# CS 4803 / 7643: Deep Learning

Topics:
- Automatic Differentiation
    - (Finish) Forward mode vs Reverse mode AD
    - Patterns in backprop

Dhruv Batra

Georgia Tech

# Administrativia

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

- HW2 out on 9/10
  - Schedule: https://www.cc.gatech.edu/classes/AY2021/cs7643_fall/

- Project discussion next class

# Recap from last time

# How do we compute gradients?

- Analytic or "Manual" Differentiation

- Symbolic Differentiation

- Numerical Differentiation

- Automatic Differentiation
  - Forward mode AD
  - Reverse mode AD
    - aka "backprop"

# Vector/Matrix Derivatives Notation

$x \to$

$y \downarrow$

| | S | V | M |
|---|---|---|---|
| S | $\frac{\partial y}{\partial x}$ | $\frac{\partial y}{\partial x}$ | $\frac{\partial y}{\partial X}$ |
| V | $\frac{\partial \vec{y}}{\partial x}$ | $\frac{\partial \vec{y}}{\partial \vec{x}}$ | |
| M | $\frac{\partial Y}{\partial x}$ | | |

tensors

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

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

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

$\downarrow num = \dfrac{dim\ 1}{rows}$

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

$\longrightarrow den = \dfrac{dim\ 2}{col}$

$\dfrac{\partial y}{\partial \vec{x}} \quad \dfrac{\partial y}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & \cdots & \frac{\partial y}{\partial x_d} \end{bmatrix}$

# Vector/Matrix Derivatives Notation

$$\frac{\partial \vec{y}}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial y_i}{\partial x_1} & \cdots & \frac{\partial y_i}{\partial x_j} & \cdots & \frac{\partial y_i}{\partial x_d} \\ & & \frac{\partial y_c}{\partial x_j} & & \end{bmatrix}$$

$c \times d$

Jacobian

# Vector Derivative Example

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

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

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

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

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

# Extension to Tensors

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

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

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

y-vec = $Y(:)$

x-vec = $X(:)$

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

Jacobian "Matrix"

# Chain Rule: Composite Functions

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

$$\frac{L(\vec{w})}{f(x)} = g_\ell(g_{\ell-1} - - - g_1(x))$$

— layers in NN

$$= (g_\ell \circ g_{\ell-1} \circ \cdots g_1)(x)$$

# Chain Rule: Scalar Case

$$x \in \mathbb{R}^1 \xrightarrow{g_1()} z \in \mathbb{R}^1 \xrightarrow{g_2()} y \in \mathbb{R}^1$$

$$y = g_2(\underbrace{g_1(x)}_{z})$$

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

Scalar mult.

# Chain Rule: Vector Case

$$\vec{x} \in \mathbb{R}^d \xrightarrow{\;\;g_1(\cdot)\;\;} \vec{z} \in \mathbb{R}^m \xrightarrow{\;\;g_2(\cdot)\;\;} \vec{y} \in \mathbb{R}^c$$

$$: \mathbb{R}^d \to \mathbb{R}^m \qquad\qquad : \mathbb{R}^m \to \mathbb{R}^c$$

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

$$J_{g_2 \circ g_1} \qquad\qquad J_{g_2} \qquad \text{Matrix Mult} \qquad J_{g_1}$$

# Chain Rule: <u>Jacobian</u> <u>view</u>

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

$$\frac{\partial \vec{y}}{\partial \vec{z}}$$

$$\frac{\partial \vec{z}}{\partial \vec{x}}$$

$$\left[ \quad \frac{\partial y_i}{\partial x_j} \quad \right] \qquad \left[ \quad \frac{\partial y_i}{\partial z_k} \quad \right] \qquad \left[ \quad \frac{\partial z_k}{\partial x_j} \quad \right]$$

$$\frac{\partial y_i}{\partial x_j} = \sum_k \frac{\partial y_i}{\partial z_k} \frac{\partial z_k}{\partial x_j}$$

# Chain Rule: Graphical view

$$\vec{x} = \vec{h}^0 \quad \boxed{g_1} \rightarrow \vec{h}^1 \in \mathbb{R}^d \quad \xrightarrow{g(\cdot)} \vec{h}^{(2)} \in \mathbb{R}^d \cdots \cdots \xrightarrow{g_{\ell \circ \cdots}} \boxed{\vec{h}^{(\ell)}} \rightarrow \boxed{L}$$

$$\begin{bmatrix} 0 \\ \vdots \\ 6 \\ 6 \end{bmatrix} \xrightarrow{g_1} \begin{bmatrix} 6 \\ 0 \\ \vdots \\ 6 \end{bmatrix} \cdots \cdots \begin{bmatrix} \vdots \\ 6 \end{bmatrix} \begin{bmatrix} \vdots \\ 1 \end{bmatrix} \rightarrow \boxed{L(\cdot)}$$

$$\frac{d}{d\vec{h}^0} = \boxed{\frac{\partial L}{\partial \vec{h}^{\ell}} \frac{\partial \vec{h}^{(\ell)}}{\partial \vec{h}^{(1)}}} = \boxed{\frac{\partial L}{\partial \vec{h}^{(\ell)}} \frac{\partial \vec{h}^{(\ell)}}{\partial \vec{h}^{\ell-1}}} \boxed{\frac{\partial \vec{h}^{\ell-1}}{\partial \vec{h}^{(\ell-2)}}} \cdots \cdots \boxed{\frac{\partial \vec{h}^{(2)}}{\partial \vec{h}^{(1)}}}$$

$$[\,] = [\quad] \begin{bmatrix} \\ \end{bmatrix}_{d\times d} = [\quad] \begin{bmatrix} \\ \end{bmatrix} \begin{bmatrix} \\ \end{bmatrix} \begin{bmatrix} \\ \end{bmatrix}_{d\times d}$$

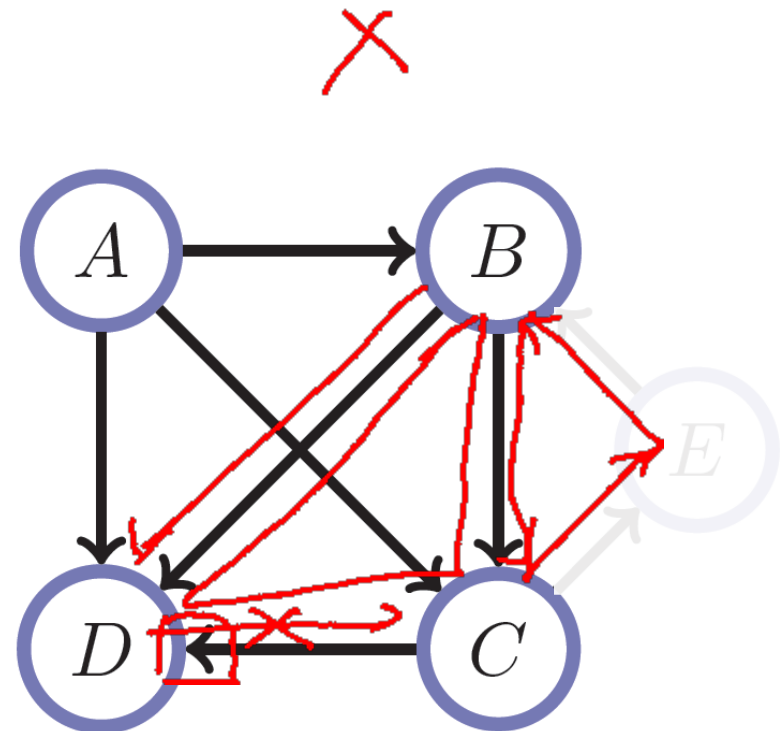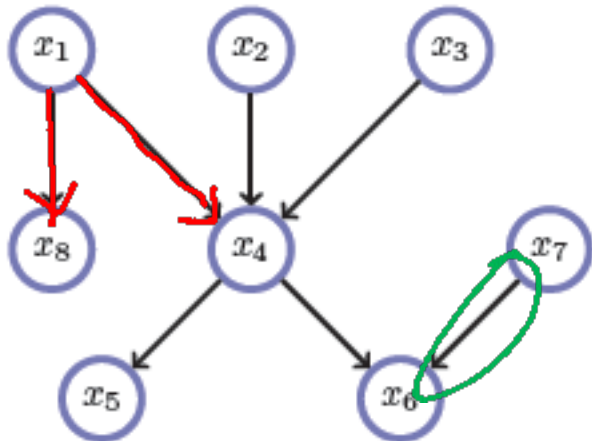$$\underbrace{O(\ell d^2)} \longleftrightarrow \qquad \underbrace{O(\ell d^3)}$$

# 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

# Directed Acyclic Graphs (DAGs)

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

DAG

# Directed Acyclic Graphs (DAGs)

- Concept
  - Topological Ordering

$$\exists \text{ bijection } \sigma : V \to \{1, \dots, n\}$$

$$\text{s.t } \forall (v_i, v_j) \in E$$

$$\sigma(v_i) < \sigma(v_j)$$

# Computational Graphs ← DAG

- Notation

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

$w_2$   $w_1$

# 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 AD

Goal: $\frac{\partial L}{\partial x}$

$z$

$\overset{\rightarrow (\ell-1)}{h}$

$a$

$\overset{\rightarrow (\ell)}{h} = g(h^{\ell-1})$

$g(\cdot)$

Output

$\frac{\partial L}{\partial \overset{\rightarrow}{h}^{\ell-1}} = \frac{\partial L}{\partial \overset{\rightarrow}{h}^{\ell}} \cdot \frac{\partial \overset{\rightarrow}{h}^{\ell}}{\partial \overset{\rightarrow}{h}^{\ell-1}}$

