

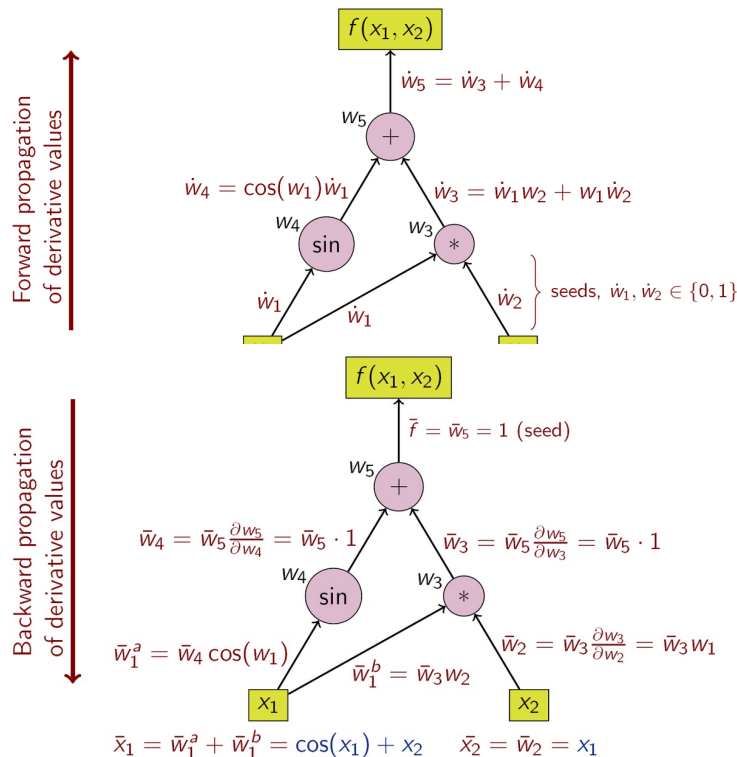
# Recitation 3.

# AutoDiff & Backprop

Sarthak Bisht, Ruimeng Chang

# Automatic Differentiation Framework

- A set of techniques to evaluate the derivative of a function by repeatedly applying chain rule
- Decompose arbitrarily complex function into simple operations
- Easy to program compared with symbolic differentiation



