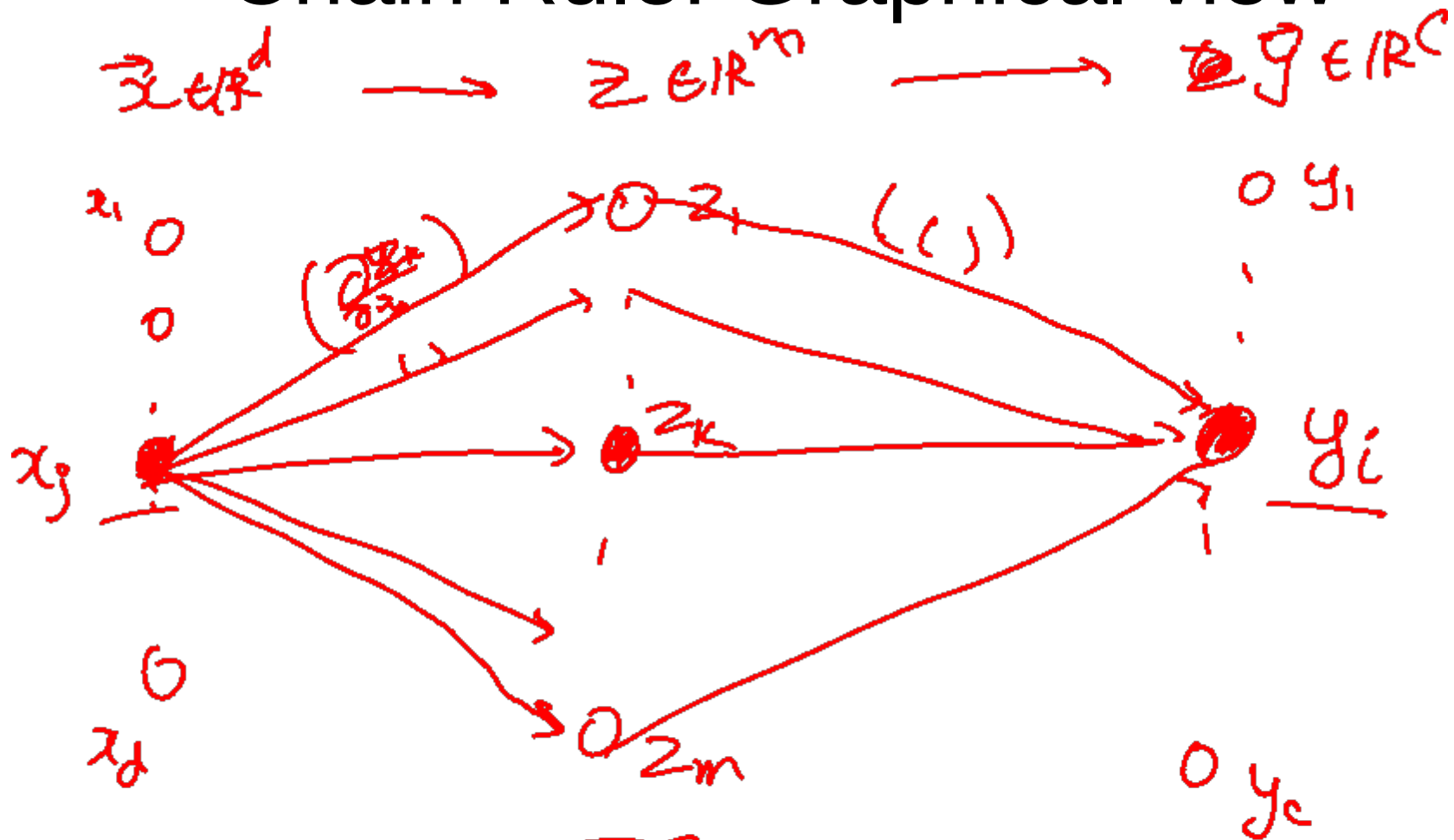
$J_{g_2 \circ g_1}$

# Chain Rule: Jacobian view

$$\begin{matrix}
 \frac{\partial y}{\partial z} \\
 \vdots \\
 \frac{\partial y_i}{\partial x_j} \\
 \vdots \\
 \frac{\partial y_j}{\partial x_i} \\
 \vdots \\
 \frac{\partial y}{\partial x}
 \end{matrix}
 \begin{matrix}
 \frac{\partial z}{\partial z} \\
 \vdots \\
 \frac{\partial z_k}{\partial z} \\
 \vdots \\
 \frac{\partial z}{\partial z} \\
 \vdots \\
 \frac{\partial z}{\partial z}
 \end{matrix}
 \begin{matrix}
 \frac{\partial z}{\partial x} \\
 \vdots \\
 \frac{\partial z_k}{\partial x_j} \\
 \vdots \\
 \frac{\partial z}{\partial x_i} \\
 \vdots \\
 \frac{\partial z}{\partial x}
 \end{matrix}$$

$$\frac{\partial y_i}{\partial x_j} = \sum_{k \in K} \frac{\partial y_i}{\partial z_k} \frac{\partial z_k}{\partial x_j}$$

# Chain Rule: Graphical view



$$\frac{\partial y_i}{\partial x_j} = \sum_{k \text{ on path}} \frac{\partial y_i}{\partial z_k} \frac{\partial z_k}{\partial x_j}$$

# Chain Rule: Cascaded

$$\underline{x} = \underline{h}^{(0)} \xrightarrow{g_1} \underline{h}^{(1)} \in \mathbb{R}^d \xrightarrow{g_2} \underline{h}^{(2)} \in \mathbb{R}^d \dots \xrightarrow{g_{\ell}} \underline{h}^{(\ell)} \rightarrow \mathcal{L}$$

$$\begin{bmatrix} 0 \\ \vdots \\ 0 \\ 0 \\ 0 \end{bmatrix} \xrightarrow{g_1} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix} \dots \dots \dots \begin{bmatrix} \cdot \\ \vdots \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} \xrightarrow{g_{\ell}} \begin{bmatrix} \cdot \\ \vdots \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} \rightarrow \mathcal{L}$$

$$\frac{\partial \mathcal{L}}{\partial \underline{h}^{(0)}} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \underline{h}^{(0)}} & \frac{\partial \underline{h}^{(\ell)}}{\partial \underline{h}^{(0)}} \end{bmatrix} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \underline{h}^{(0)}} & \frac{\partial \underline{h}^{(\ell)}}{\partial \underline{h}^{(1)}} \cdot \frac{\partial \underline{h}^{(1)}}{\partial \underline{h}^{(0)}} \end{bmatrix} \dots \dots \dots \begin{bmatrix} \frac{\partial \underline{h}^{(\ell-1)}}{\partial \underline{h}^{(\ell-2)}} \dots \frac{\partial \underline{h}^{(2)}}{\partial \underline{h}^{(1)}} \end{bmatrix}$$

$$\mathcal{O}(d) \times \mathcal{O}(d) \times \mathcal{O}(d) = \underbrace{\mathcal{O}(d^2)}_{\left[ \begin{matrix} \cdot \\ \vdots \\ \cdot \\ \cdot \\ \cdot \end{matrix} \right]} \times \underbrace{\mathcal{O}(d^3)}_{\left[ \begin{matrix} \cdot \\ \vdots \\ \cdot \\ \cdot \\ \cdot \end{matrix} \right]} \times \mathcal{O}(d)$$



# Chain Rule: How should we multiply?

# Example: Logistic Regression

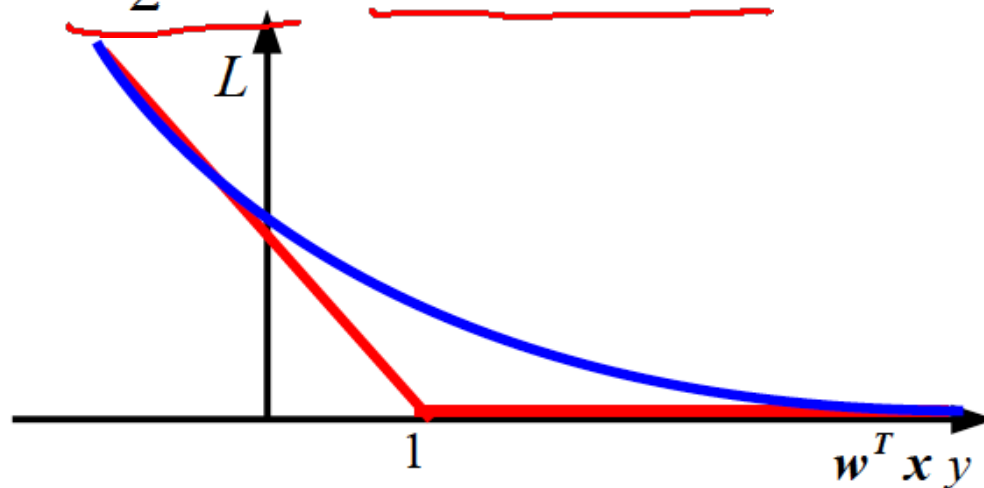
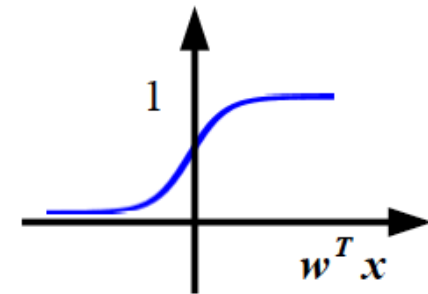
Input:  $\mathbf{x} \in \mathbb{R}^D$

Binary label:  $y \in \{-1, +1\}$

Parameters:  $\mathbf{w} \in \mathbb{R}^D$

Output prediction:  $p(y=1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$

Loss:  $L = \frac{1}{2} \|\mathbf{w}\|^2 - \lambda \log(p(y|\mathbf{x}))$



Log Loss

# Logistic Regression Derivatives

$$L = -\log\left(\frac{1}{1+e^{-w^T x}}\right) = \log\left(\frac{1+e^{-w^T x}}{P}\right)$$

$$\left[ \frac{\partial L}{\partial \vec{w}} = \left(\frac{1}{P}\right) \left(e^{-w^T x}\right) \cdot \left(-x^T\right) \right]$$

$$[ ]$$

# Logistic Regression Derivatives

# Convolutional network (AlexNet)

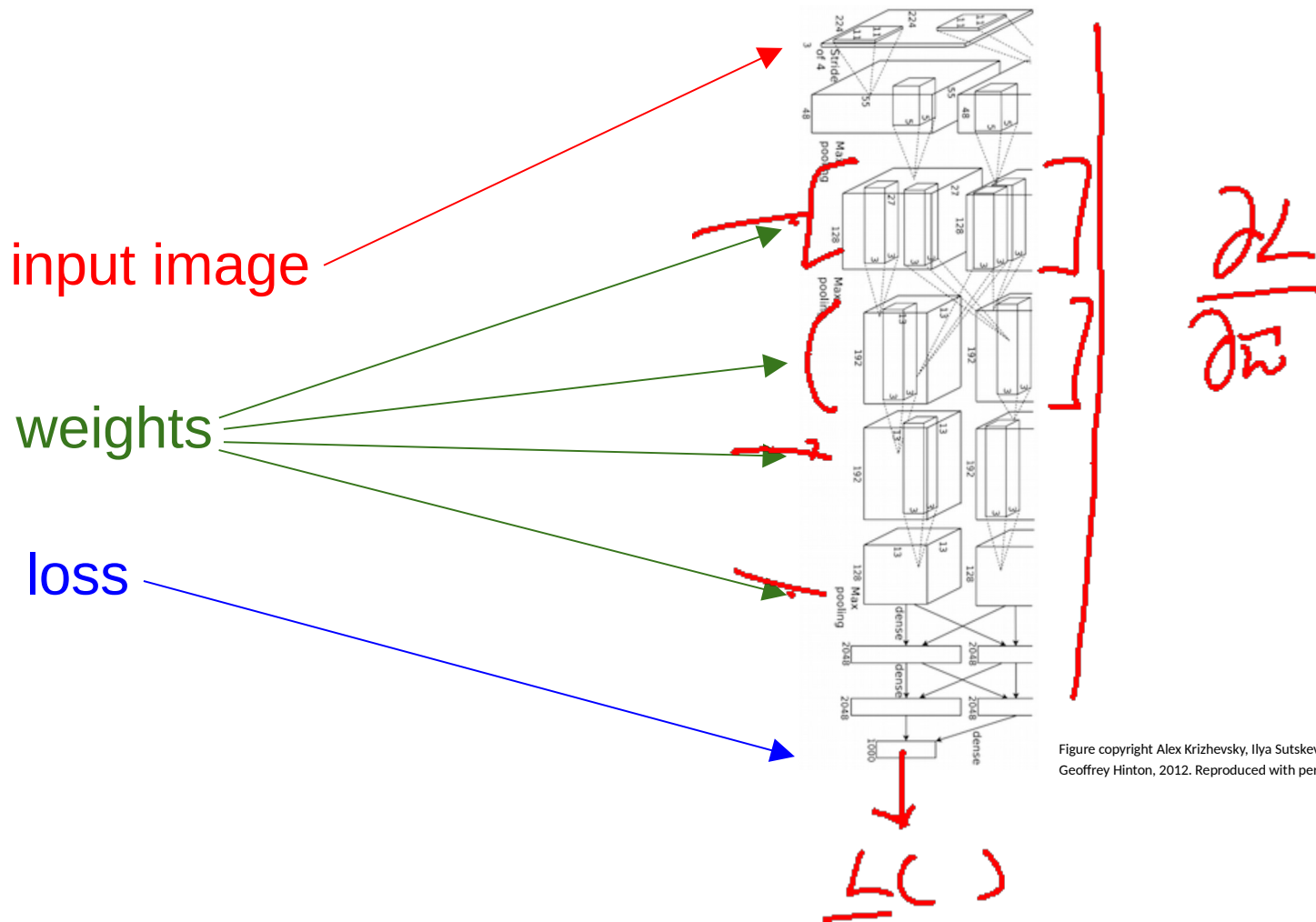


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

# Neural Turing Machine

input image

loss

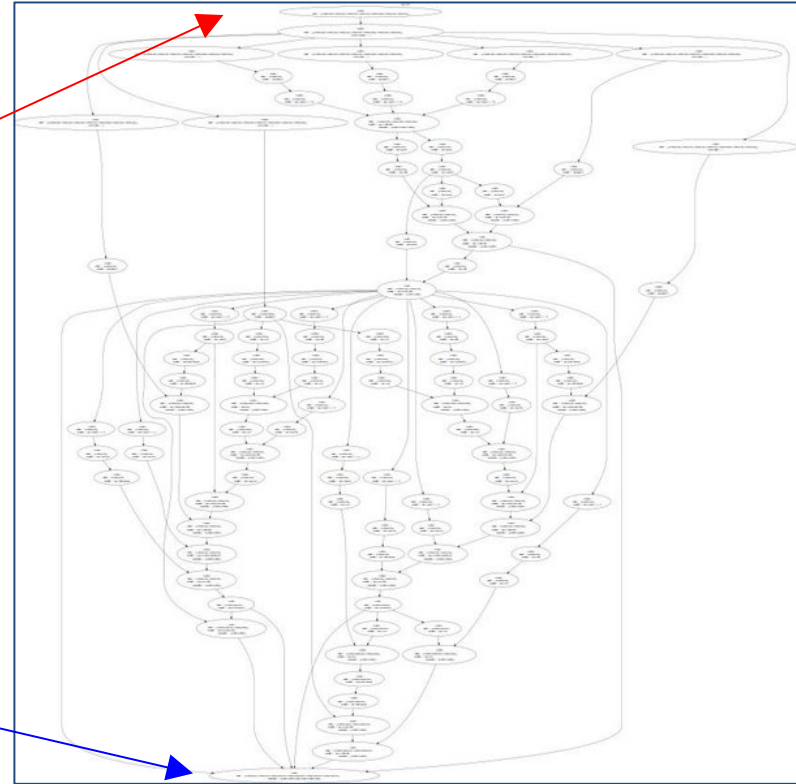


Figure reproduced with permission from a [Twitter post](#) by Andrej Karpathy.