Input: $\frac{\partial L}{\partial \overset{\rightarrow}{h}^{\ell}}$

Input    Layer Jacobian

# Plan for Today

- Automatic Differentiation
  - (Finish) Forward mode vs Reverse mode AD
  - Backprop
  - Patterns in backprop

# Example: Forward mode AD

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

$$\frac{\partial f}{\partial x} = \left[ \frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \right]$$

$$\frac{\partial f}{\partial x_1}$$
$$= \cos x_1 + x_2$$

$$w_1 = \sin(x_1)$$

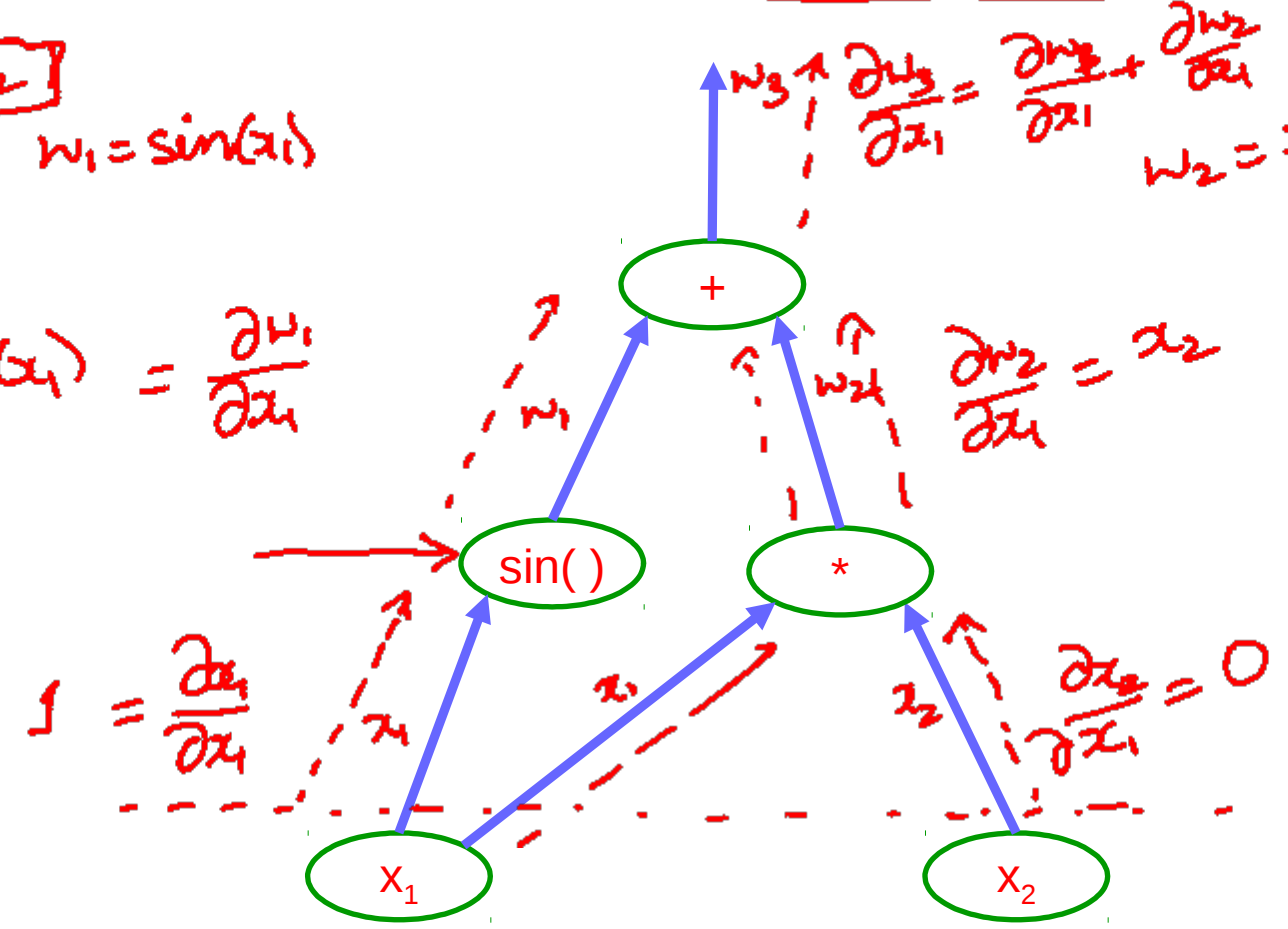$$w_3 \uparrow \quad \frac{\partial w_3}{\partial x_1} = \frac{\partial w_1}{\partial x_1} + \frac{\partial w_2}{\partial x_1}$$

$$w_2 = x_1 x_2$$

$$\cos(x_1) = \frac{\partial w_1}{\partial x_1}$$

$$+$$

$$\frac{\partial w_2}{\partial x_1} = x_2$$

$$w_1 \qquad w_2 \uparrow$$

$$\text{sin( )} \qquad *$$

$$1 = \frac{\partial x_1}{\partial x_1}$$

$$x_1 \qquad x_1 \qquad x_2 \qquad \frac{\partial x_2}{\partial x_1} = 0$$

$$X_1 \qquad X_2$$

# Example: Forward mode AD

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

$\dfrac{\partial f}{\partial x_2}$

$= x_1$

$\dfrac{\partial v_3}{\partial x_2} = x_1$

$\dfrac{\partial u}{\partial x_2} = 0$

$\dfrac{\partial w_1}{\partial x_2} = x_1$

+

sin( )    *

$0 = \dfrac{\partial x_1}{\partial x_2}$    $x_1$    $\dfrac{\partial x_2}{\partial x_2} = 1$

$x_1$    $x_2$

# Example: Forward mode AD

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix} f(x_1, x_2) = \sin(x_1) + x_1 x_2$$
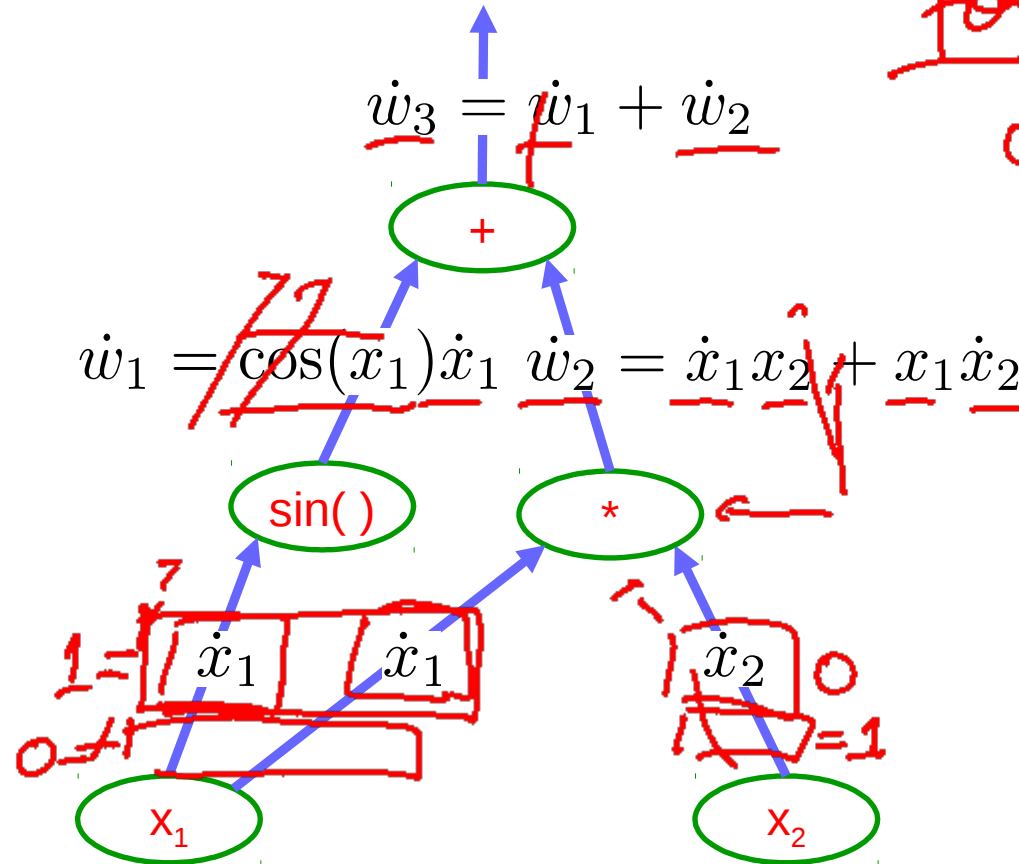
for a ind $x_1, x_2$