# Autodiff Example

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

Goal: Given  $x_1, x_2$  as inputs, calculate  $\partial y / \partial x_1$

# Autodiff Example

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2)$$

Add intermediate variables

$$v_1 = x_1$$

$$v_2 = x_2$$

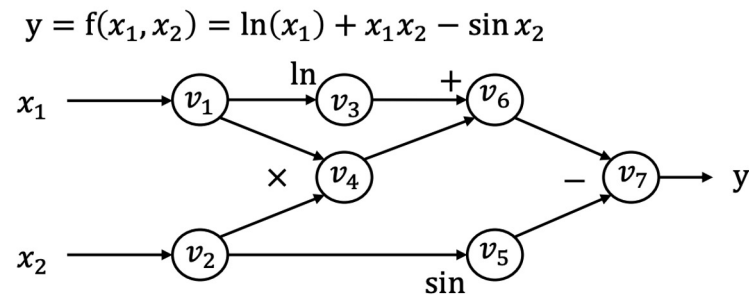
$$v_3 = \ln(v_1)$$

$$v_4 = v_1 v_2$$

$$v_5 = \sin(v_2)$$

$$v_6 = v_3 + v_4$$

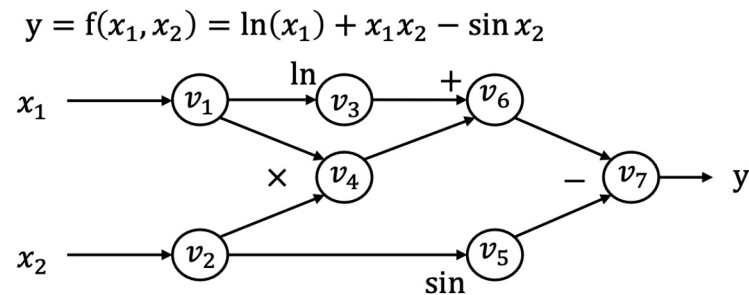
$$y = v_7 = v_6 - v_5$$



Goal: Given  $x_1, x_2$  as inputs, calculate  $\partial y / \partial x_1$

# Autodiff Example (Forward)

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2)$$



Add intermediate variables

Goal: Given  $x_1, x_2$  as inputs, calculate  $\partial y / \partial x_1$

$$v_1 = x_1$$

$$v_2 = x_2$$

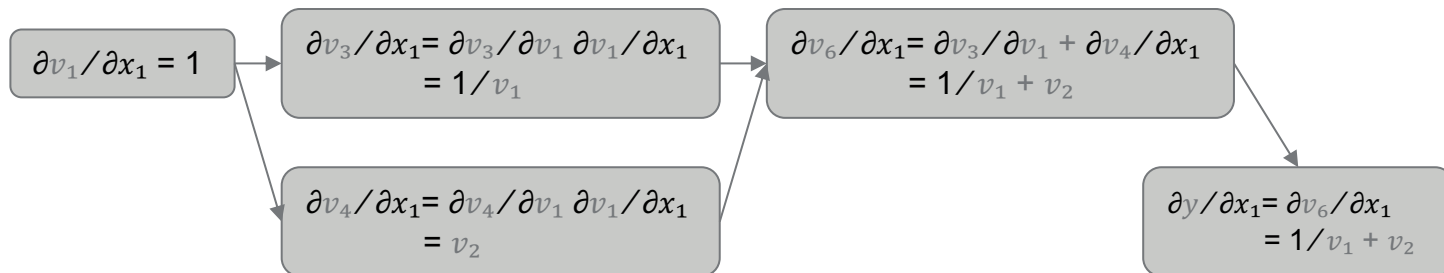
$$v_3 = \ln(v_1)$$

$$v_4 = v_1 v_2$$

$$v_5 = \sin(v_2)$$

$$v_6 = v_3 + v_4$$

$$y = v_7 = v_6 - v_5$$



# Autodiff Example (Reversed)

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2)$$

Add intermediate variables

$$v_1 = x_1$$

$$v_2 = x_2$$

$$v_3 = \ln(v_1)$$

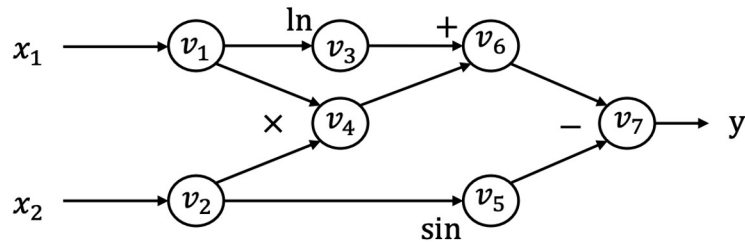
$$v_4 = v_1 v_2$$

$$v_5 = \sin(v_2)$$

$$v_6 = v_3 + v_4$$

$$y = v_7 = v_6 - v_5$$

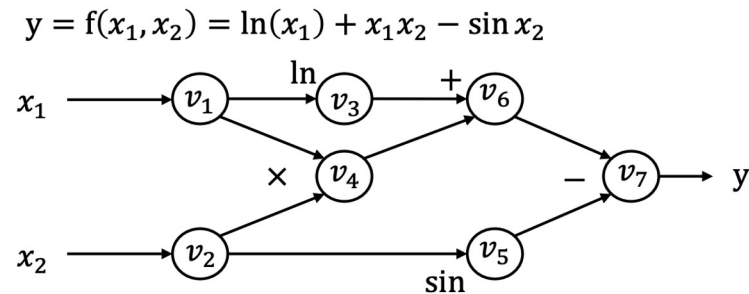
$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin x_2$$