# How do we compute gradients?

- Analytic or “Manual” Differentiation
- Symbolic Differentiation
- Numerical Differentiation

- Automatic Differentiation
  - Forward mode AD
  - Reverse mode AD
    - aka “backprop”

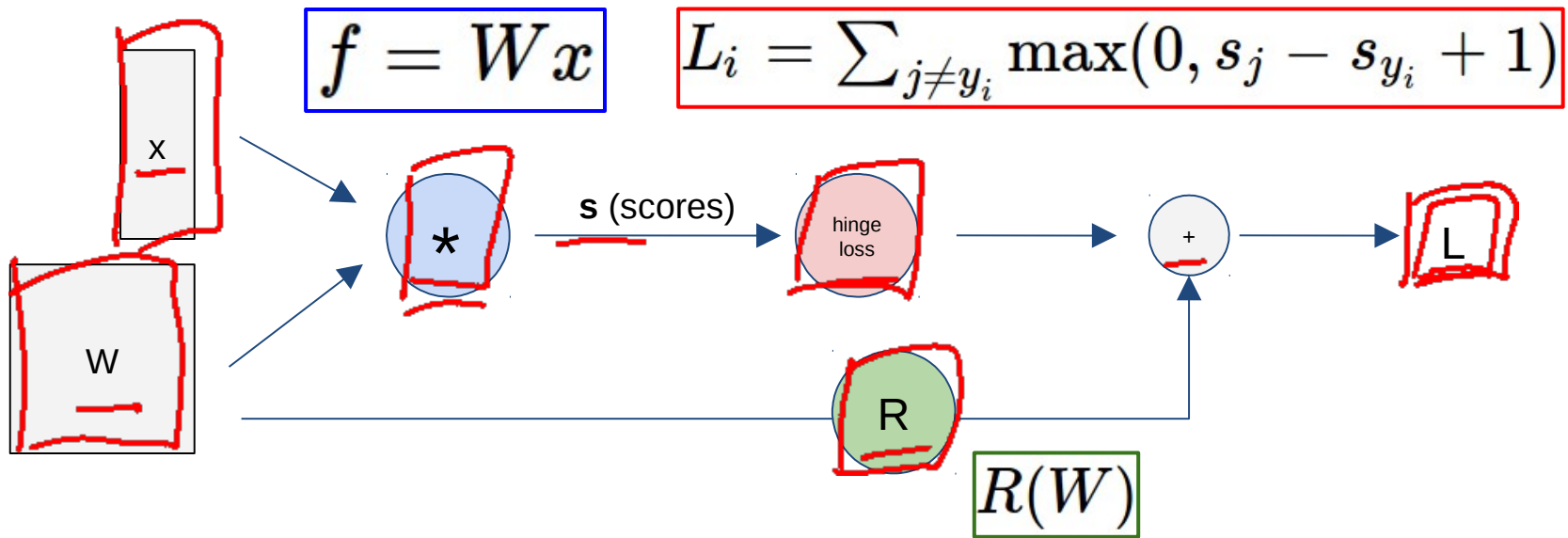
# Deep Learning = Differentiable Programming

- Computation = Graph
  - Input = Data + Parameters
  - Output = Loss
  - Scheduling = Topological ordering

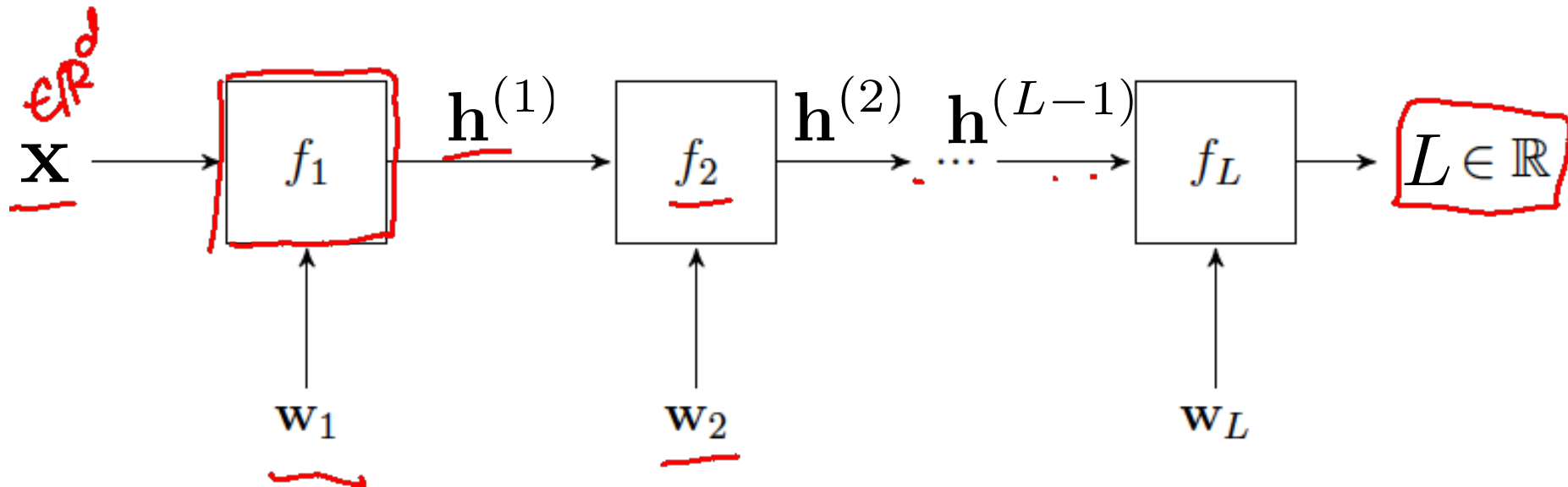
- Auto-Diff ←
  - A family of algorithms for implementing chain-rule on computation graphs



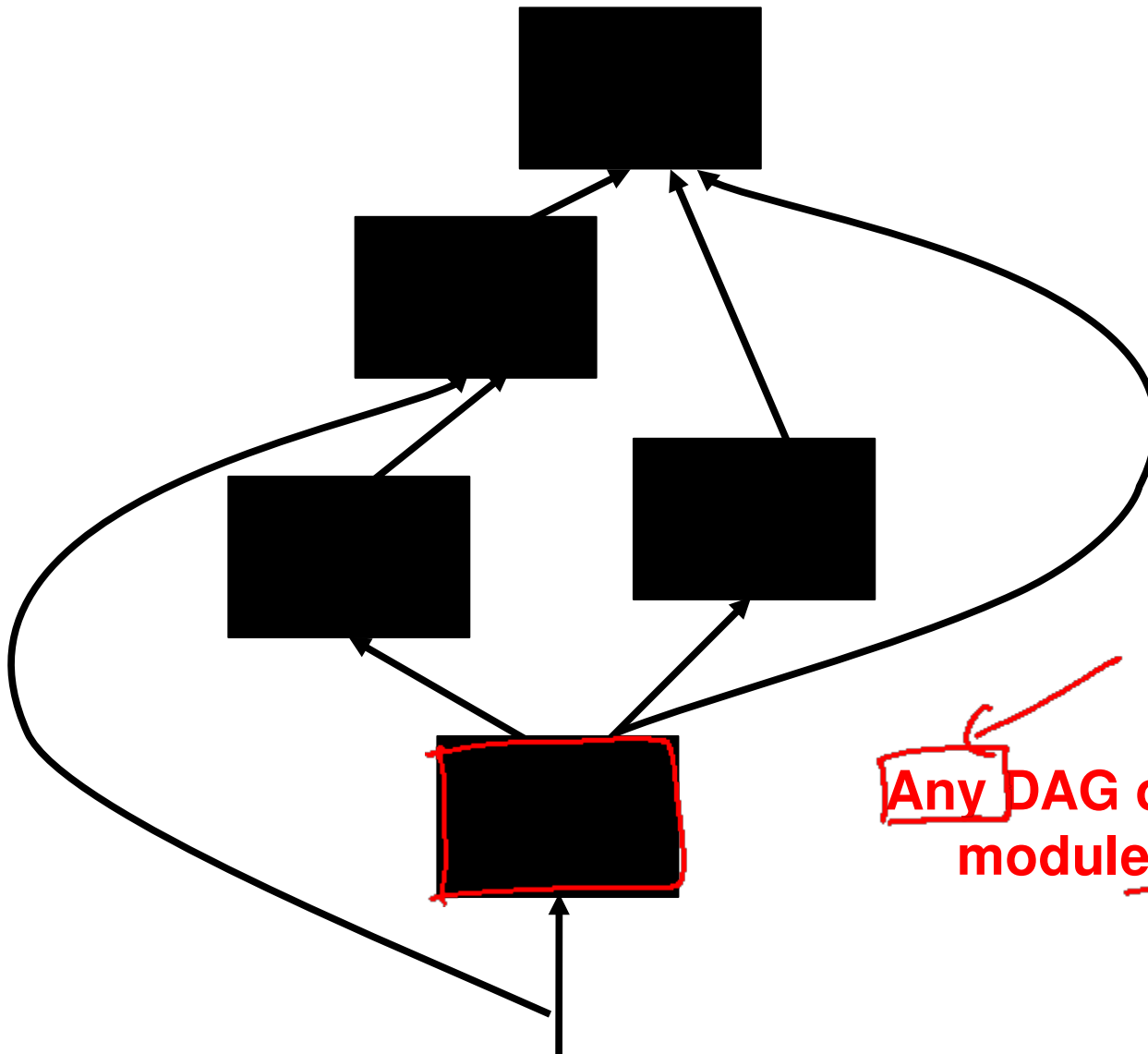
# Computational Graph



# Neural Network Computation Graph



# Computational Graph

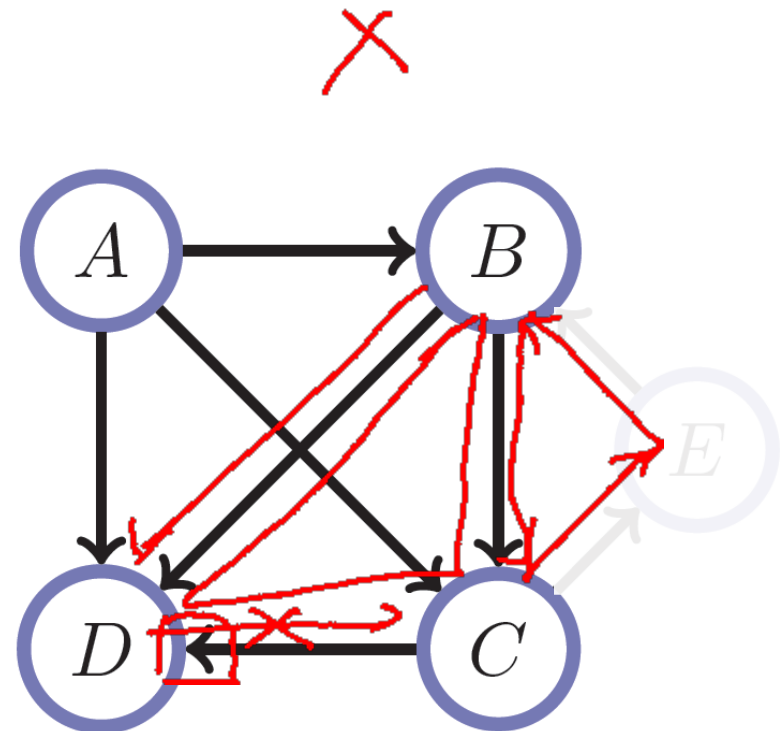
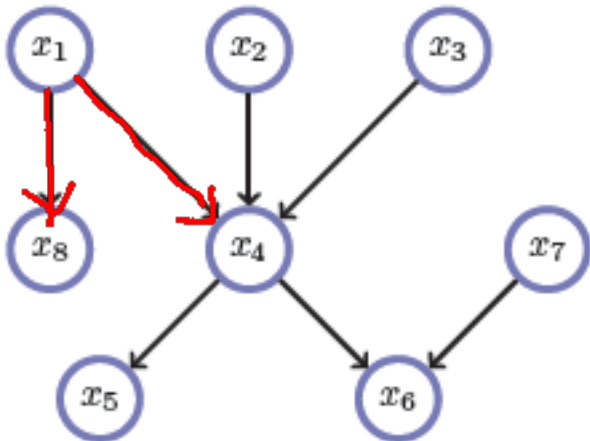


Any DAG of differentiable modules is allowed!