compute/ evaluate

$\dot{w}_3 = \frac{\partial w_3}{\partial a}$

$\dot{w}_2 = \frac{\partial w_2}{\partial a}$

$\dot{w}_1 = \frac{\partial w_1}{\partial a}$

$\dot{x}_1 = \frac{\partial x_1}{\partial a}$

$\dot{x}_2 = \frac{\partial x_2}{\partial a}$

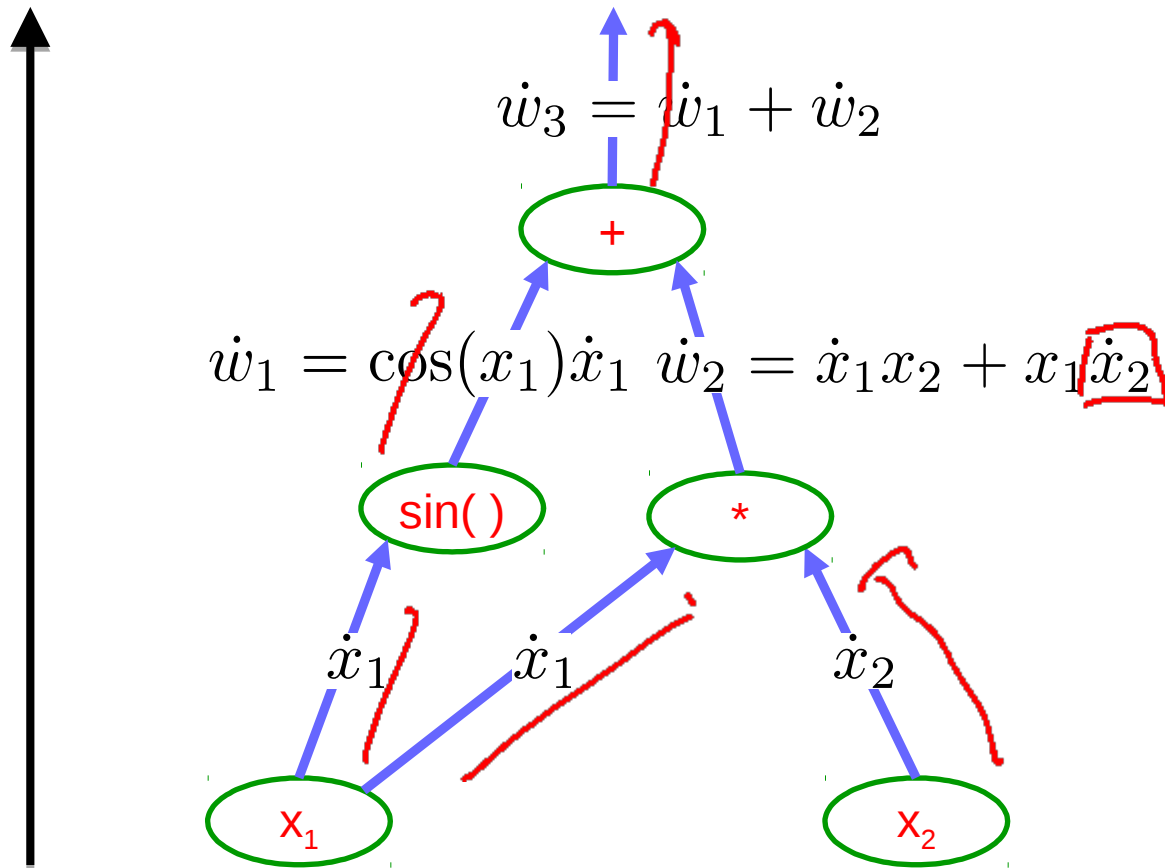$\dot{w}_3 = \dot{w}_1 + \dot{w}_2$

(+)

$\dot{w}_1 = \cos(x_1)\dot{x}_1 \quad \dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$

sin( )     *

$\dot{x}_1$     $\dot{x}_1$     $\dot{x}_2$

$1 =$     $0 =$     $= 1$

$X_1$     $X_2$

# Example: Forward mode AD

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

$$\frac{\partial f}{\partial x}$$

$$\dot{w}_3 = \dot{w}_1 + \dot{w}_2$$

$$+$$

$$\dot{w}_1 = \cos(x_1)\dot{x}_1 \quad \dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$$

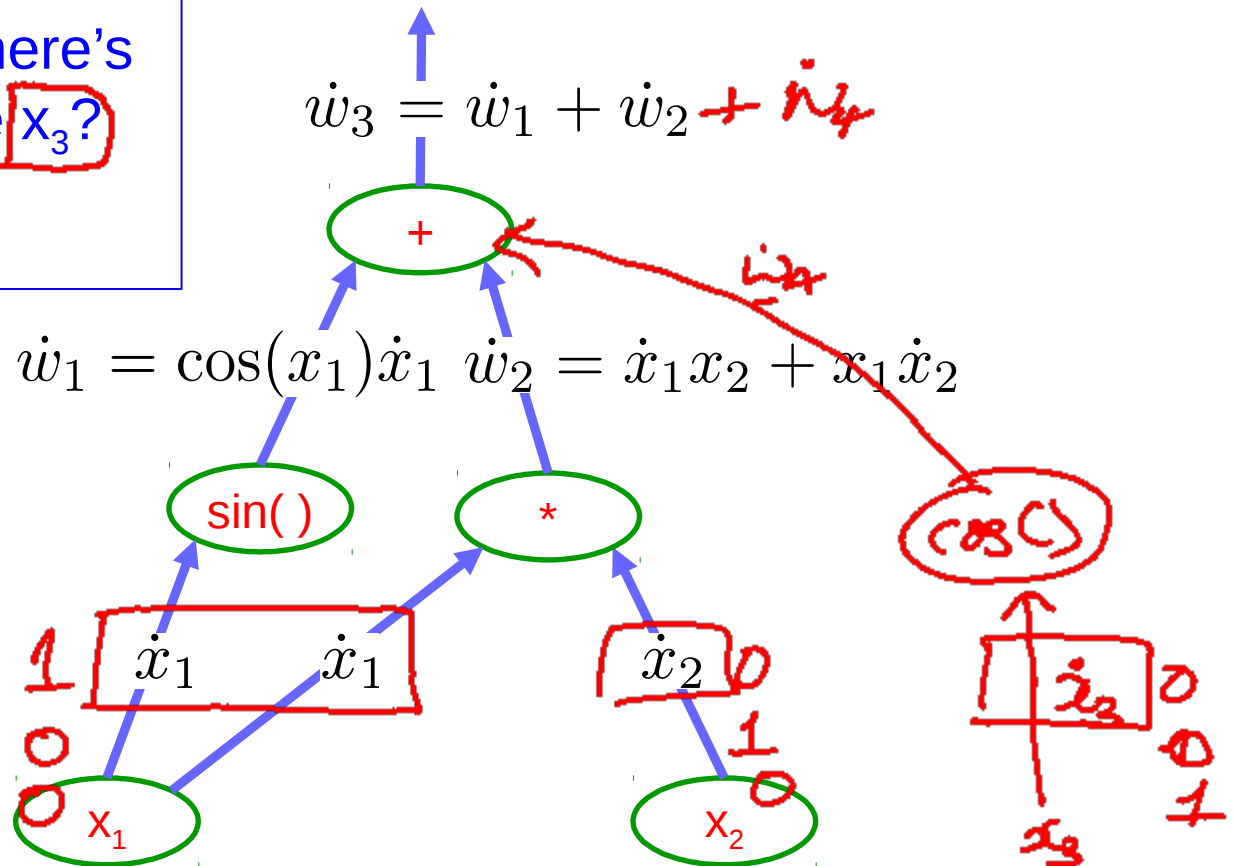sin( )   *

$$\dot{x}_1 \qquad \dot{x}_1 \qquad \dot{x}_2$$

$X_1$   $X_2$

# Example: Forward mode AD

$$f(x_1, x_2) = \sin(x_1) + x_1 x_2 + \cos(x_3)$$

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

$$\dot{w}_3 = \dot{w}_1 + \dot{w}_2 + \dot{w}_4$$

$$\dot{w}_1 = \cos(x_1)\dot{x}_1 \quad \dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$$

for a in $\{x_1, x_2, x_3\}$

$+$

$\dot{w}_4$

sin( )   *   cos( )

$\dot{x}_1$   $\dot{x}_1$   $\dot{x}_2$   $0$   $\dot{x}_3$   $0$
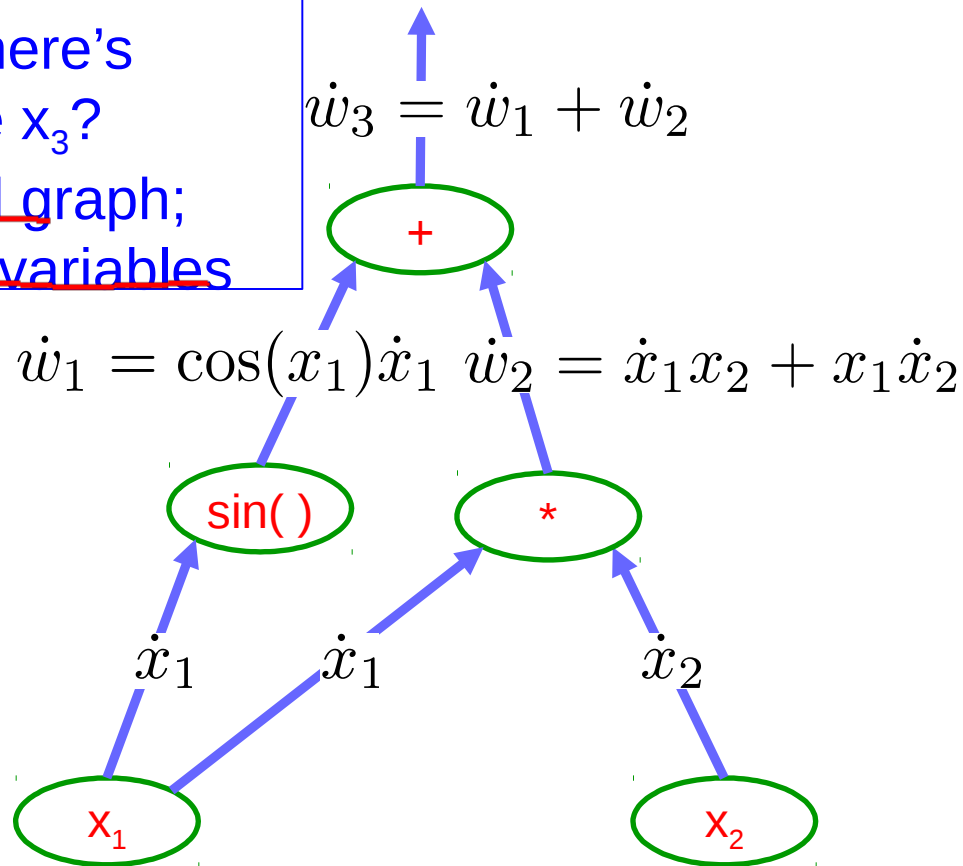
$1$   $0$   $1$   $0$