Goal: Given  $x_1, x_2$  as inputs, calculate  $\partial y / \partial x_1$

# Autodiff Example (Reversed)

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2)$$



Add intermediate variables

$$v_1 = x_1$$

$$v_2 = x_2$$

$$v_3 = \ln(v_1)$$

$$v_4 = v_1 v_2$$

$$v_5 = \sin(v_2)$$

$$v_6 = v_3 + v_4$$

$$y = v_7 = v_6 - v_5$$

$$\frac{\partial y}{\partial v_1} = \frac{\partial y}{\partial v_3} \frac{\partial v_3}{\partial v_1} + \frac{\partial y}{\partial v_4} \frac{\partial v_4}{\partial v_1} = 1/v_1 + v_2$$

$$\frac{\partial y}{\partial v_3} = \frac{\partial y}{\partial v_6} \frac{\partial v_6}{\partial v_3} = 1$$

$$\frac{\partial y}{\partial v_4} = \frac{\partial y}{\partial v_6} \frac{\partial v_6}{\partial v_4} = 1$$

$$\frac{\partial y}{\partial v_6} = \frac{\partial y}{\partial v_7} \frac{\partial v_7}{\partial v_6} = 1$$

$$\frac{\partial y}{\partial v_7} = 1$$

# Reference

<https://dlsyscourse.org/slides/4-automatic-differentiation.pdf>

[https://en.wikipedia.org/wiki/Automatic\\_differentiation](https://en.wikipedia.org/wiki/Automatic_differentiation)

<https://deeplearning.cs.cmu.edu/F22/document/recitation/Recitation3/Computing%20Derivatives%20and%20Autograd%20-%20Fall22.pdf>



# HW1 Autograd

- Optional bonus HW
- Alternative implementation of MyTorch (HW1P1) based on the reversed version of Autodiff framework

What you'll accomplish

- Implement an engine Autograd that stores every operation in sequence (computation graph!)
- Calculate and update gradients for each operation in computation graph
- Build activations, losses, layers using Autograd
- Run MLP using Autograd

# Backpropagation

- A special case of Reverse Autodiff where the function output is scalar (i.e. the loss value)
- Part of neural network training
  1. Forward Propagation with current parameters
  2. Calculate the loss
  - 3. Backward Propagation to calculate the gradients of the parameters**
  4. Step to update the parameters with gradients

# Backpropagation of Loss

- “The function we want to minimize or maximize is called the objective function or criterion. When we are minimizing it, we may also call it the cost function, loss function, or error function.”
- Loss is the starting point of the backpropagation
- Backpropagation is the method to minimize the loss function by adjusting network's weights and biases. The level of adjustment is determined by the **gradients of the loss function w.r.t. those parameters**.

(Remember: gradient is the transpose of derivative)

# Derivative

The derivative of a function  $f$  measures the sensitivity of change of the function value w.r.t. a change in its input value  $x$ .

$$\frac{dy}{dx} = \lim_{\Delta x \rightarrow 0} \frac{\Delta y}{\Delta x}$$

The derivative of  $y$  with respect to  $x$ :

$$\Delta y = \nabla_x y \Delta x$$

Useful formula to check derivative shape!