# Directed Acyclic Graphs (DAGs)

- Exactly what the name suggests
  - Directed edges
  - No (directed) cycles
  - Underlying undirected cycles okay

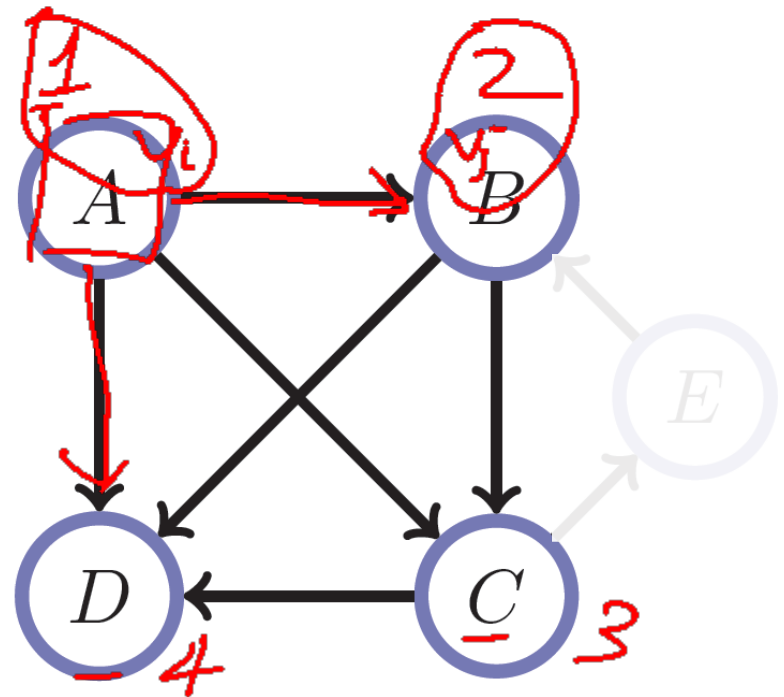
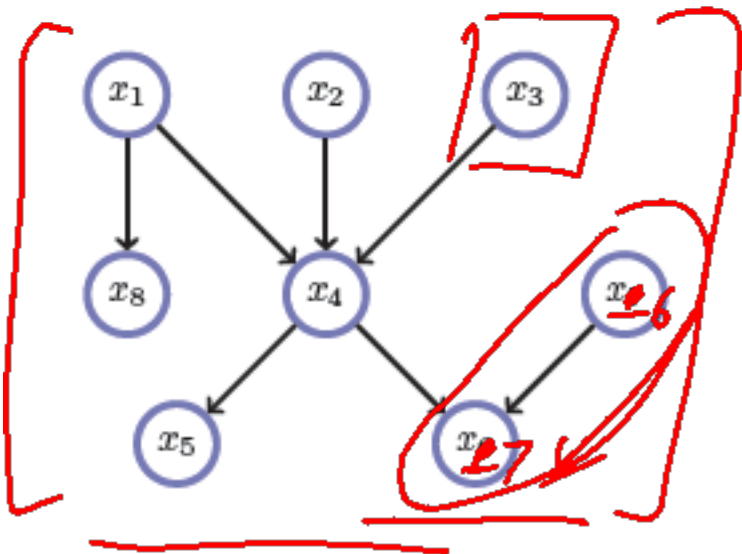
DAG ✓



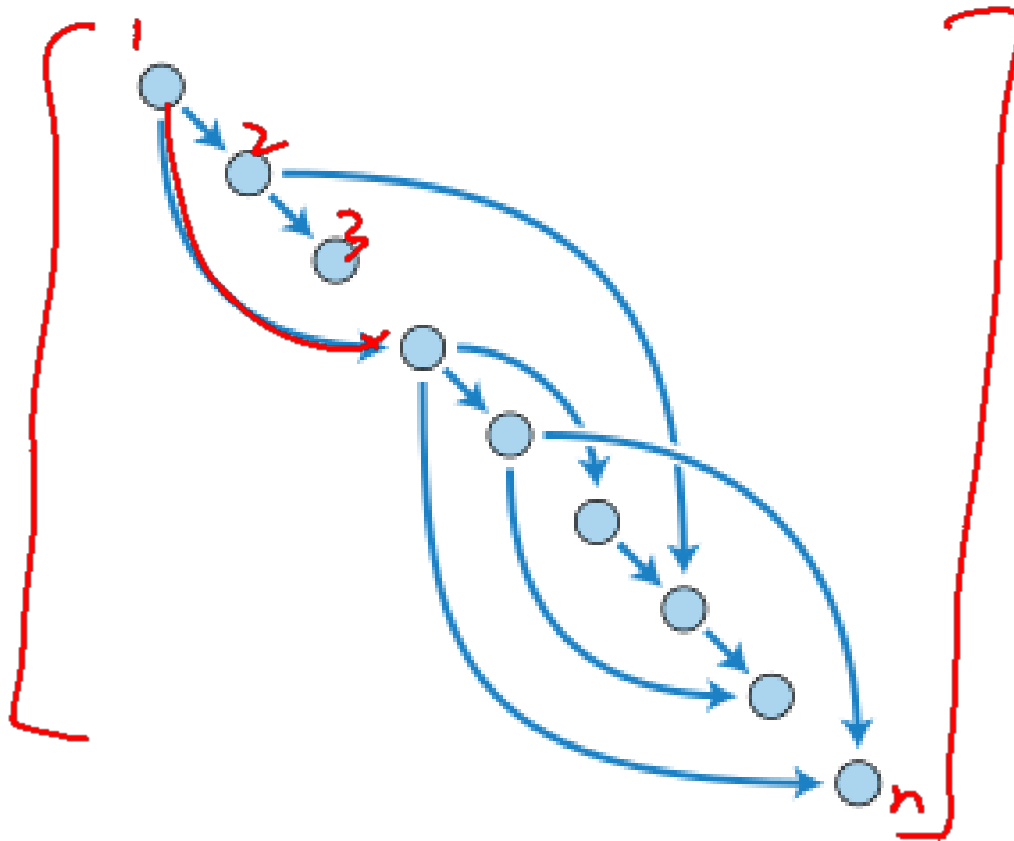
# Directed Acyclic Graphs (DAGs)

$\exists$  bijections:  $V \rightarrow \{1, \dots, n\}$   
s.t.  $\forall (v_i, v_j) \in E$   
 $\sigma(v_i) < \sigma(v_j)$

- Concept
  - Topological Ordering



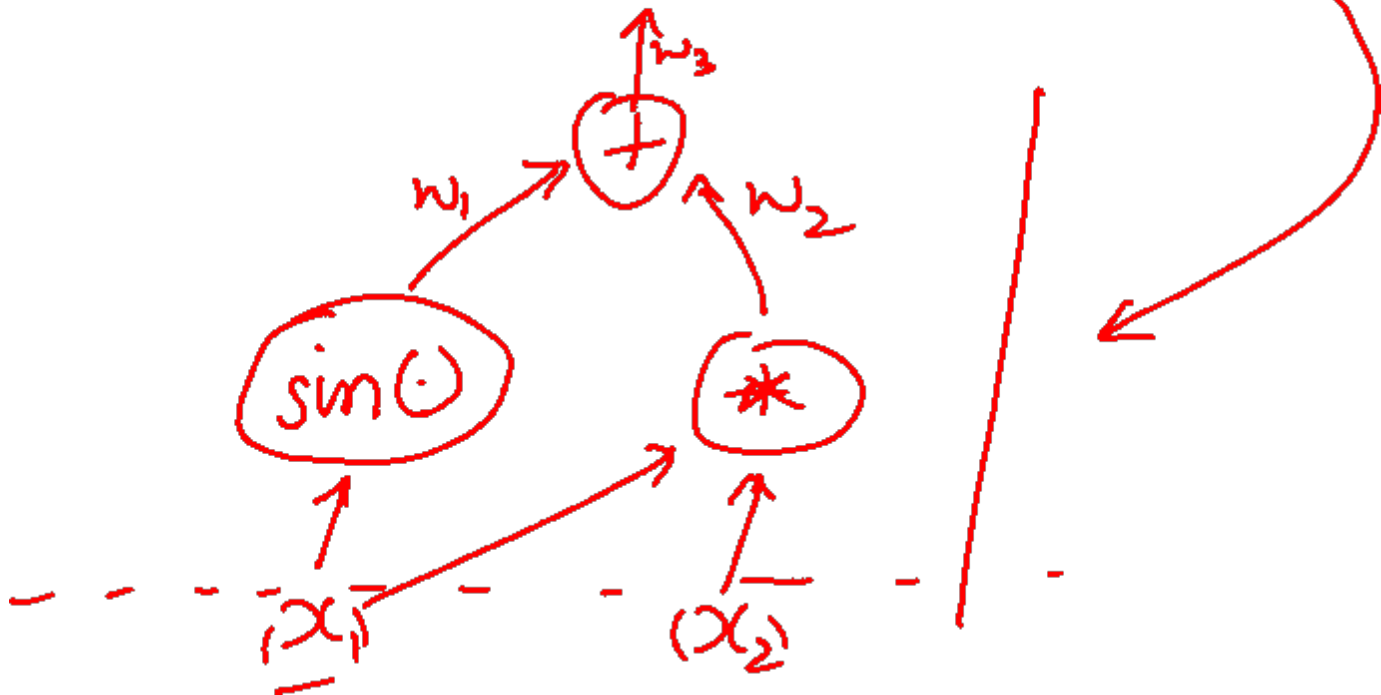
# Directed Acyclic Graphs (DAGs)



# Computational Graphs

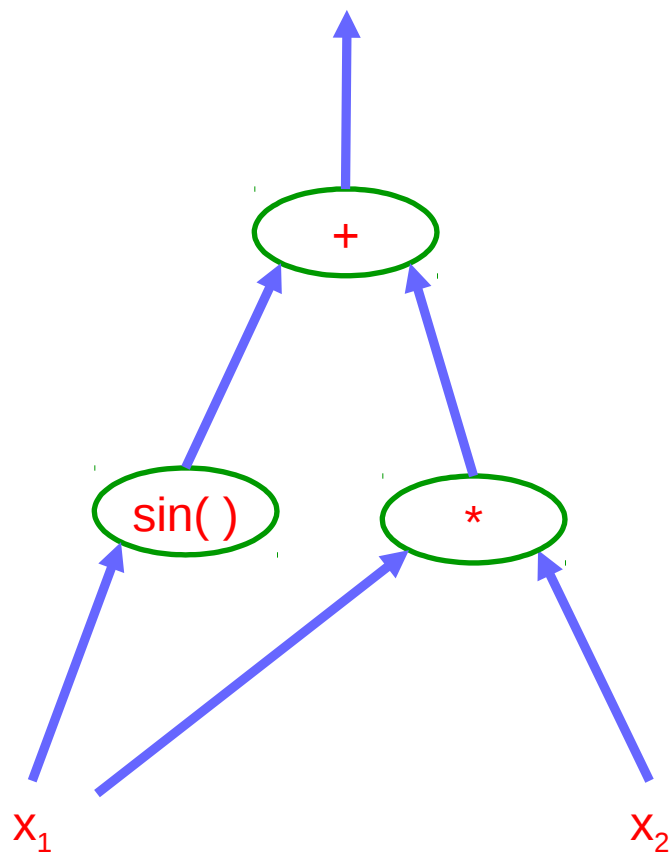
- Notation

$$f(\underline{x}_1, \underline{x}_2) = \overbrace{x_1 x_2}^{w_2} + \overbrace{\sin(x_1)}^{w_1}$$



# Example

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$



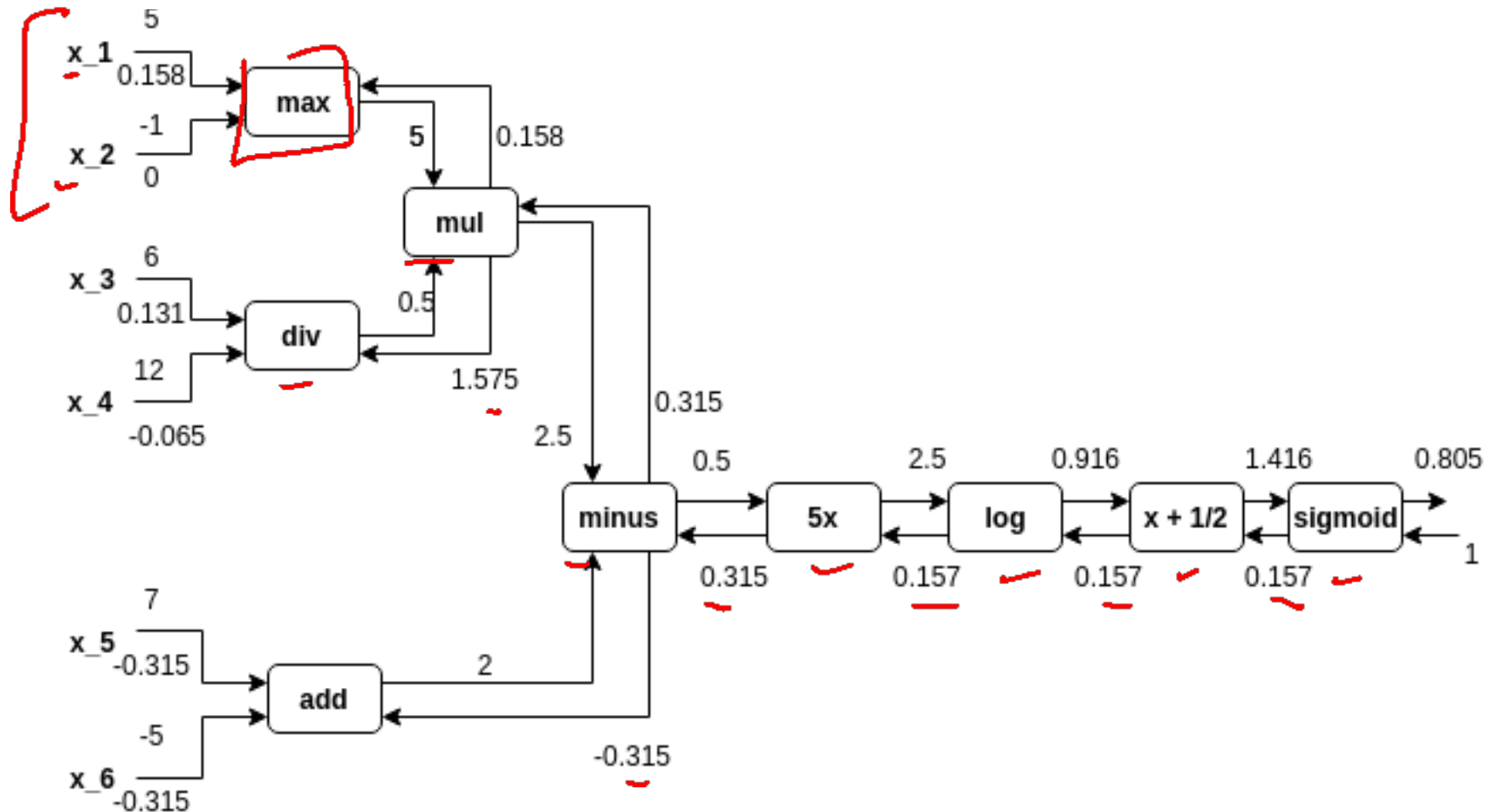


# HWO

$$\underline{f(\mathbf{x})} = \underline{\sigma} \left( \underline{\log} \left( \underline{5} \left( \underline{\max\{x_1, x_2\}} \cdot \underline{\frac{x_3}{x_4}} - \underline{(x_5 + x_6)} \right) \right) + \underline{\frac{1}{2}} \right)$$