$0$   $0$

$x_1$   $x_2$   $x_3$

# Example: Forward mode AD

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

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

A: more sophisticated graph; d+1 "passes" for d+1 variables

$$\dot{w}_3 = \dot{w}_1 + \dot{w}_2$$

+

$$\dot{w}_1 = \cos(x_1)\dot{x}_1 \quad \dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$$

sin( )     *

$$\dot{x}_1 \quad\quad \dot{x}_1 \quad\quad\quad \dot{x}_2$$

$x_1$       $x_2$

# Example: Forward mode AD

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

$\dfrac{\partial \vec{f}}{\partial \vec{x}}$

$f_2 = \cos(x_2)$

$f_2 = w_4$

$\dot{w}_4$

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

$\dot{w}_3 = \dot{w}_1 + \dot{w}_2$

$\boxed{+}$

$\dot{w}_1 = \cos(x_1)\dot{x}_1 \quad \dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$

for a in $(x_1, x_2)$

$\sin(\ )$     $*$     $\cos(\ )$

1
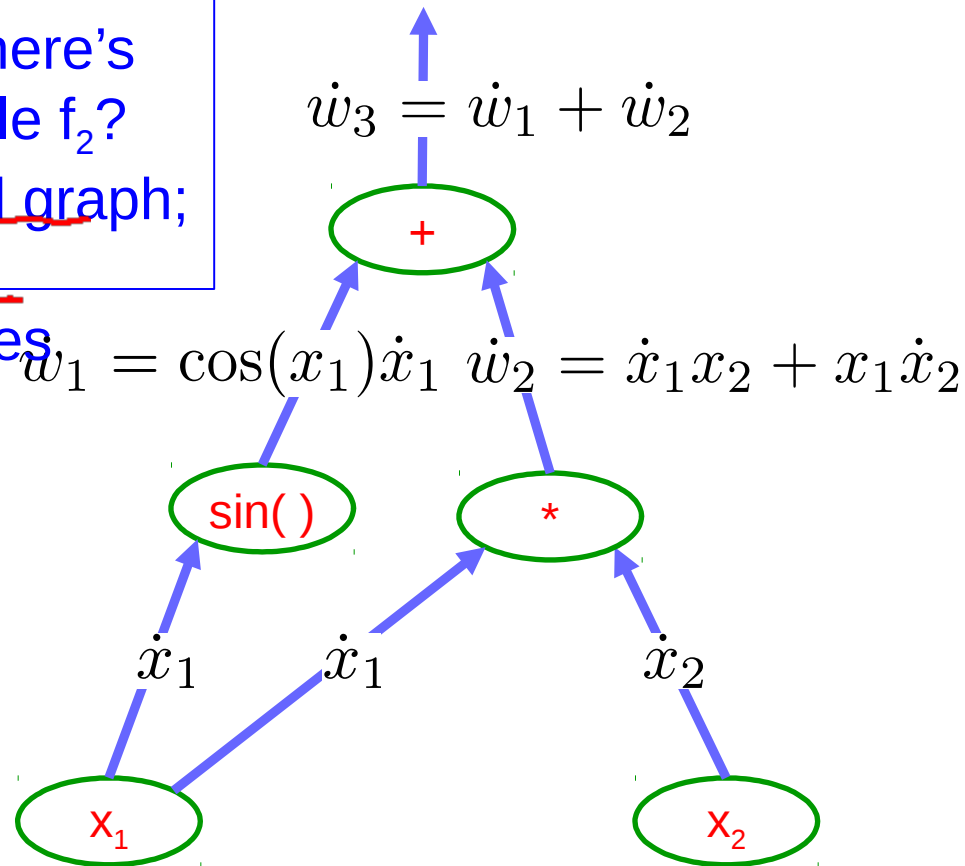0    $\dot{x}_1$    $\dot{x}_1$    0   $\dot{x}_2$    $\dot{x}_2$

1

$X_1$       $X_2$

# Example: Forward mode AD

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

Q: What happens if there's another output variable f₂?
A: more sophisticated graph;

d "passes" for variables

$$\dot{w}_3 = \dot{w}_1 + \dot{w}_2$$

+

$$\dot{w}_1 = \cos(x_1)\dot{x}_1 \quad \dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$$

sin( )   *

$$\dot{x}_1 \quad \dot{x}_1 \quad \dot{x}_2$$

X₁   X₂

# Example: Reverse mode AD

Goal:

$$\frac{\partial f}{\partial x} = \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2}\right]$$

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

$$w_3 = w_1 + w_2$$

$$\frac{\partial w_3}{\partial w_1} = 1$$

$$\bar{w}_3 = \frac{\partial f}{\partial w_3} = \frac{\partial w_3}{\partial w_3} = 1$$

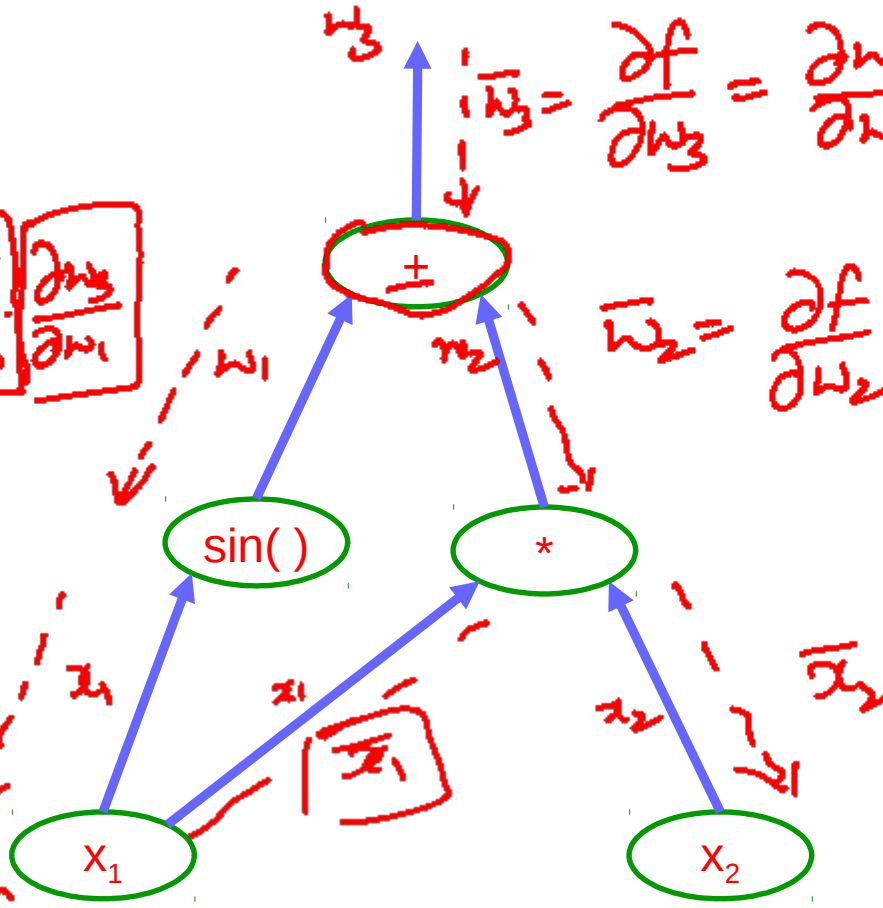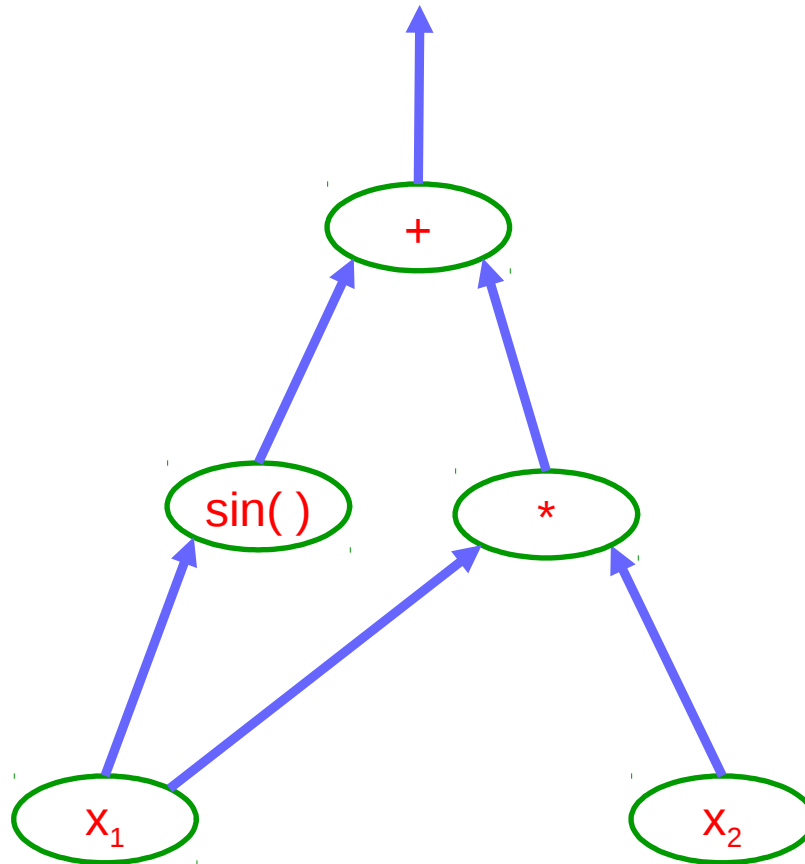$$\bar{w}_1 = \frac{\partial f}{\partial w_1} = \boxed{\frac{\partial f}{\partial w_3}} \boxed{\frac{\partial w_3}{\partial w_1}}$$