Note: For scalar  $y$ , the derivative wrt  $x$  will have the shape as transpose to  $x$

# Computing Derivatives – Scalar Chain Rule

$$L = f(z)$$

$$z = g(x)$$

All terms are scalars

$\frac{\partial L}{\partial z}$  is given

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x} = \frac{\partial L}{\partial z} g'(x)$$

# Computing Derivatives – Scalar Multiplication

---

$$L = f(z)$$

$$z = Wx$$

All terms are scalars

$\frac{\partial L}{\partial z}$  is given

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x} = \frac{\partial L}{\partial z} W$$

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial W} = x \frac{\partial L}{\partial z}$$

# Computing Derivatives – Multivariable Chain Rule

$$L = f(z)$$

$$z = g(x)$$

$x$  is  $D \times 1$  vector,  $z$  is  $K \times 1$  vector

$\nabla_z L$  is given  $(1 \times K)$  matrix

$$\nabla_x L = \nabla_z L \nabla_x z$$

$$M \times D \quad M \times K \quad K \times D$$

# Computing Derivatives - Multivariate Matrix Multiplication

$$L = f(z)$$

$$z = Wx$$

$x$  is a  $D \times 1$  vector

$z$  is a  $K \times 1$  vector

$W$  is a  $K \times D$  matrix

$\nabla_z L$  is given  $(1 \times K)$  vector

$$\nabla_x L = \nabla_z L \nabla_x Z = (\nabla_z L)W \quad 1 \times D$$

$$\nabla_W L = \nabla_z L \nabla_W Z = x(\nabla_z L) \quad D \times K$$



# Computing Derivatives - Multivariable Generalized Chain Rule

$$L = f(z)$$

$$z = z_1 + z_2 + \cdots + z_n = g_1(x) + g_2(x) + \cdots + g_n(x)$$

$x$  is a  $D \times 1$  vector

$z$  is a  $K \times 1$  vector

$\nabla_z L$  is given  $(1 \times K)$  matrix

$$\nabla_x L = \nabla_z L \nabla_x Z = \nabla_z L (\nabla_x Z_1 + \nabla_x Z_2 + \cdots + \nabla_x Z_n)$$

# Questions?

# Backward Pass

## Why?

- Output layer ( $N$ ) :

- For  $i = 1 \dots D_N$

- $\frac{\partial Div}{\partial y_i^{(N)}} = \frac{\partial Div(Y, d)}{\partial y_i}$

- $\frac{\partial Div}{\partial z_i^{(N)}} = \frac{\partial Div}{\partial y_i^{(N)}} f'_N(z_i^{(N)})$

- $\frac{\partial Div}{\partial w_{ij}^{(N)}} = y_i^{(N-1)} \frac{\partial Div}{\partial z_j^{(N)}} \text{ for } j = 0 \dots D_{N-1}$

Called "**Backpropagation**" because the derivative of the loss is propagated "backwards" through the network

- For layer  $k = N - 1$  *downto* 1

Very analogous to the forward pass:

- For  $i = 1 \dots D_k$

- $\frac{\partial Div}{\partial y_i^{(k)}} = \sum_j w_{ij}^{(k+1)} \frac{\partial Di}{\partial z_j^{(k+1)}}$