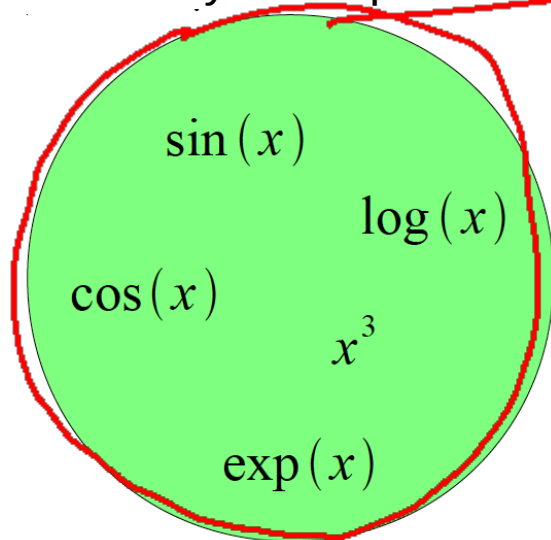
# HW0 Submission by Samyak Datta

$$f(\mathbf{x}) = \sigma \left( \log \left( 5 \left( \max\{x_1, x_2\} \cdot \frac{x_3}{x_4} - (x_5 + x_6) \right) \right) + \frac{1}{2} \right)$$



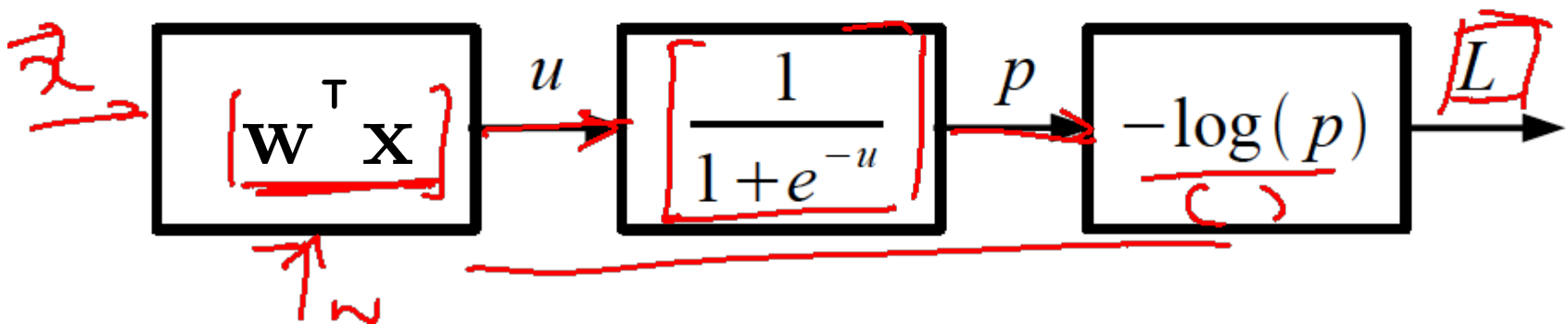
# Logistic Regression as a Cascade

Given a library of simple functions



Compose into a  
complicate function

$$-\log\left(\frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}\right)$$



# Deep Learning = Differentiable Programming

- Computation = Graph
  - Input = Data + Parameters
  - Output = Loss
  - Scheduling = Topological ordering
- Auto-Diff
  - A family of algorithms for implementing chain-rule on computation graphs

# [Forward mode] vs [Reverse Mode] AD

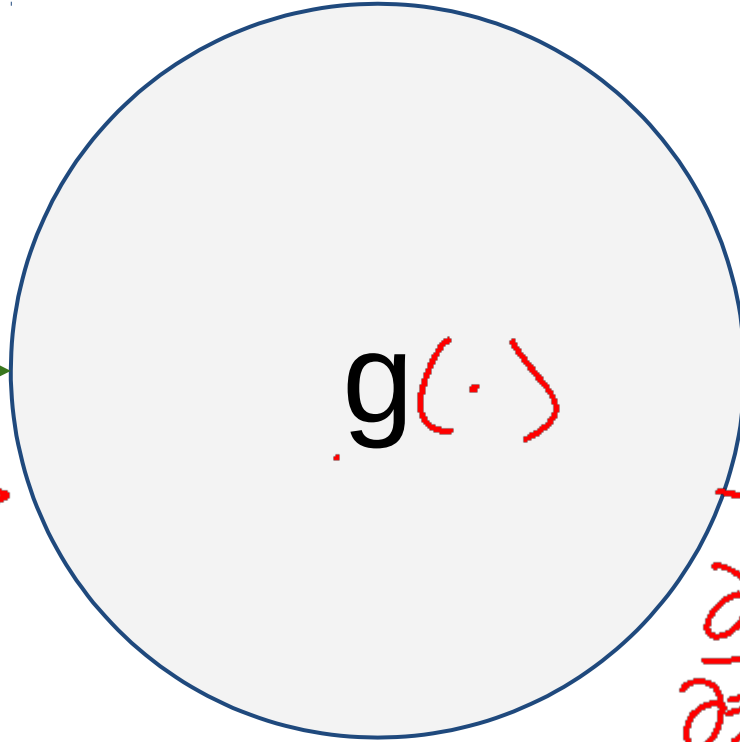
- Key Computations

# Forward mode AD

Goal:  $\frac{\partial L}{\partial \theta} \left[ \frac{\partial L}{\partial x} \right]$

$\vec{h}^{(l-1)}$

$\vec{h}^{(l)}$   
 $\vec{h} = g(\vec{h}^{(l-1)})$



$g(\cdot)$

Input:

$$\frac{\partial \vec{h}^{(l-1)}}{\partial \vec{x}}$$

$\uparrow$

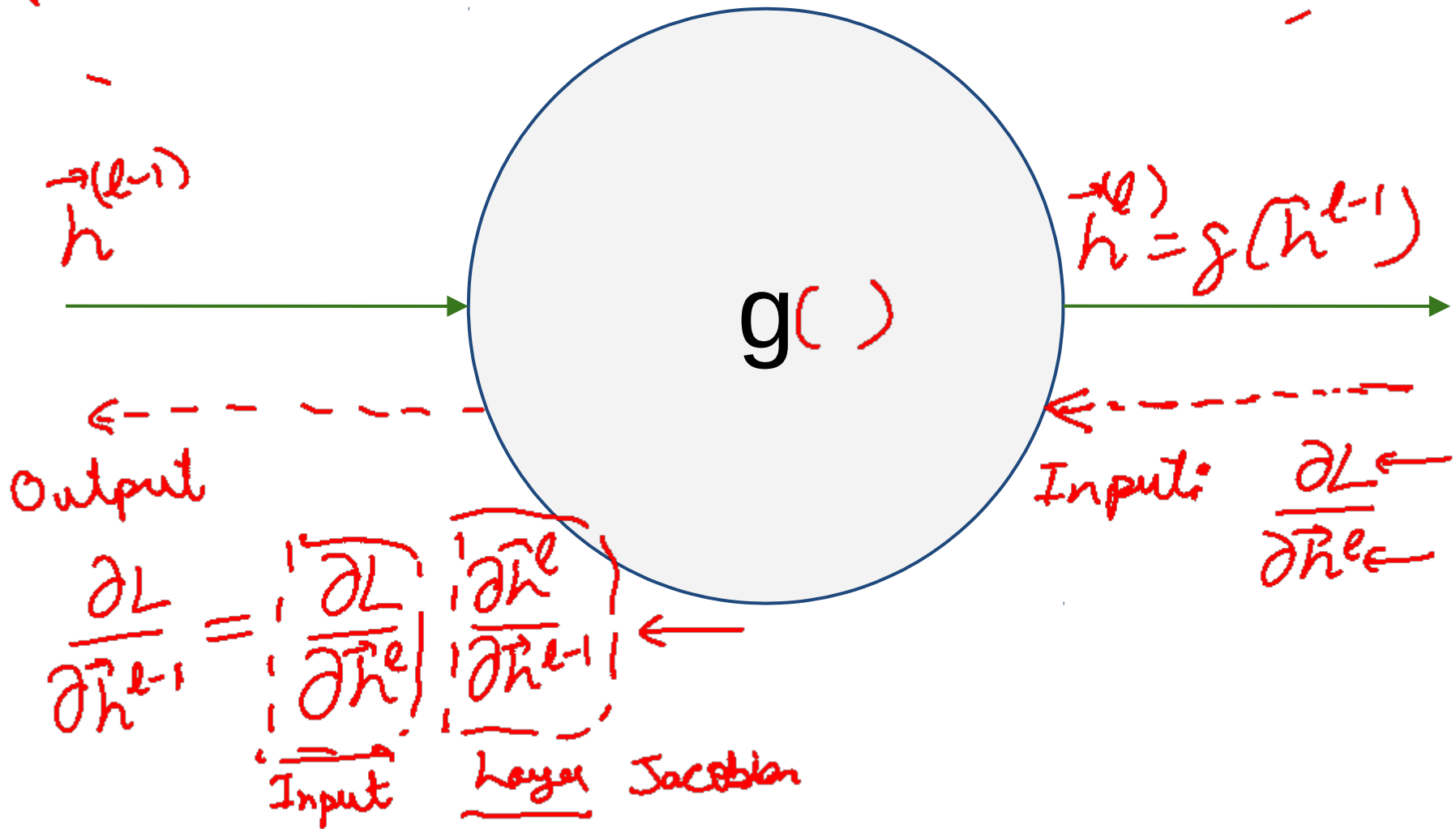
$$\frac{\partial \vec{h}^{(l)}}{\partial \vec{x}}$$

$$= \frac{\partial \vec{h}^{(l)}}{\partial \vec{h}^{(l-1)}} \frac{\partial \vec{h}^{(l-1)}}{\partial \vec{x}}$$