$$\bar{w}_2 = \frac{\partial f}{\partial w_2} = \boxed{\frac{\partial f}{\partial w_3}}$$

by def. $= \frac{\partial f}{\partial w_3}$

$$\boxed{\bar{x}_1} = \frac{\partial f}{\partial w_1} = \boxed{\frac{\partial f}{\partial w_1}}\boxed{\frac{\partial w_1}{\partial u_1}}$$

input Jacobian $\frac{\partial}{\partial w_1}$ cos($x_1$)

$w_3$

$\pm$

$w_1$ $w_2$

sin( )  *

$x_1$  $x_1$  $x_2$  $\bar{x}_2$

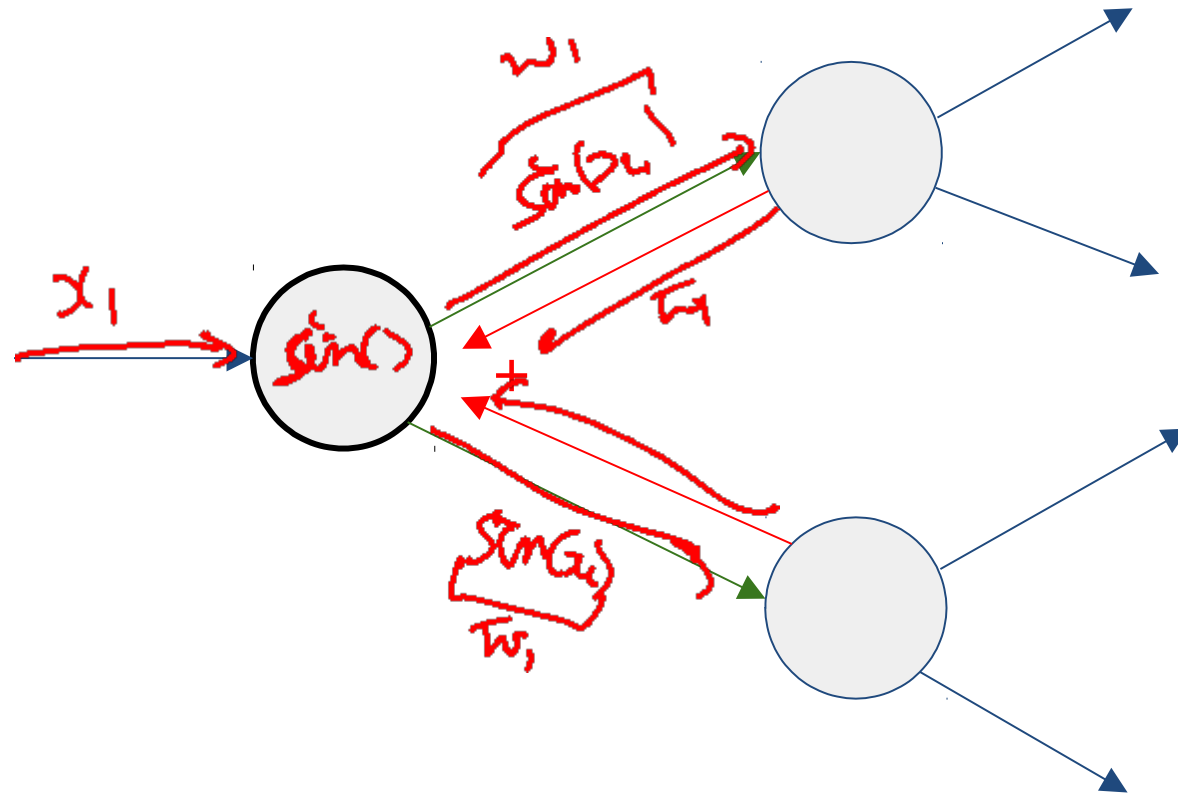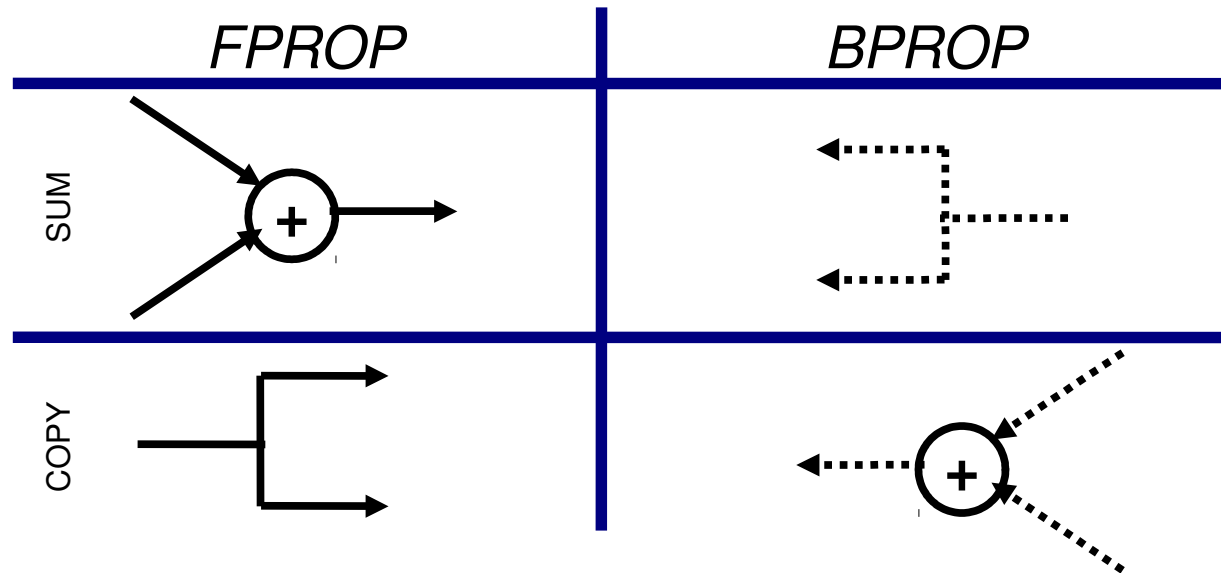$\bar{x}_1$