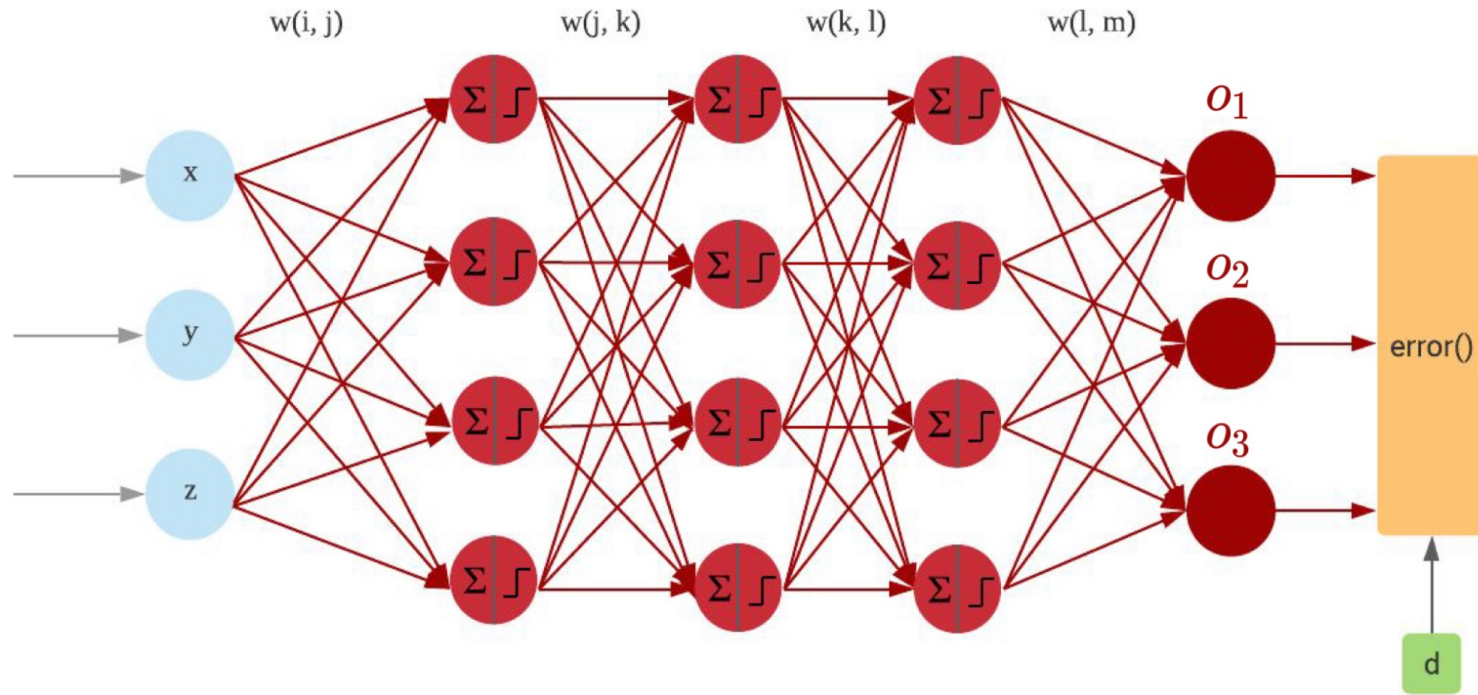
Backward weighted combination of next layer

- $\frac{\partial Div}{\partial z_i^{(k)}} = \frac{\partial Div}{\partial y_i^{(k)}} f'_k(z_i^{(k)})$

Backward equivalent of activation

- $\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}} \text{ for } j = 0 \dots D_{k-1}$

# Why calculate partial derivatives in backward direction?



# AutoDiff vs Numerical Differentiation?

Directly compute the partial gradient by definition

$$\frac{\partial f(\theta)}{\partial \theta_i} = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon e_i) - f(\theta)}{\epsilon}$$

A more numerically accurate way to approximate the gradient

$$\frac{\partial f(\theta)}{\partial \theta_i} = \frac{f(\theta + \epsilon e_i) - f(\theta - \epsilon e_i)}{2\epsilon} + o(\epsilon^2)$$

Suffer from numerical error, less efficient to compute

# AutoDiff vs Symbolic Differentiation?

Write down the formulas, derive the gradient by sum, product and chain rules

$$\frac{\partial(f(\theta)+g(\theta))}{\partial\theta} = \frac{\partial f(\theta)}{\partial\theta} + \frac{\partial g(\theta)}{\partial\theta}$$

$$\frac{\partial(f(\theta)g(\theta))}{\partial\theta} = g(\theta) \frac{\partial f(\theta)}{\partial\theta} + f(\theta) \frac{\partial g(\theta)}{\partial\theta}$$

$$\frac{\partial f(g(\theta))}{\partial\theta} = \frac{\partial f(g(\theta))}{\partial g(\theta)} \frac{\partial g(\theta)}{\partial\theta}$$

Naively do so can result in wasted computations

Example:  $f(\theta) = \prod_{i=0}^n \theta_i$        $\frac{f(\theta)}{\partial\theta_k} = \prod_{j \neq k}^n \theta_j$