Layer Jacobian

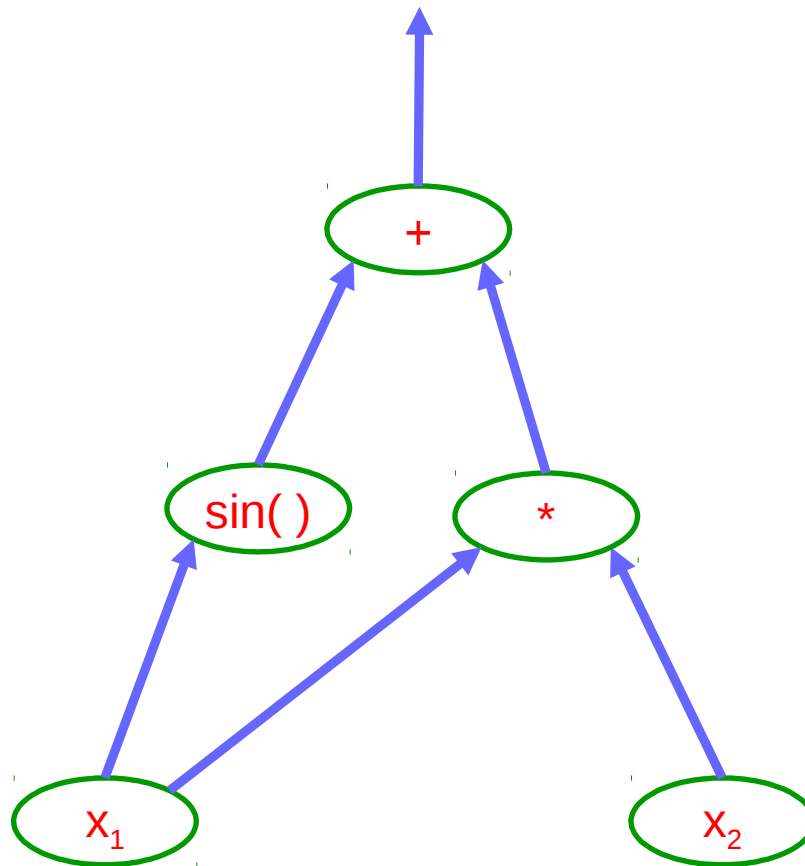
Input FM-A

# Reverse mode AD Goal: $\frac{\partial L}{\partial x}$



# Example: Forward mode AD

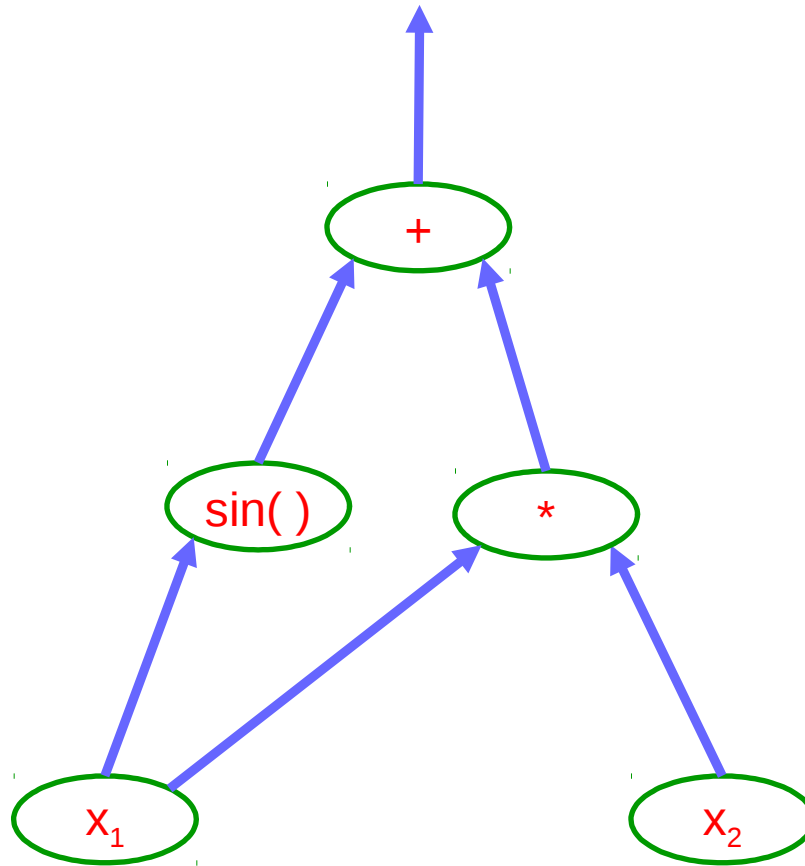
$$f(x_1, x_2) = \sin(x_1) + x_1x_2$$





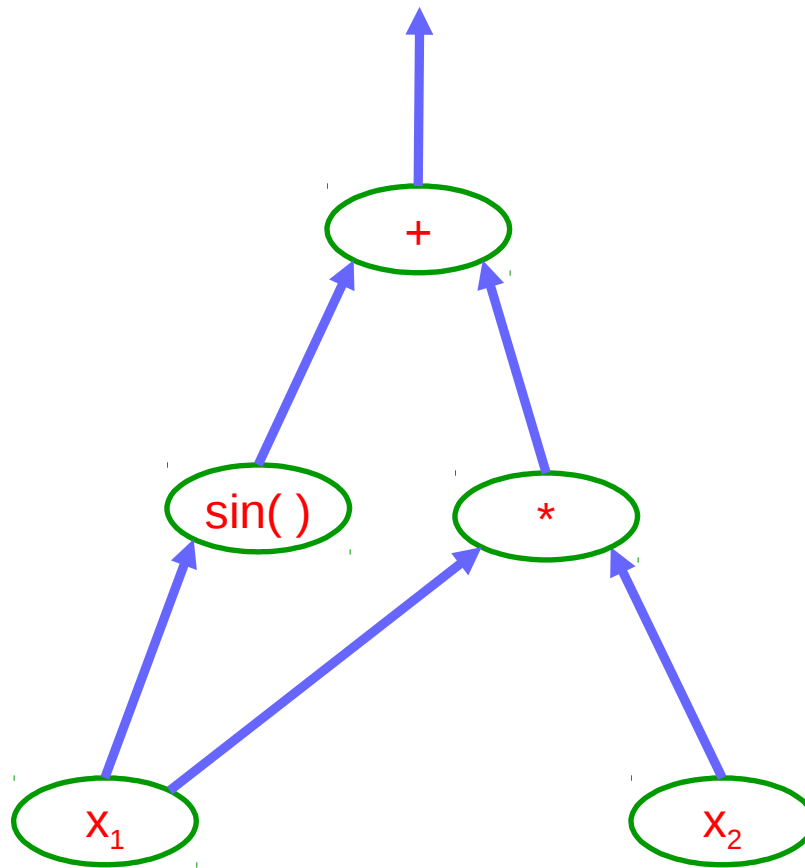
# Example: Forward mode AD

$$f(x_1, x_2) = \sin(x_1) + x_1x_2$$



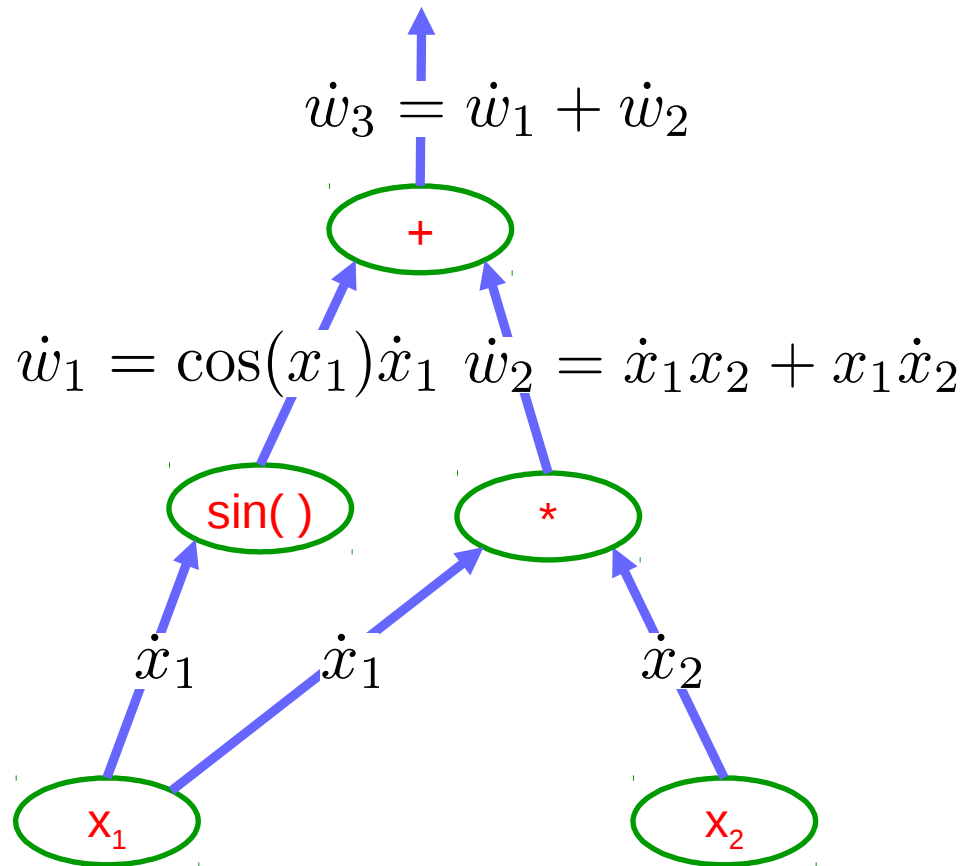
# Example: Forward mode AD

$$f(x_1, x_2) = \sin(x_1) + x_1x_2$$



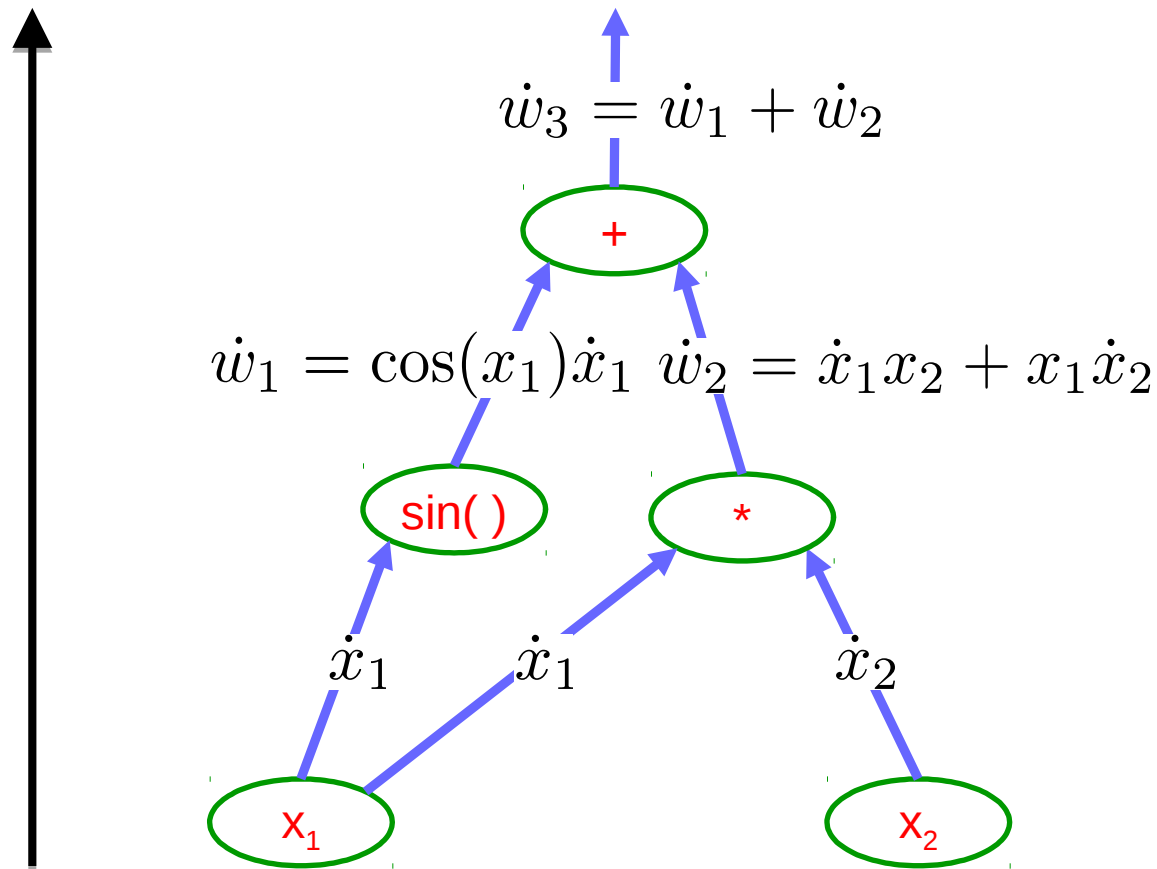
# Example: Forward mode AD

$$f(x_1, x_2) = \sin(x_1) + x_1x_2$$



# Example: Forward mode AD

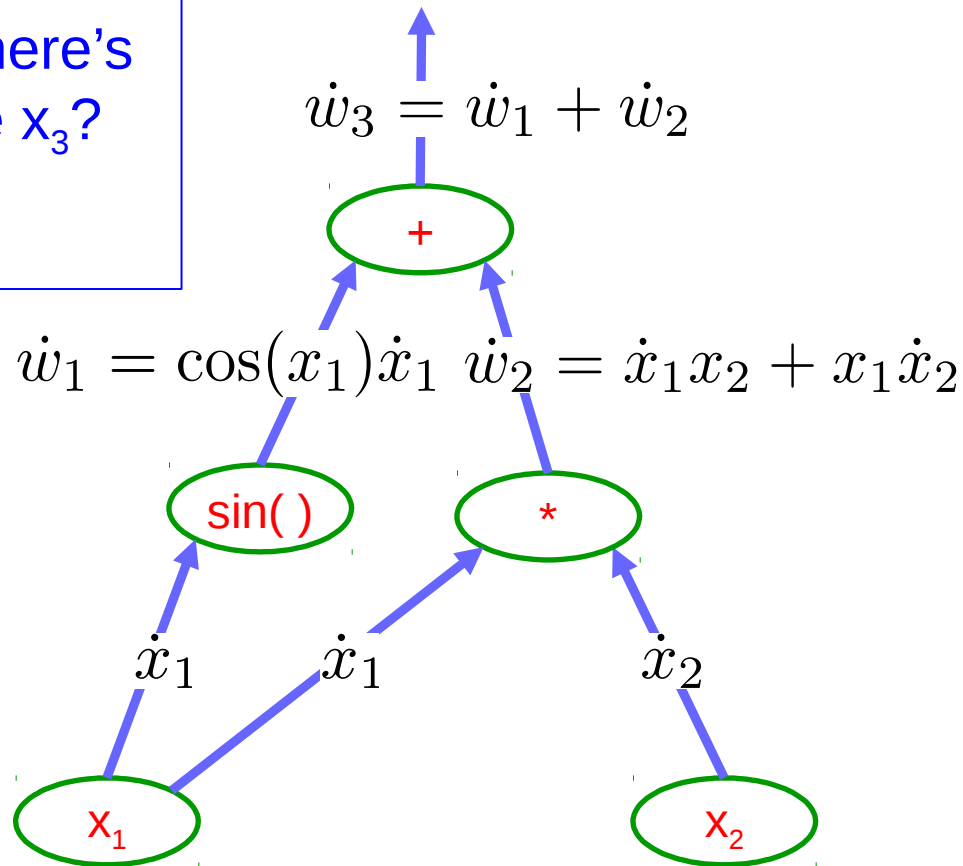
$$f(x_1, x_2) = \sin(x_1) + x_1x_2$$



# Example: Forward mode AD

$$f(x_1, x_2) = \sin(x_1) + x_1x_2$$

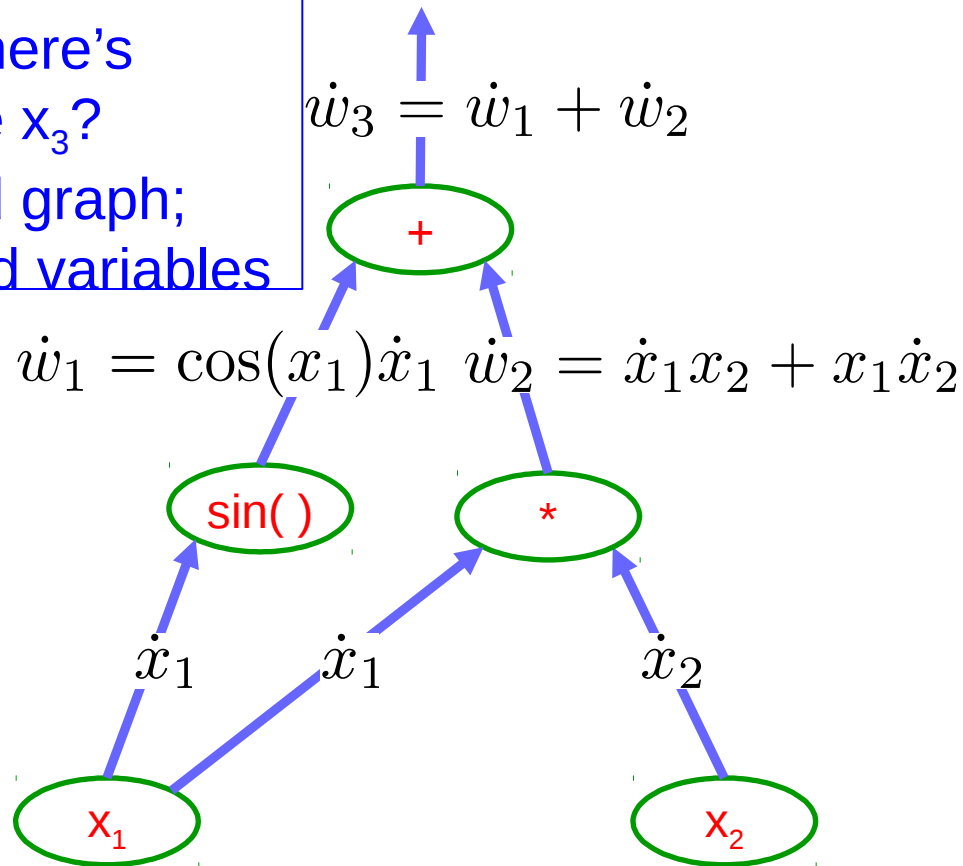
Q: What happens if there's another input variable  $x_3$ ?