$X_1$  $X_2$

# Example: Reverse mode AD

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

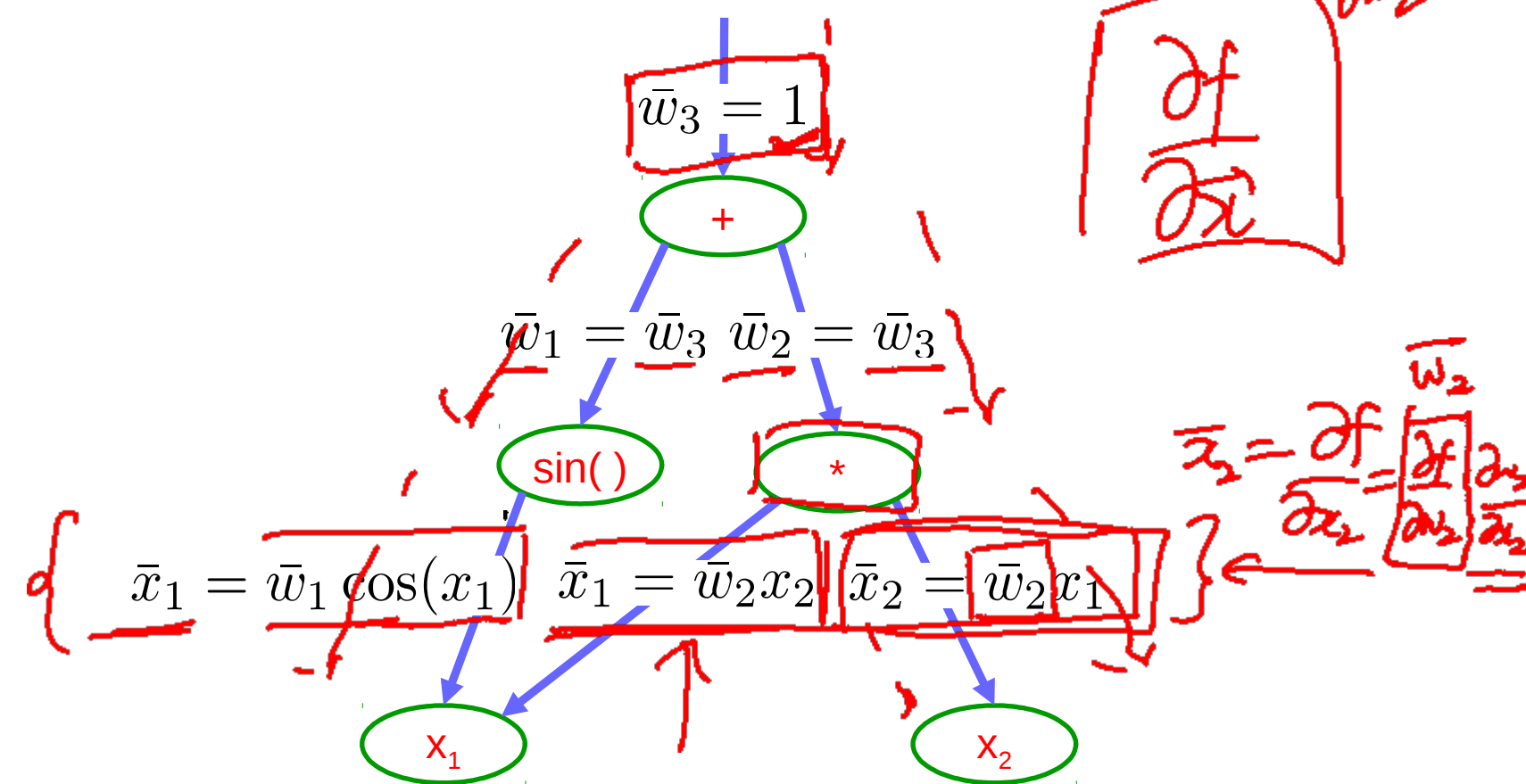# Gradients add at branches

# Duality in Fprop and Bprop

# Example: Reverse mode AD

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



$$\bar{w}_3 = 1$$

$$\bar{w}_1 = \bar{w}_3 \quad \bar{w}_2 = \bar{w}_3$$

sin( )       *

$$\bar{x}_1 = \bar{w}_1 \cos(x_1) \quad \bar{x}_1 = \bar{w}_2 x_2 \quad \bar{x}_2 = \bar{w}_2 x_1$$

$x_1$       $x_2$

$w_3 = x_1 x_2$

$$\frac{\partial w_3}{\partial x_2} = x_1$$

$$\frac{\partial f}{\partial x}$$

$$\bar{w}_2$$

$$\bar{x}_2 = \frac{\partial f}{\partial x_2} = \frac{\partial f}{\partial w_3} \frac{\partial w_3}{\partial x_2}$$

# Example: Reverse mode AD

$$f(x_1, x_2) = \sin(x_1) + x_1 x_2 \; {\color{red}+ \cos(x_3)}$$

$${\color{red}w_3 = w_1 + w_2 + w_4}$$

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

$$\bar{w}_3 = 1$$

$+$

$$\bar{w}_1 = \bar{w}_3 \quad \bar{w}_2 = \bar{w}_3$$

${\color{red}w_4}$

${\color{red}w_4 =}$

sin( )        *        ${\color{red}\cos(\ )}$

$$\bar{x}_1 = \bar{w}_1 \cos(x_1) \quad \bar{x}_1 = \bar{w}_2 x_2 \quad \bar{x}_2 = \bar{w}_2 x_1$$

${\color{red}\bar{x}_3}$

$x_1$        $x_2$        ${\color{red}x_3}$

# Example: Reverse mode AD

$$f(x_1, x_2) = \sin(x_1) + x_1 x_2$$
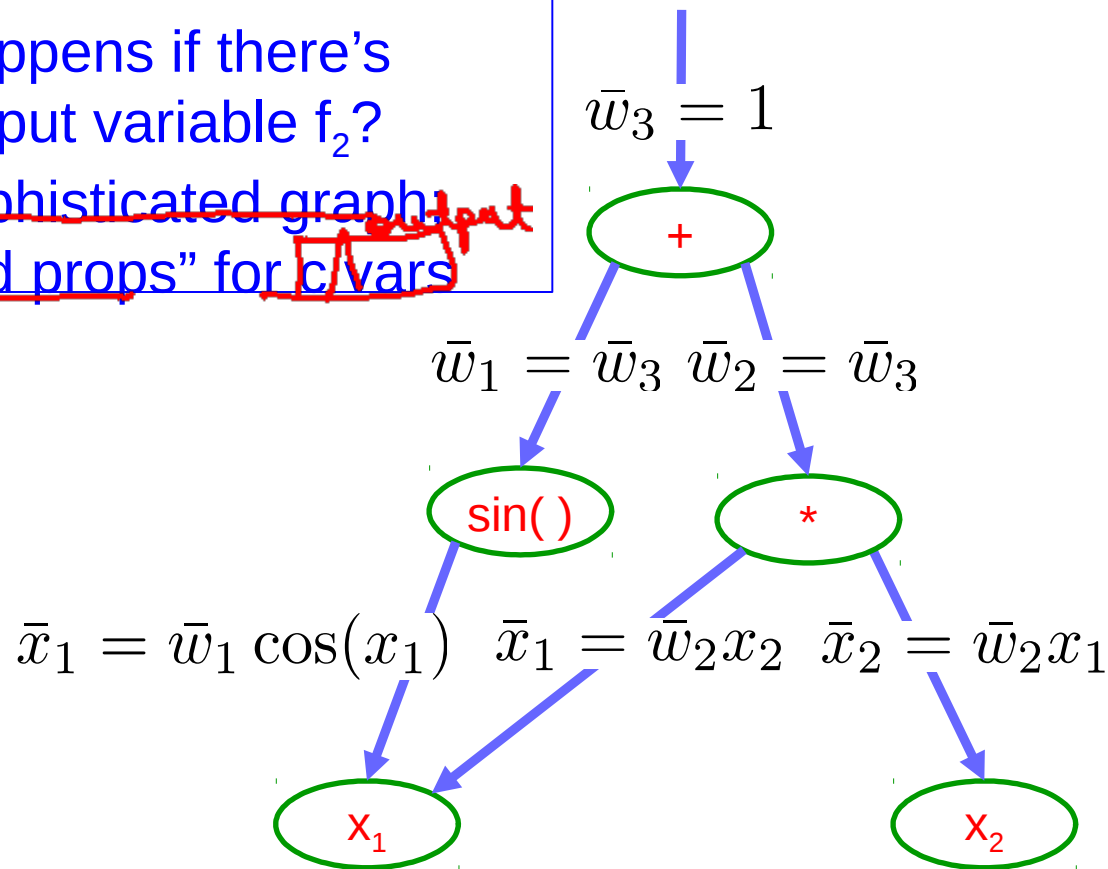
Q: What happens if there's another input variable $x_3$?

A: more sophisticated graph;

single "backward prop"

$$\bar{w}_3 = 1$$

( + )

$$\bar{w}_1 = \bar{w}_3 \quad \bar{w}_2 = \bar{w}_3$$

( sin( ) )   ( * )

$$\bar{x}_1 = \bar{w}_1 \cos(x_1) \quad \bar{x}_1 = \bar{w}_2 x_2 \quad \bar{x}_2 = \bar{w}_2 x_1$$

( $x_1$ )    ( $x_2$ )

# Example: Reverse mode AD

$$f_1(x_1, x_2) = \sin(x_1) + x_1 x_2 \qquad f_2 = \cos(x_2)$$

Q: What happens if there's another output variable f₂?

$$\bar{w}_3 = 1$$

$\bar{w}_4 = f_2$

$\bar{w}_4$

(+)

$$\bar{w}_1 = \bar{w}_3 \quad \bar{w}_2 = \bar{w}_3$$

for a in $\{f_1, f_2\}$

$\bar{w}_4 = \dfrac{\partial a}{\partial w_4}$

$\bar{w}_3 = \dfrac{\partial a}{\partial w_3}$

sin( )          *

(cos())

$$\bar{x}_1 = \bar{w}_1 \cos(x_1) \quad \bar{x}_1 = \bar{w}_2 x_2 \quad \bar{x}_2 = \bar{w}_2 x_1$$

$\bar{x}_1$

X₁          X₂

$$\bar{x}_1 = \dfrac{\partial f}{\partial x_1} = \dfrac{\partial f_2}{\partial w_2} \dfrac{\partial w_2}{\partial x_1} = \bar{w}_2 x_2$$