Cost  $n(n-1)$  multiplies to compute all partial gradients

# What's the difference?

## Forward Mode AutoDiff

**Define**  $\dot{v}_i = \frac{\partial v_i}{\partial x_1}$

We can then compute the  $\dot{v}_i$  iteratively in the forward topological order of the computational graph

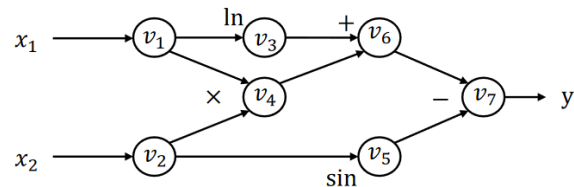
Forward AD trace

$$\begin{aligned}\dot{v}_1 &= 1 \\ \dot{v}_2 &= 0 \\ \dot{v}_3 &= \dot{v}_1/v_1 = 0.5 \\ \dot{v}_4 &= \dot{v}_1 v_2 + \dot{v}_2 v_1 = 1 \times 5 + 0 \times 2 = 5 \\ \dot{v}_5 &= \dot{v}_2 \cos v_2 = 0 \times \cos 5 = 0 \\ \dot{v}_6 &= \dot{v}_3 + \dot{v}_4 = 0.5 + 5 = 5.5 \\ \dot{v}_7 &= \dot{v}_6 - \dot{v}_5 = 5.5 - 0 = 5.5\end{aligned}$$

**Now we have**  $\frac{\partial y}{\partial x_1} = \dot{v}_7 = 5.5$

9

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin x_2$$



## Reverse Mode AutoDiff

**Define adjoint**  $\bar{v}_i = \frac{\partial y}{\partial v_i}$

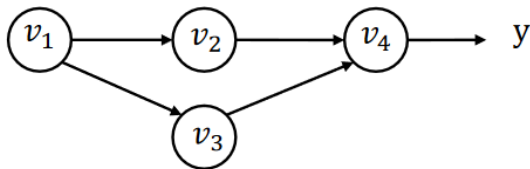
We can then compute the  $\bar{v}_i$  iteratively in the **reverse** topological order of the computational graph

Reverse AD evaluation trace

$$\begin{aligned}\bar{v}_7 &= \frac{\partial y}{\partial v_7} = 1 \\ \bar{v}_6 &= \bar{v}_7 \frac{\partial v_7}{\partial v_6} = \bar{v}_7 \times 1 = 1 \\ \bar{v}_5 &= \bar{v}_7 \frac{\partial v_7}{\partial v_5} = \bar{v}_7 \times (-1) = -1 \\ \bar{v}_4 &= \bar{v}_6 \frac{\partial v_6}{\partial v_4} = \bar{v}_6 \times 1 = 1 \\ \bar{v}_3 &= \bar{v}_6 \frac{\partial v_6}{\partial v_3} = \bar{v}_6 \times 1 = 1 \\ \bar{v}_2 &= \bar{v}_5 \frac{\partial v_5}{\partial v_2} + \bar{v}_4 \frac{\partial v_4}{\partial v_2} = \bar{v}_5 \times \cos v_2 + \bar{v}_4 \times v_1 = -0.284 + 2 = 1.716 \\ \bar{v}_1 &= \bar{v}_4 \frac{\partial v_4}{\partial v_1} + \bar{v}_3 \frac{\partial v_3}{\partial v_1} = \bar{v}_4 \times v_2 + \bar{v}_3 \frac{1}{v_1} = 5 + \frac{1}{2} = 5.5\end{aligned}$$

# Multiple Pathway Case

$v_1$  is being used in multiple pathways ( $v_2$  and  $v_3$ )