# Example: Forward mode AD

$$f(x_1, x_2) = \sin(x_1) + x_1x_2$$

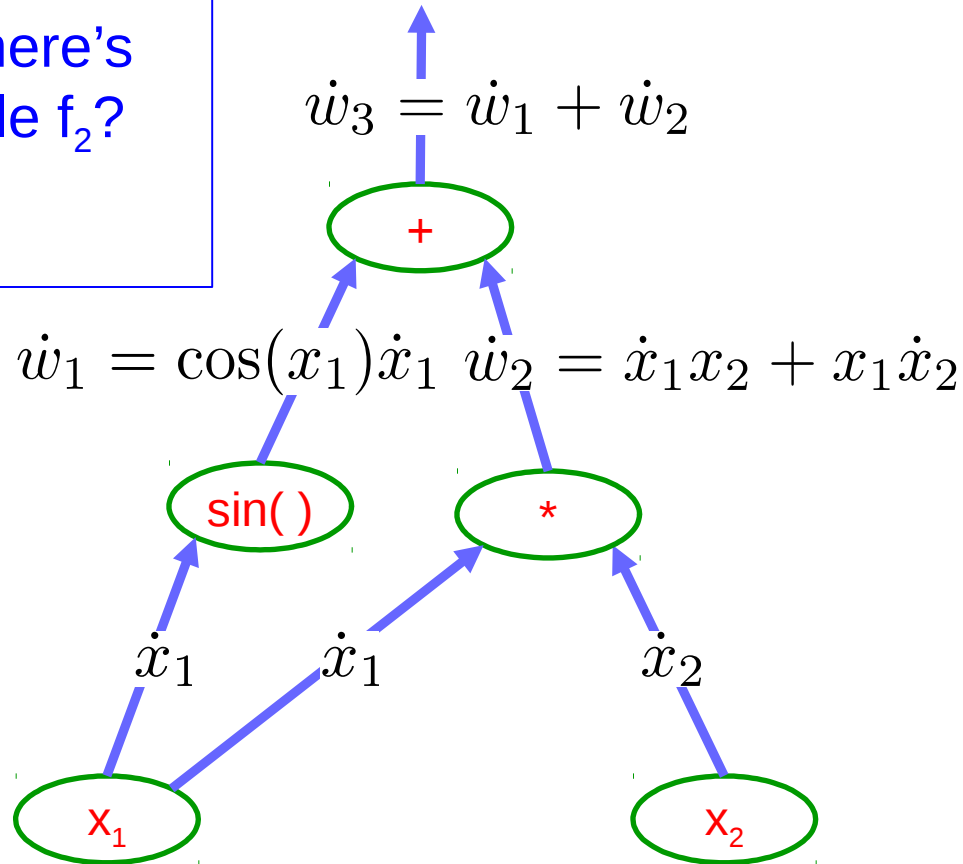
Q: What happens if there's another input variable  $x_3$ ?  
A: more sophisticated graph; d "forward props" for d variables



# Example: Forward mode AD

$$f(x_1, x_2) = \sin(x_1) + x_1x_2$$

Q: What happens if there's another output variable  $f_2$ ?



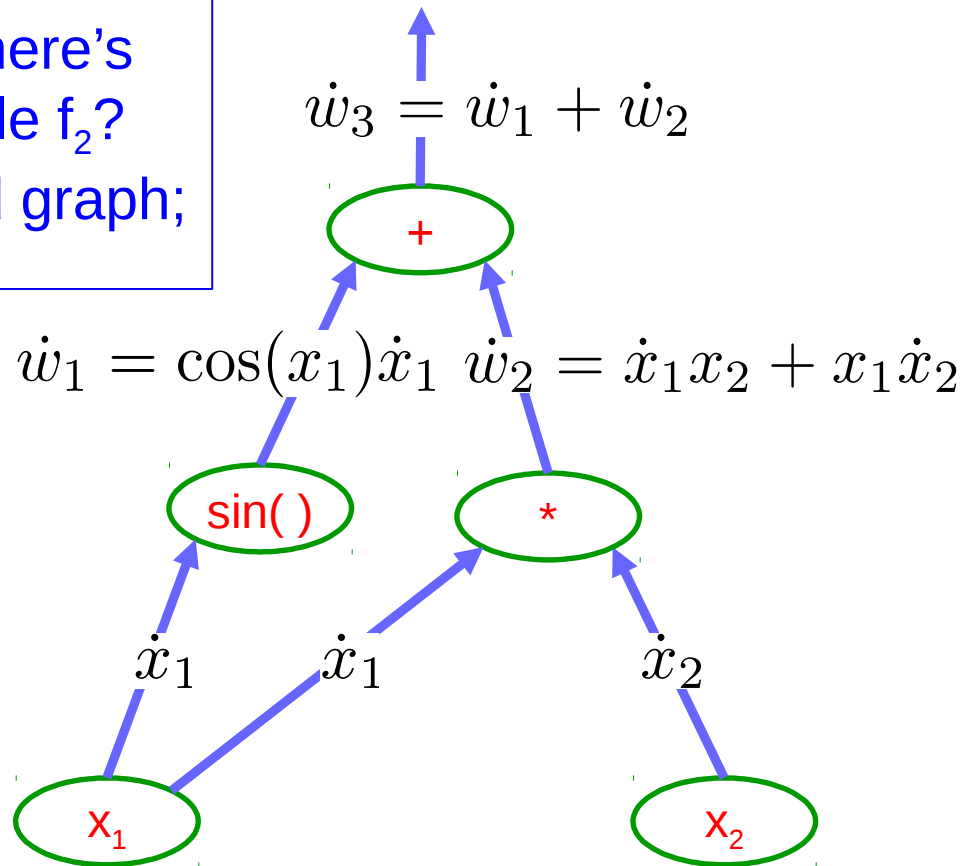
# Example: Forward mode AD

$$f(x_1, x_2) = \sin(x_1) + x_1x_2$$

Q: What happens if there's another output variable  $f_2$ ?

A: more sophisticated graph;

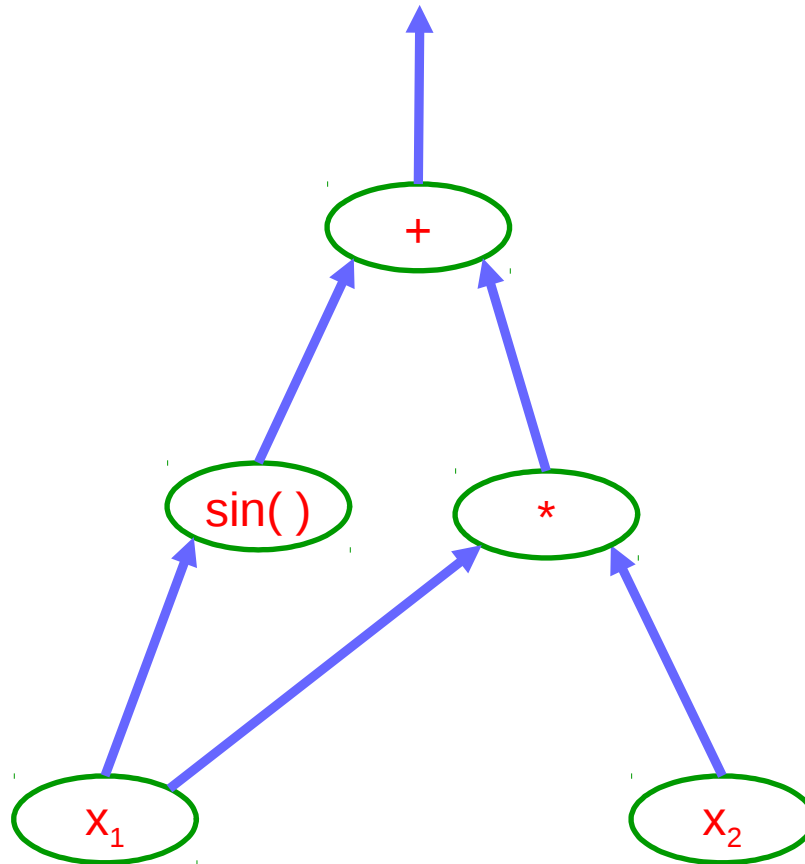
single "forward prop"