# Example: Reverse mode AD

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

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

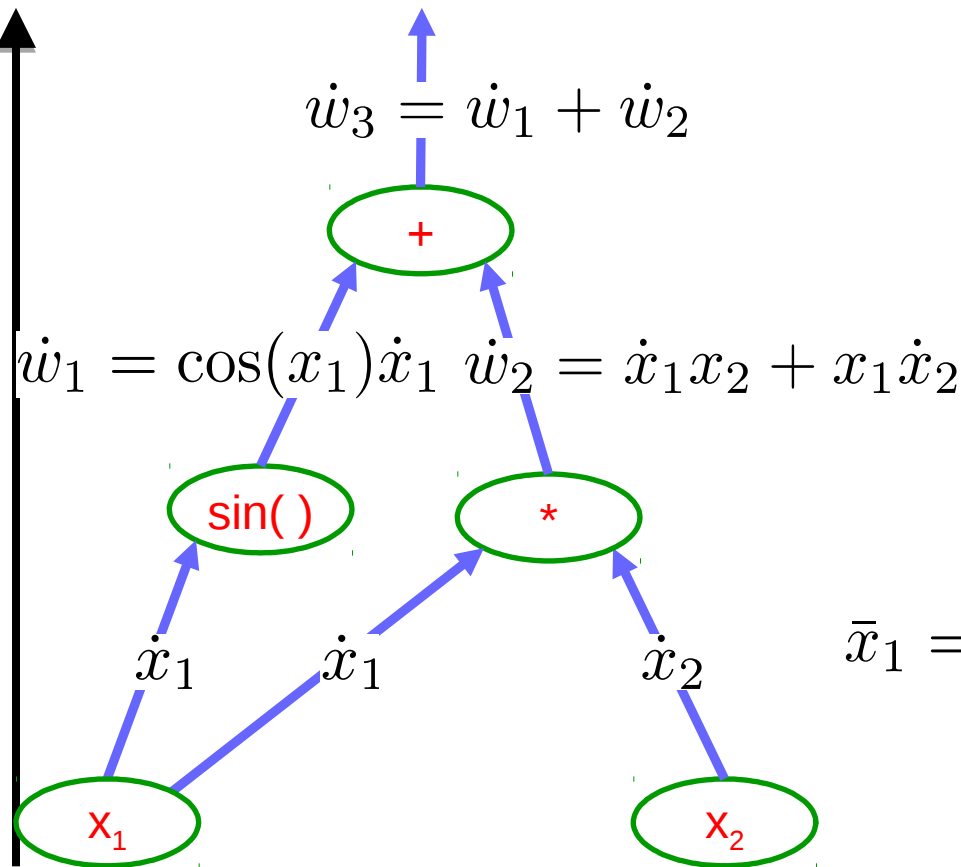A: more sophisticated graph; c "backward props" for c vars

$$\bar{w}_3 = 1$$

+

$$\bar{w}_1 = \bar{w}_3 \quad \bar{w}_2 = \bar{w}_3$$

sin( )   *

$$\bar{x}_1 = \bar{w}_1 \cos(x_1) \quad \bar{x}_1 = \bar{w}_2 x_2 \quad \bar{x}_2 = \bar{w}_2 x_1$$

$x_1$   $x_2$

# Forward mode AD vs Reverse Mode AD

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



+

sin( )      *

$x_1$      $x_2$

$\frac{\partial f}{\partial x}$

$\dot{w}_3 = \dot{w}_1 + \dot{w}_2$

+

$\dot{w}_1 = \cos(x_1)\dot{x}_1$   $\dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$

sin( )      *

$\dot{x}_1$   $\dot{x}_1$   $\dot{x}_2$

$x_1$      $x_2$

$\bar{w}_3 = 1$

$\frac{\partial f}{\partial x}$

+

$\bar{w}_1 = \bar{w}_3$   $\bar{w}_2 = \bar{w}_3$

sin( )      *

$\bar{x}_1 = \bar{w}_1 \cos(x_1)$   $\bar{x}_1 = \bar{w}_2 x_2$   $\bar{x}_2 = \bar{w}_2 x_1$

$x_1$      $x_2$

40

# Forward mode vs Reverse Mode

- What are the differences?



$$\dot{w}_3 = \dot{w}_1 + \dot{w}_2$$

$$\dot{w}_1 = \cos(x_1)\dot{x}_1 \quad \dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$$

$$\bar{w}_3 = 1$$

$$\bar{w}_1 = \bar{w}_3 \quad \bar{w}_2 = \bar{w}_3$$

$$\bar{x}_1 = \bar{w}_1 \cos(x_1) \quad \bar{x}_1 = \bar{w}_2 x_2 \quad \bar{x}_2 = \bar{w}_2 x$$

sin( )    *    +    $x_1$    $x_2$    $\dot{x}_1$    $\dot{x}_1$    $\dot{x}_2$

# Forward mode vs Reverse Mode

- What are the differences?

- Which one is faster to compute?
  - Forward or backward?

*Depends* *Is $c > d$ or $d > c$?*

# Forward mode vs Reverse Mode

- x ✉ Graph ✉ L
- Intuition of Jacobian



$$\frac{\partial \vec{L}}{\partial \vec{x}} = \begin{bmatrix} i & & \\ & & \\ & \frac{\partial L_i}{\partial x_j} & \\ & & \end{bmatrix}$$

$\vec{z} \in \mathbb{R}^d$   $\vec{L} \in \mathbb{R}^c$

Graph / NN / Diff Program

for i in 1:c
RM-AD output in 1 pass

c × d

for j in 1:d
FM-AD outputs in 1-pass

# Forward mode vs Reverse Mode

- What are the differences?

- Which one is <u>faster</u> to compute?
  - Forward or backward?

- Which one is more memory efficient (less storage)?
  - Forward or backward?

# Plan for Today

- Automatic Differentiation
  - (Finish) Forward mode vs Reverse mode AD
  - Backprop
  - Patterns in backprop

# Neural Network Computation Graph



$$h^{(1)} = max(0, W_1 x)$$

$$\frac{\partial L}{\partial w_{(i)}}$$

# Backprop



$B \times 256$
$\in |R$

$\mathbf{x} \longrightarrow \boxed{f_1} \longrightarrow \mathbf{h}^{(1)} \longrightarrow \boxed{f_2} \longrightarrow \mathbf{h}^{(2)} \longrightarrow \cdots \longrightarrow \mathbf{h}^{(L-1)} \longrightarrow \boxed{f_L} \longrightarrow L$

$\mathbf{w}_1 \qquad \mathbf{w}_2 \qquad \mathbf{w}_L$

$d\mathbf{x} \longleftarrow \boxed{df_1} \longleftarrow d\mathbf{h}^{(1)} \longleftarrow \boxed{df_2} \longleftarrow d\mathbf{h}^{(2)} \longleftarrow \cdots \longleftarrow d\mathbf{h}^{(L-1)} \longleftarrow \boxed{df_L} \longleftarrow dL = 1$

$d\mathbf{w}_1 \qquad d\mathbf{w}_2 \qquad d\mathbf{w}_L$

Figure Credit: Andrea Vedaldi

# Computational Graph

**Any DAG of differentiable modules is allowed!**

Slide Credit: Marc'Aurelio Ranzato

# Key Computation: Forward-Prop

# Key Computation: Back-Prop

# Neural Network Training

- Step 1: Compute Loss on mini-batch [F-Pass]

# Neural Network Training

- Step 1: Compute Loss on mini-batch     [F-Pass]

# Neural Network Training

- Step 1: Compute Loss on mini-batch    [F-Pass]

# Neural Network Training

- Step 1: Compute Loss on mini-batch      [F-Pass]
- Step 2: Compute gradients wrt parameters  [B-Pass]

# Neural Network Training

- Step 1: Compute Loss on mini-batch     [F-Pass]
- Step 2: Compute gradients wrt parameters   [B-Pass]

# Neural Network Training

- Step 1: Compute Loss on mini-batch     [F-Pass]
- Step 2: Compute gradients wrt parameters   [B-Pass]



Slide Credit: Marc'Aurelio Ranzato, Yann LeCun

# Neural Network Training

- Step 1: Compute Loss on mini-batch     [F-Pass]
- Step 2: Compute gradients wrt parameters  [B-Pass]
- Step 3: Use gradient to update parameters



$$\theta \leftarrow \theta - \eta \frac{dL}{d\theta}$$

# Backpropagation: a simple example

# Backpropagation: a simple example

# Patterns in backward flow

# Patterns in backward flow

Q: What is an **add** gate?
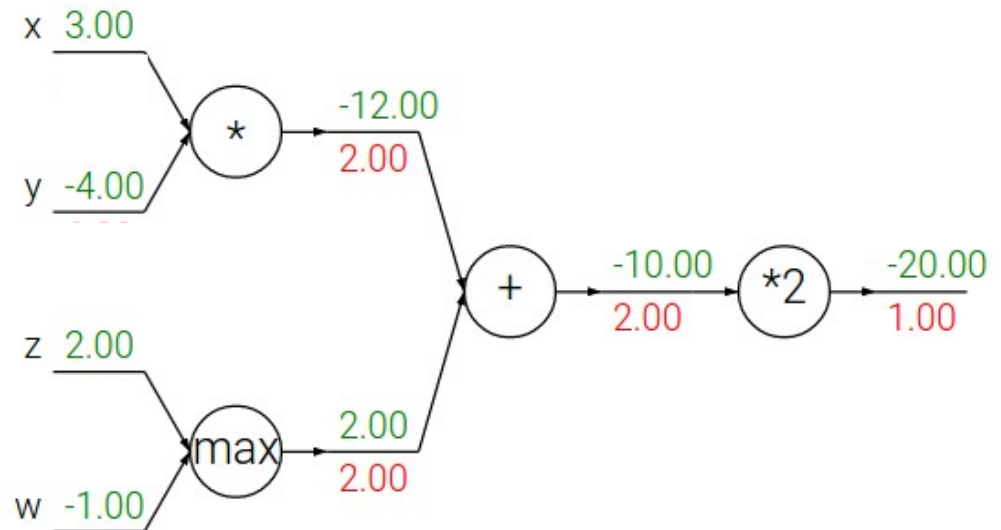
# Patterns in backward flow