$y$  can be written in the form of  $y = f(v_2, v_3)$

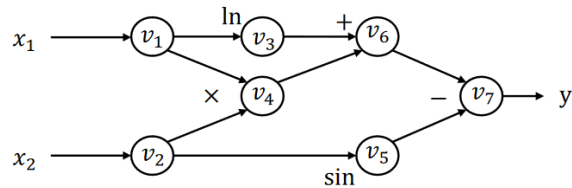
$$\overline{v_1} = \frac{\partial y}{\partial v_1} = \frac{\partial f(v_2, v_3)}{\partial v_2} \frac{\partial v_2}{\partial v_1} + \frac{\partial f(v_2, v_3)}{\partial v_3} \frac{\partial v_3}{\partial v_1} = \overline{v_2} \frac{\partial v_2}{\partial v_1} + \overline{v_3} \frac{\partial v_3}{\partial v_1}$$

Define partial adjoint  $\overline{v_{i \rightarrow j}} = \overline{v_j} \frac{\partial v_j}{\partial v_i}$  for each input output node pair  $i$  and  $j$

$$\overline{v_i} = \sum_{j \in \text{next}(i)} \overline{v_{i \rightarrow j}}$$

We can compute partial adjoints separately then sum them together

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin x_2$$





# Reverse AutoDiff Algorithm

```
def gradient(out):
    node_to_grad = {out: [1]}

    for i in reverse_topo_order(out):
         $\overline{v_i} = \sum_j \overline{v_{i \rightarrow j}} = \text{sum}(\text{node\_to\_grad}[i])$ 

        for  $k \in \text{inputs}(i)$ :
            compute  $\overline{v_{k \rightarrow i}} = \overline{v_i} \frac{\partial v_i}{\partial v_k}$ 
            append  $\overline{v_{k \rightarrow i}}$  to  $\text{node\_to\_grad}[k]$ 

    return adjoint of input  $\overline{v_{\text{input}}}$ 
```

Dictionary that records a list of partial adjoints of each node

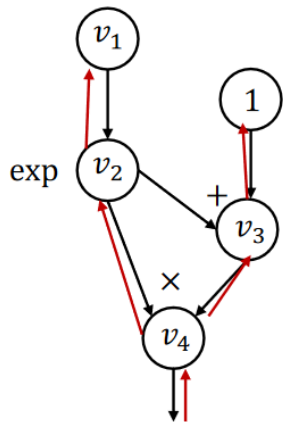
Sum up partial adjoints

“Propagates” partial adjoint to its input



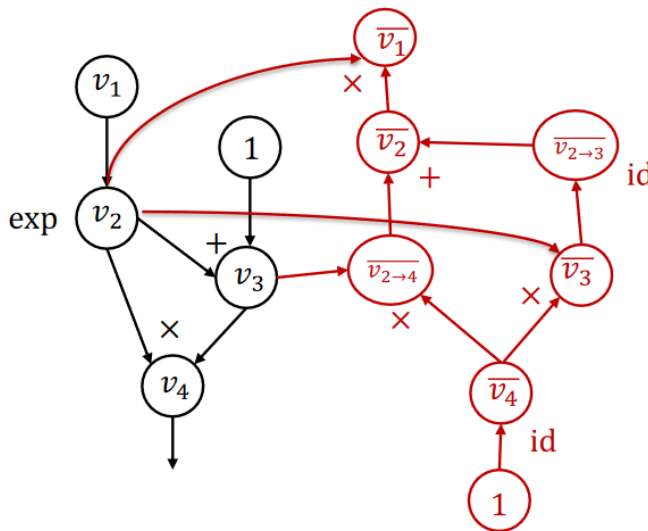
# Backprop vs Reverse mode AutoDiff

**Backprop**



- Run backward operations the same forward graph
- Used in first generation deep learning frameworks (caffe, cuda-convnet)

**Reverse mode AD by  
extending computational graph**



- Construct separate graph nodes for adjoints
- Used by modern deep learning frameworks