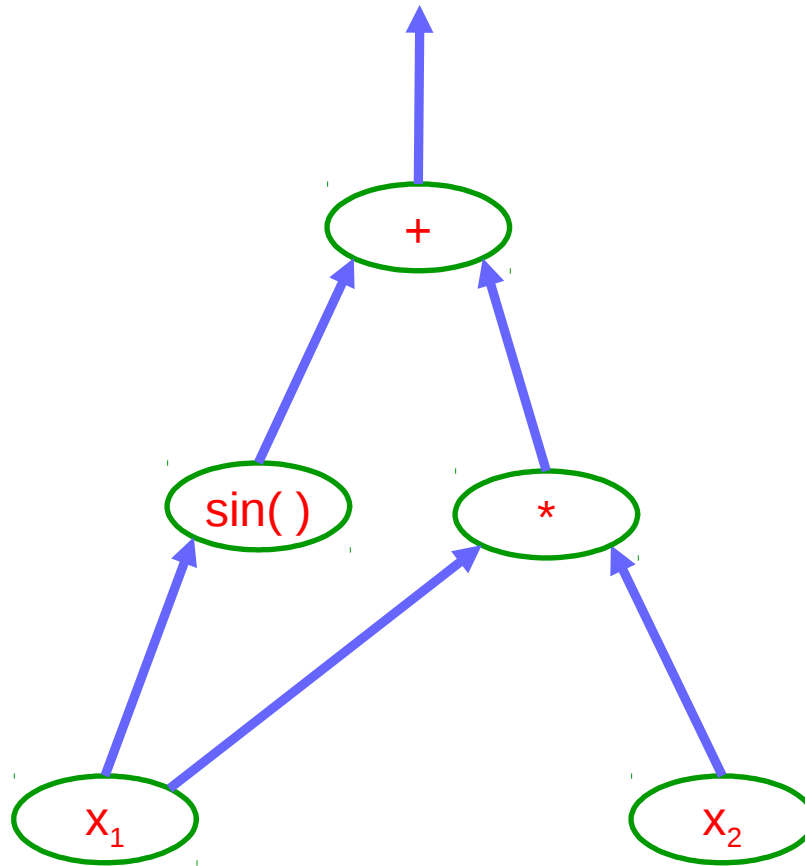
# Example: Reverse mode AD

$$f(x_1, x_2) = \sin(x_1) + x_1x_2$$



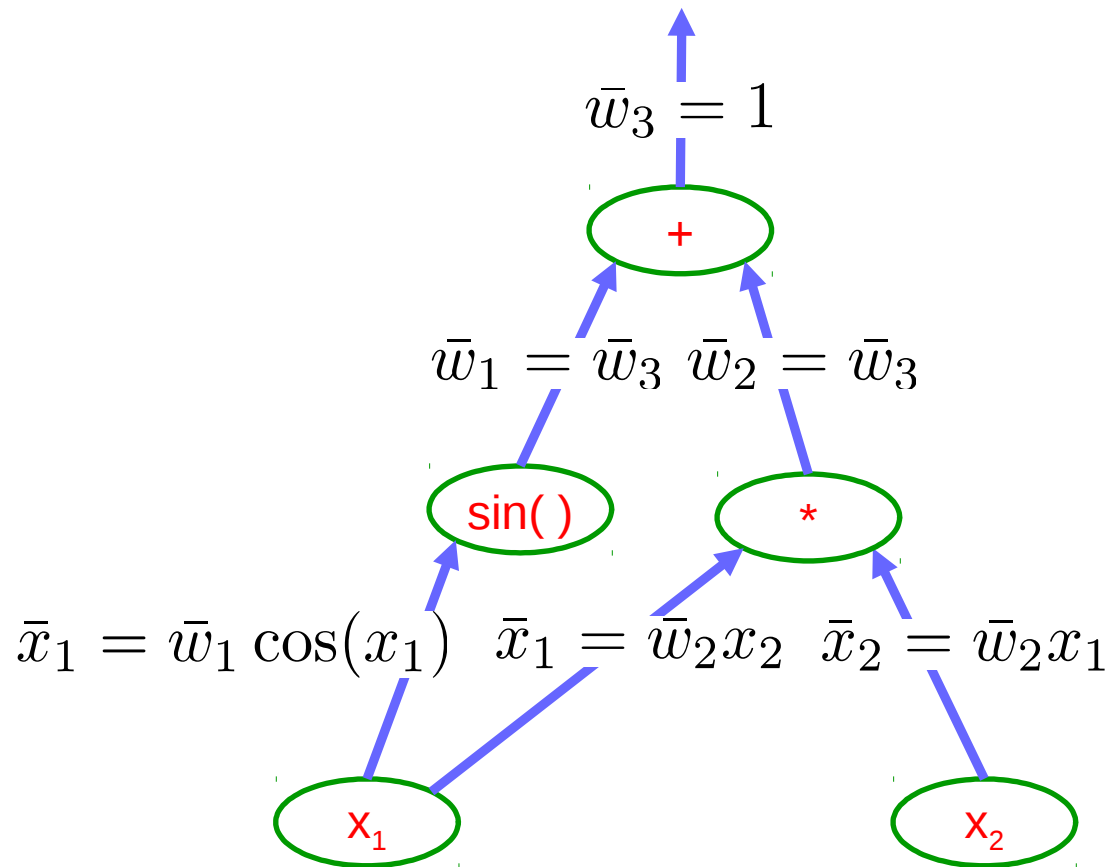
# Example: Reverse mode AD

$$f(x_1, x_2) = \sin(x_1) + x_1x_2$$

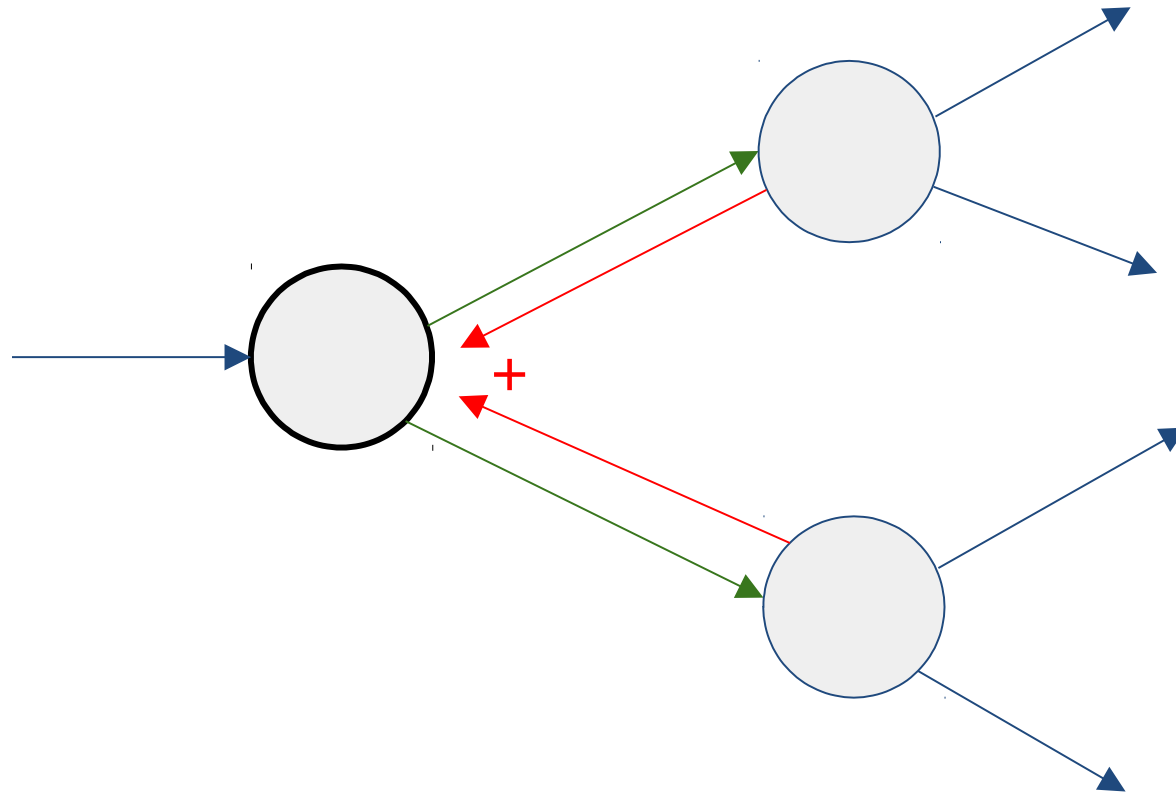


# Example: Reverse mode AD

$$f(x_1, x_2) = \sin(x_1) + x_1x_2$$



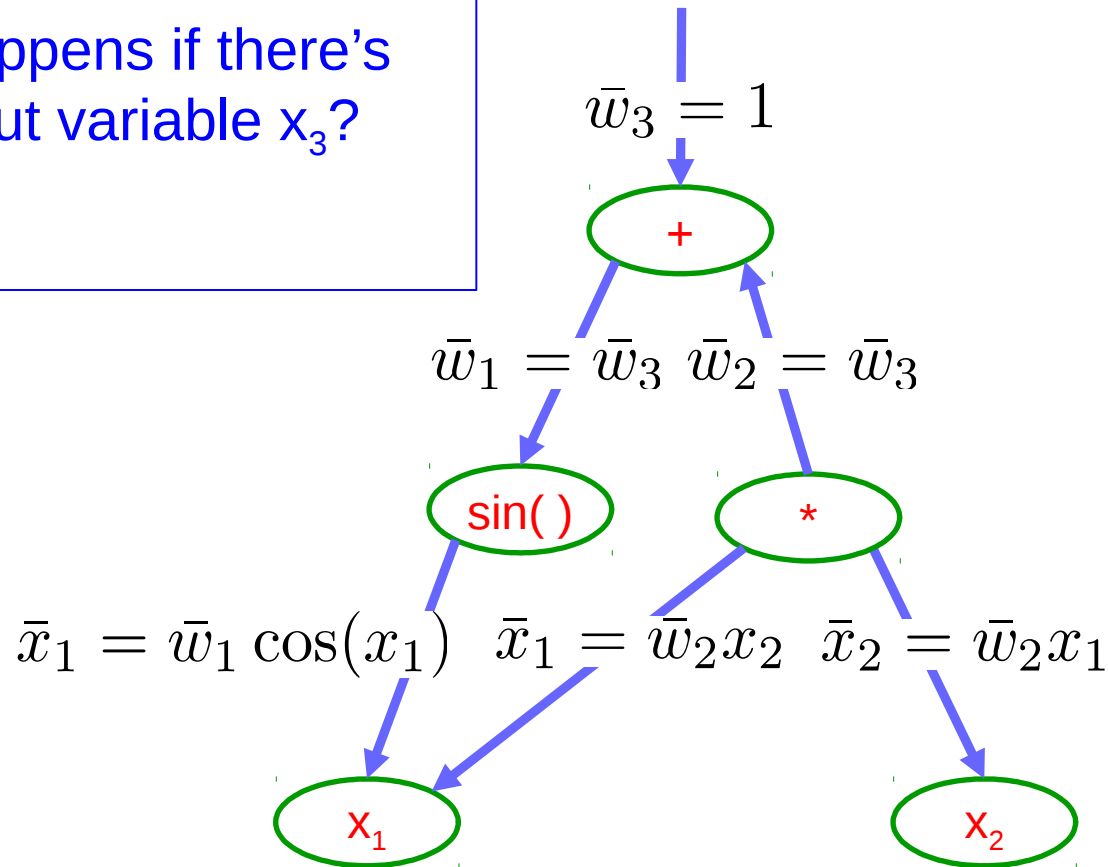
# Gradients add at branches



# Example: Reverse mode AD

$$f(x_1, x_2) = \sin(x_1) + x_1x_2$$

Q: What happens if there's another input variable  $x_3$ ?

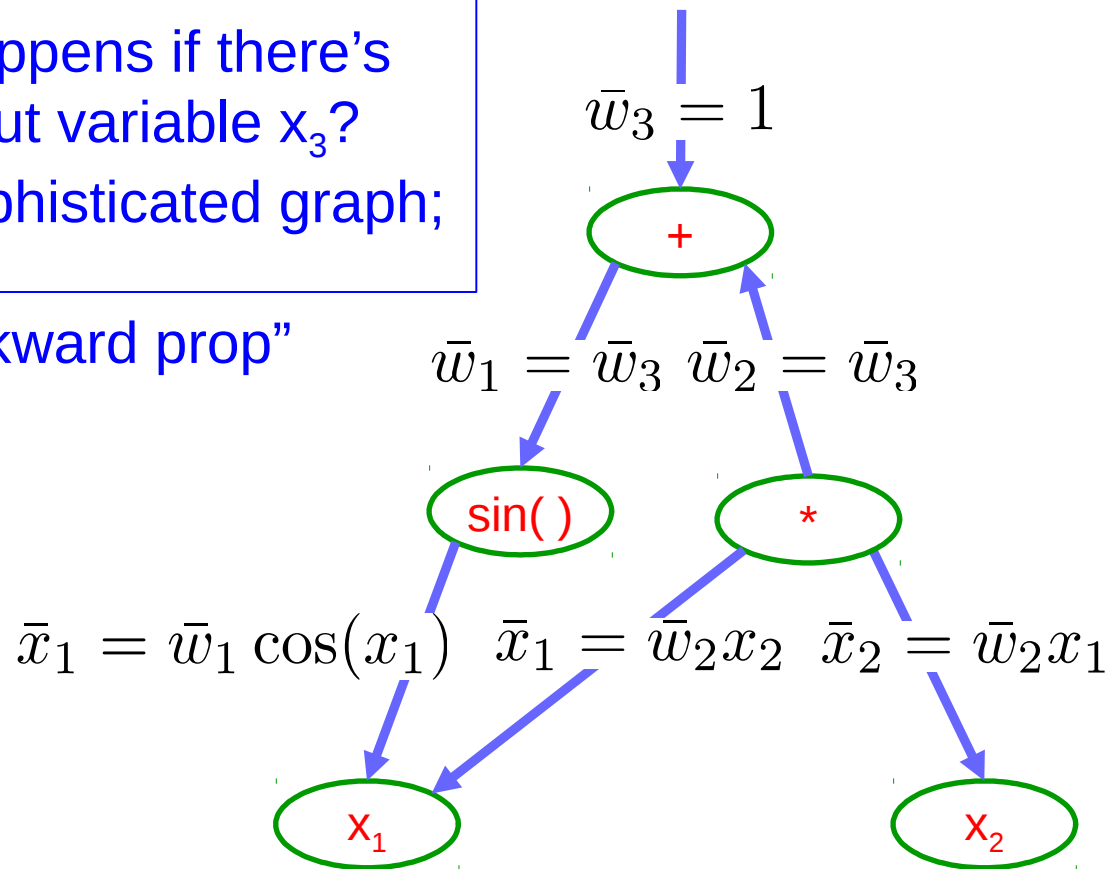


# Example: Reverse mode AD

$$f(x_1, x_2) = \sin(x_1) + x_1x_2$$

Q: What happens if there's another input variable  $x_3$ ?  
A: more sophisticated graph;

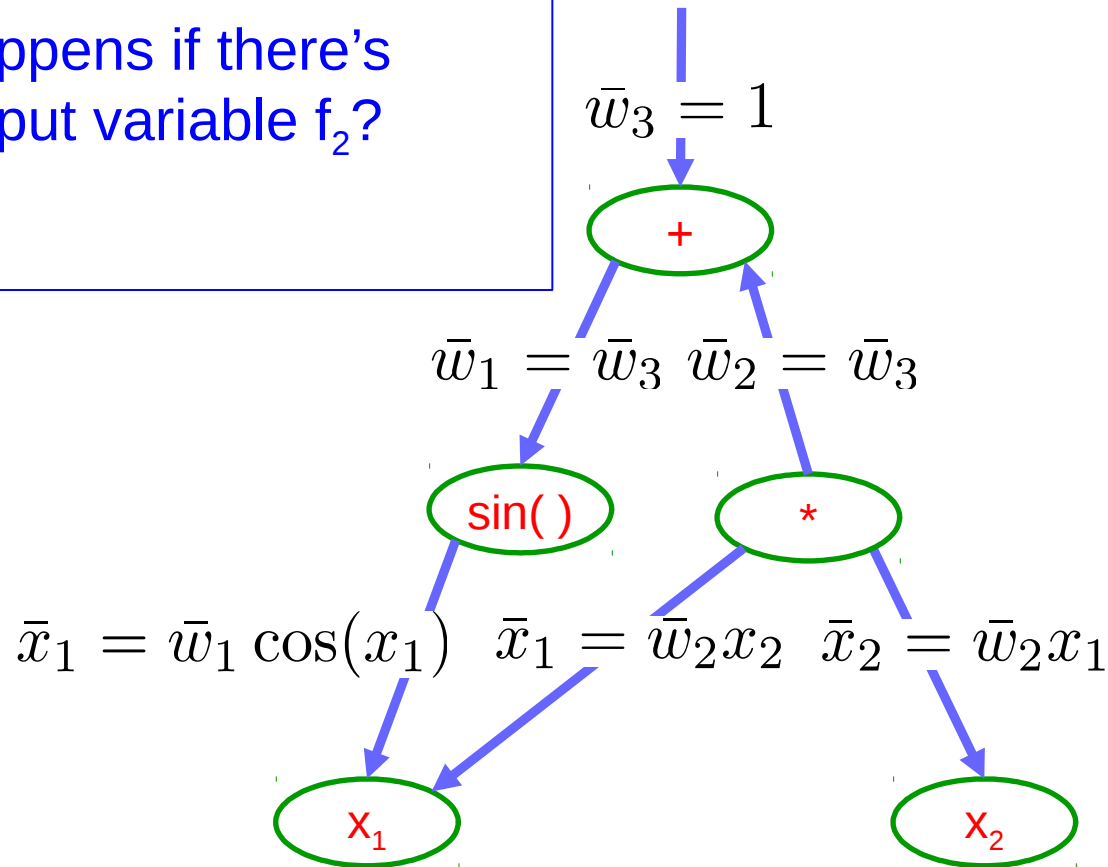
single "backward prop"