**add** gate: gradient distributor

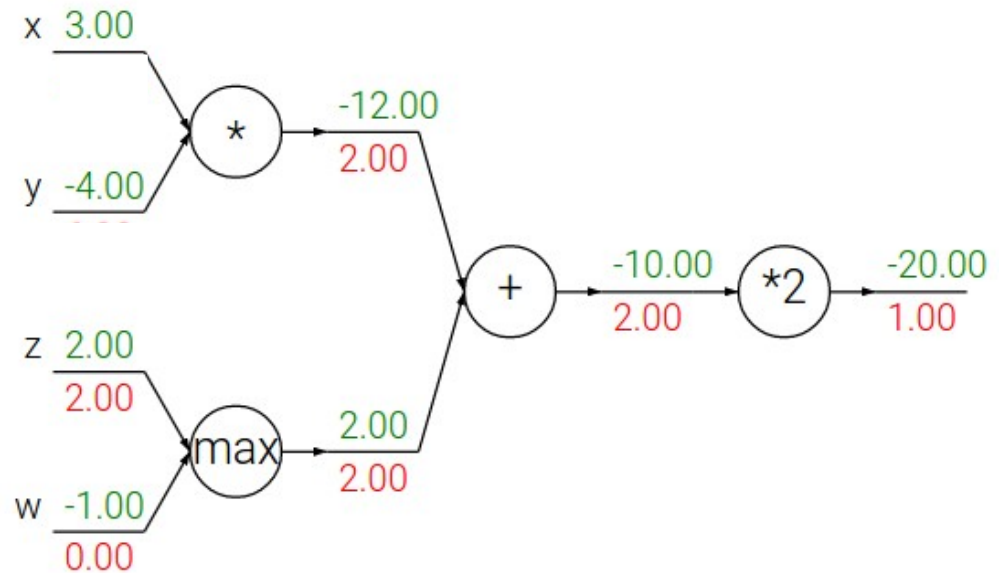# Patterns in backward flow

**add** gate: gradient distributor

Q: What is a **max** gate?

# Patterns in backward flow

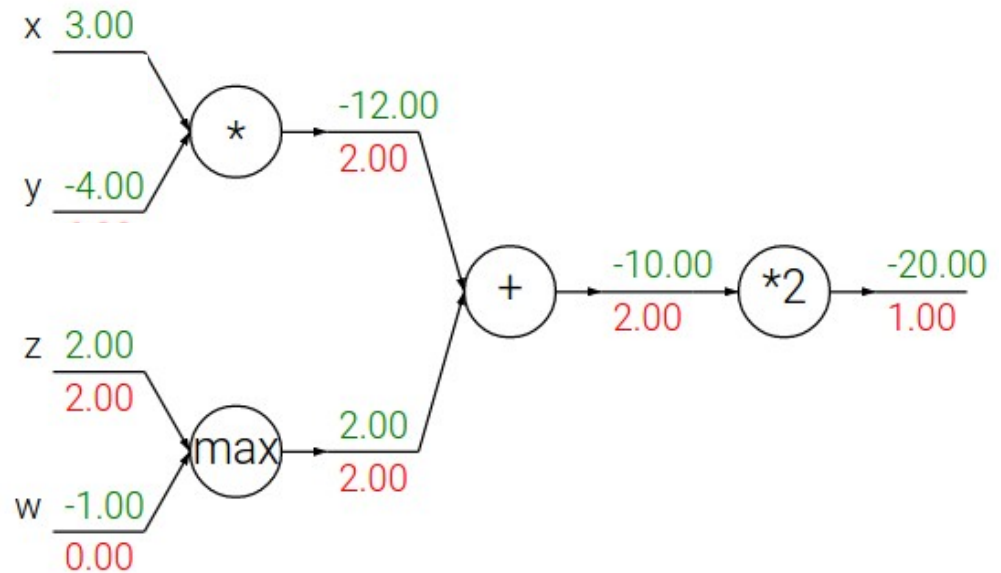**add** gate: gradient distributor

**max** gate: gradient router

# Patterns in backward flow

**add** gate: gradient distributor
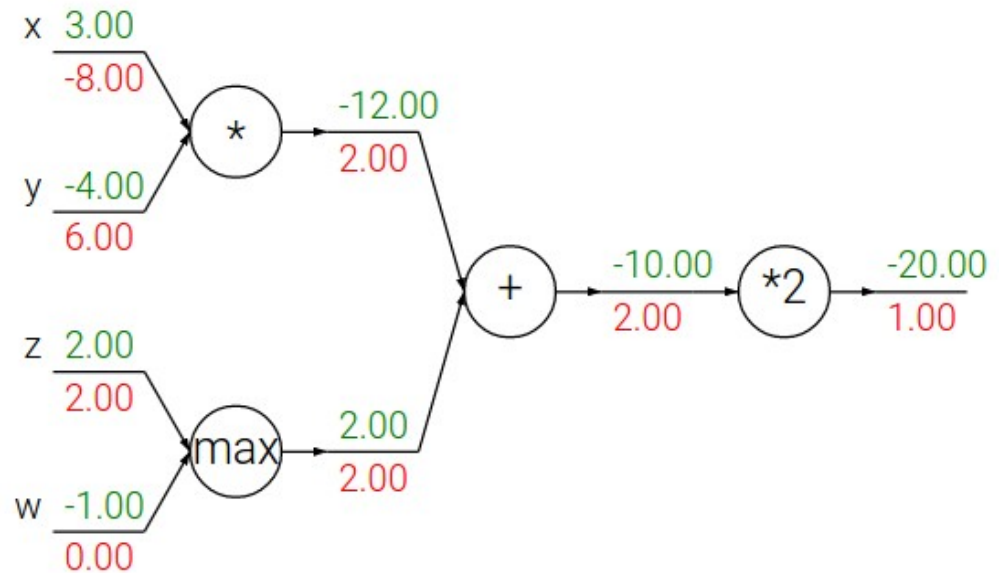
**max** gate: gradient router
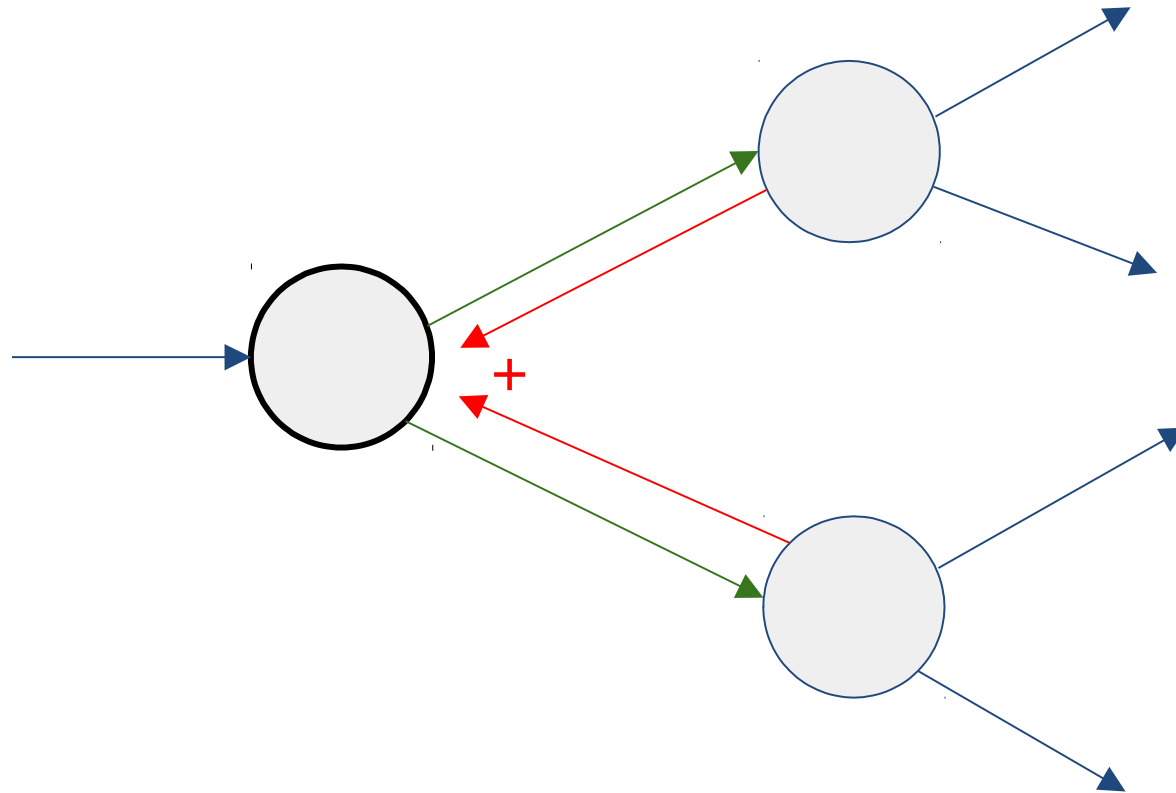
Q: What is a **mul** gate?

# Patterns in backward flow

**add** gate: gradient distributor

**max** gate: gradient router

**mul** gate: gradient switcher

# Gradients add at branches

# Duality in Fprop and Bprop

*FPROP*                    *BPROP*

SUM

COPY