# Example: Reverse mode AD

$$f(x_1, x_2) = \sin(x_1) + x_1x_2$$

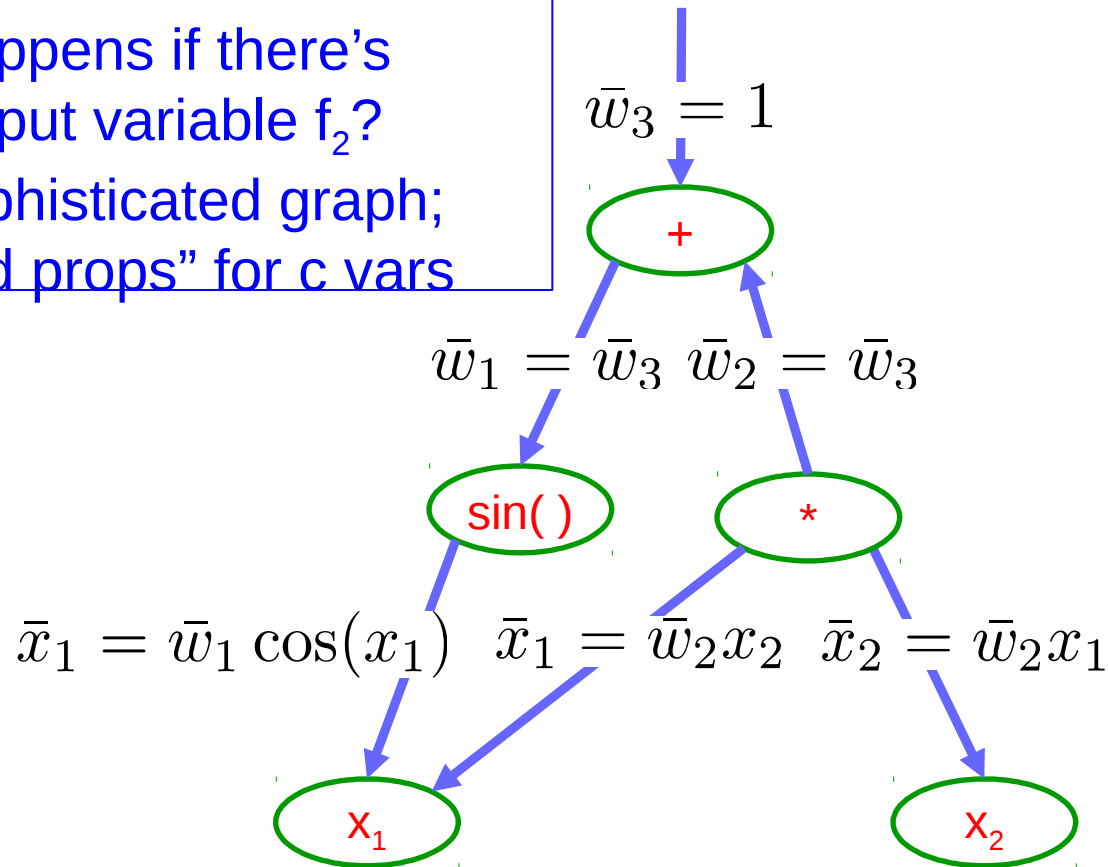
Q: What happens if there's another output variable  $f_2$ ?



# Example: Reverse mode AD

$$f(x_1, x_2) = \sin(x_1) + x_1x_2$$

Q: What happens if there's another output variable  $f_2$ ?  
A: more sophisticated graph; c "backward props" for c vars

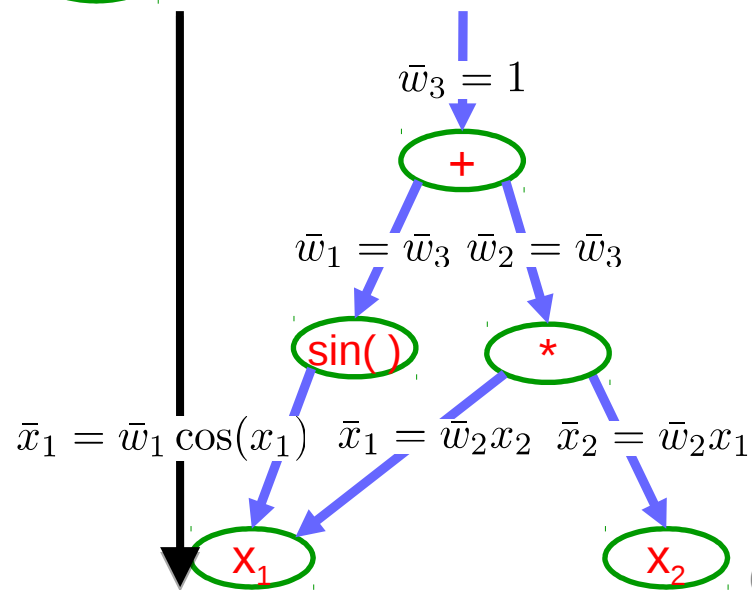
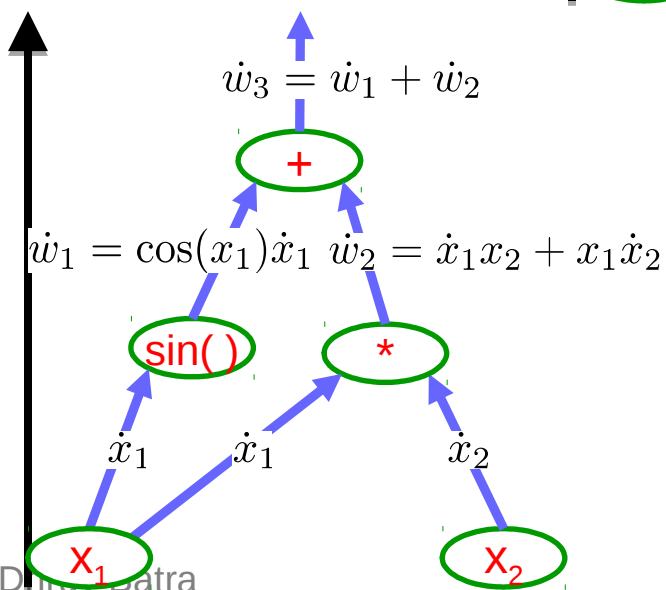
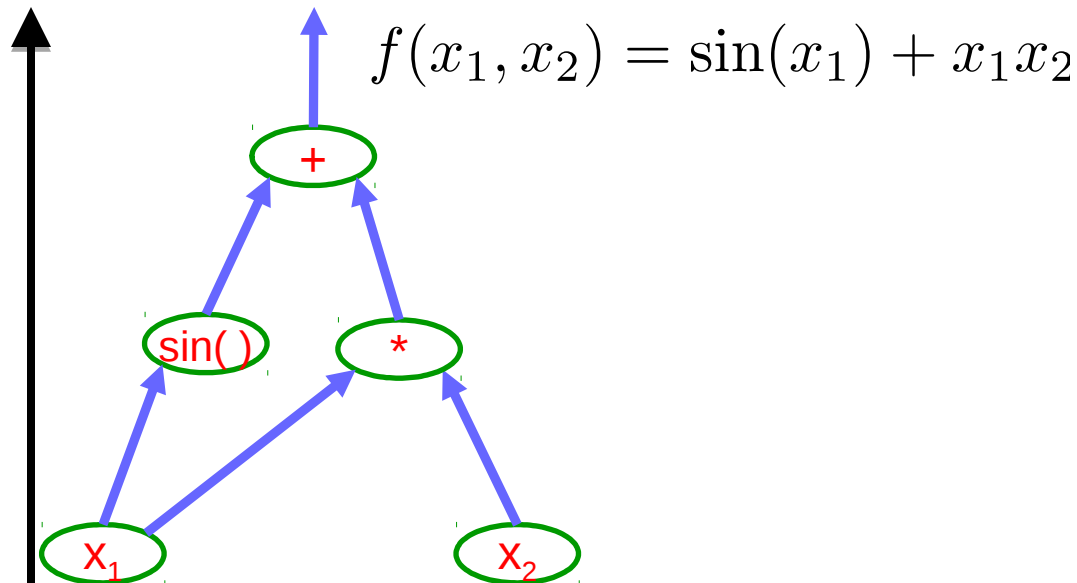




# Forward mode vs Reverse Mode

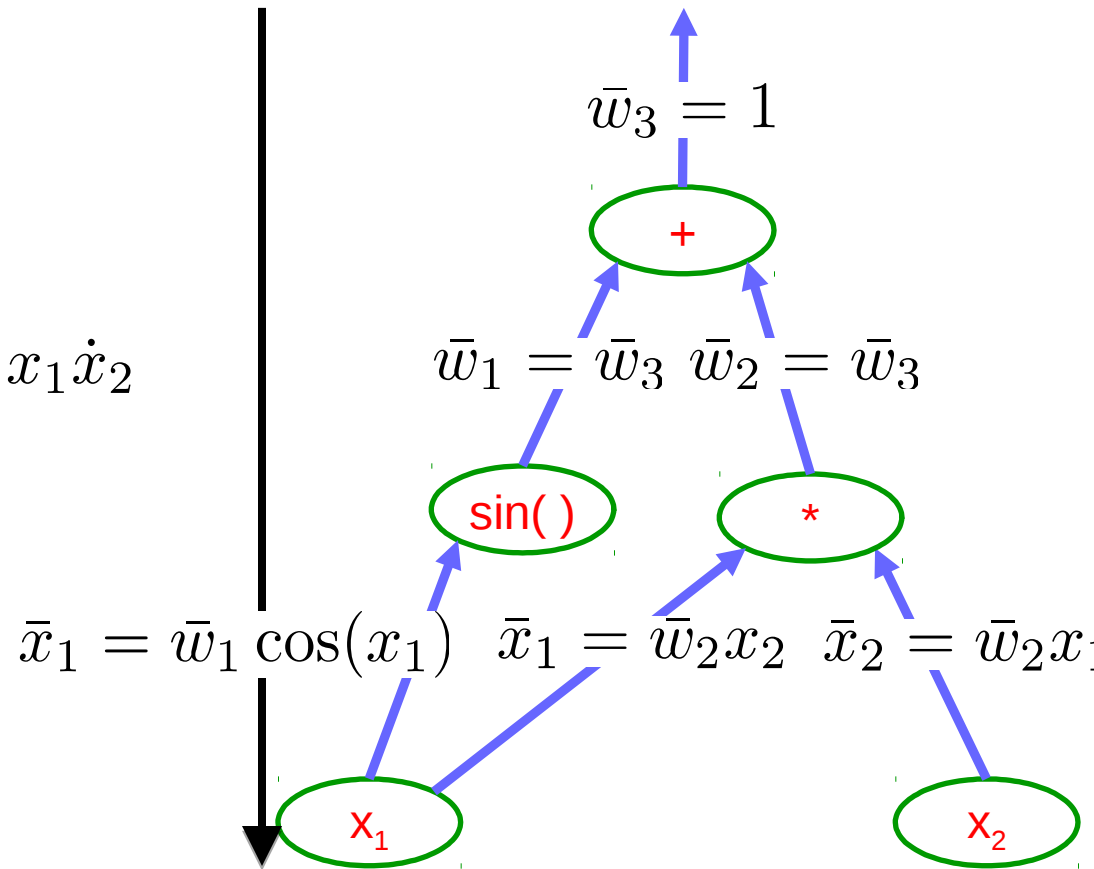
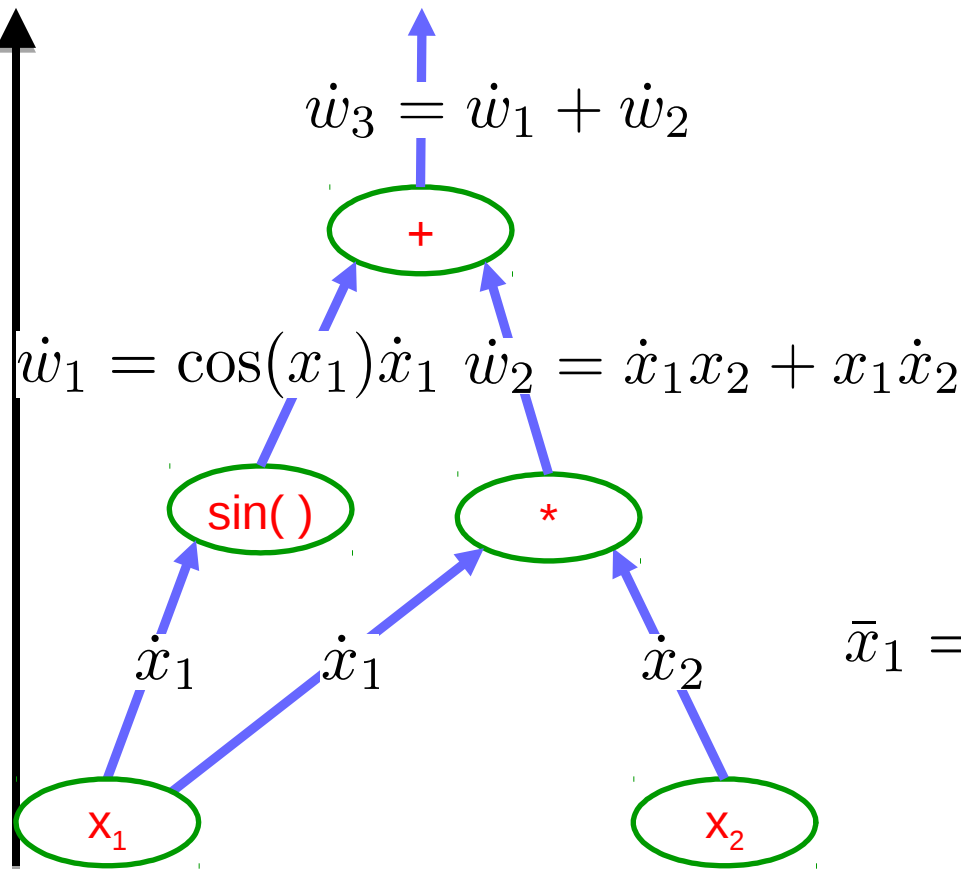
- $x \rightarrow \text{Graph} \rightarrow L$
- Intuition of Jacobian

# Forward Pass vs Forward mode AD vs Reverse Mode AD



# Forward mode vs Reverse Mode

- What are the differences?



# Forward mode vs Reverse Mode

- What are the differences?
- Which one is faster to compute?
  - Forward or backward?

# Forward mode vs Reverse Mode

- What are the differences?
- Which one is faster to compute?
  - Forward or backward?
- Which one is more memory efficient (less storage)?
  - Forward or backward?