

CS 4803 / 7643: Deep Learning

Topics:

- (Finish) Computational Graphs
 - Notation + example
- (Finish) Computing Gradients
 - Forward mode vs Reverse mode AD
 - Patterns in backprop
 - Backprop in FC+ReLU NNs

Dhruv Batra
Georgia Tech

Administrativa

- HW1 Reminder
 - Due: 10/02, 11:55pm
 - https://www.cc.gatech.edu/classes/AY2019/cs7643_fall/assets/hw1.pdf
 - https://www.cc.gatech.edu/classes/AY2019/cs7643_fall/hw1-q6/

Recap from last time

Strategy: Follow the slope

$$\min_w \underline{L(w; D)}$$



Gradient Descent

```
# Vanilla Gradient Descent
```

```
while True:
```

```
    [weights_grad = evaluate_gradient(loss_fun, data, weights)] ← backprop  
    weights += - step_size * weights_grad # perform parameter update
```

$$w^{(0)} = \text{init}$$

for $t=1 \dots \text{times}$

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

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

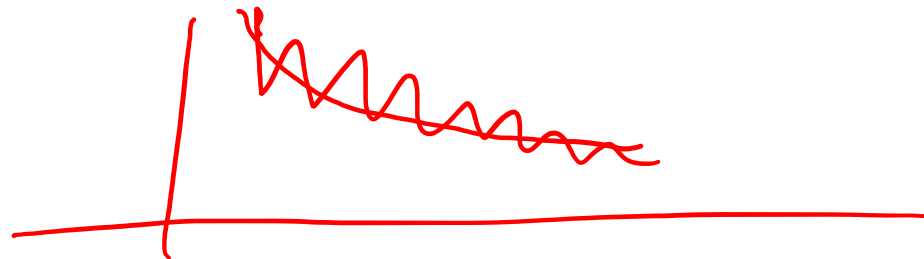
```
# Vanilla Minibatch Gradient Descent
```

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

$$l_1 = x$$

$$l_{n+1} = 4l_n(1 - l_n)$$

$$f(x) = l_4 = 64x(1-x)(1-2x)^2(1-8x+8x^2)^2$$

Manual
Differentiation

$$f'(x) = 128x(1-x)(-8+16x)(1-2x)^2(1-8x+8x^2) + 64(1-x)(1-2x)^2(1-8x+8x^2)^2 - 64x(1-2x)^2(1-8x+8x^2)^2 - 256x(1-x)(1-2x)(1-8x+8x^2)^2$$

Coding

```
f(x):
v = x
for i = 1 to 3
    v = 4v(1 - v)
v
or, in closed-form,
f(x):
    64x (1-x) (1-2x)^2 (1-8x+8x^2)^2
```

Coding

```
f'(x):
128x(1-x)(-8+16x)(1-2
x)^2(1-8x+8x^2)+64(1
-x)(1-2x)^2(1-8x+8
x^2)^2-64x(1-2x)^2(1-8
x+8x^2)^2-256x(1-x)(1-
2x)(1-8x+8x^2)^2
f'(x_0) = f'(x_0)
Exact
```

Symbolic
Differentiation
of the Closed-form

Automatic
Differentiation

```
f'(x):
(v,v') = (x,1)
for i = 1 to 3
    (v,v') = (4v(1-v), 4v'-8vv')
```

$$f'(x_0) = f'(x_0)$$

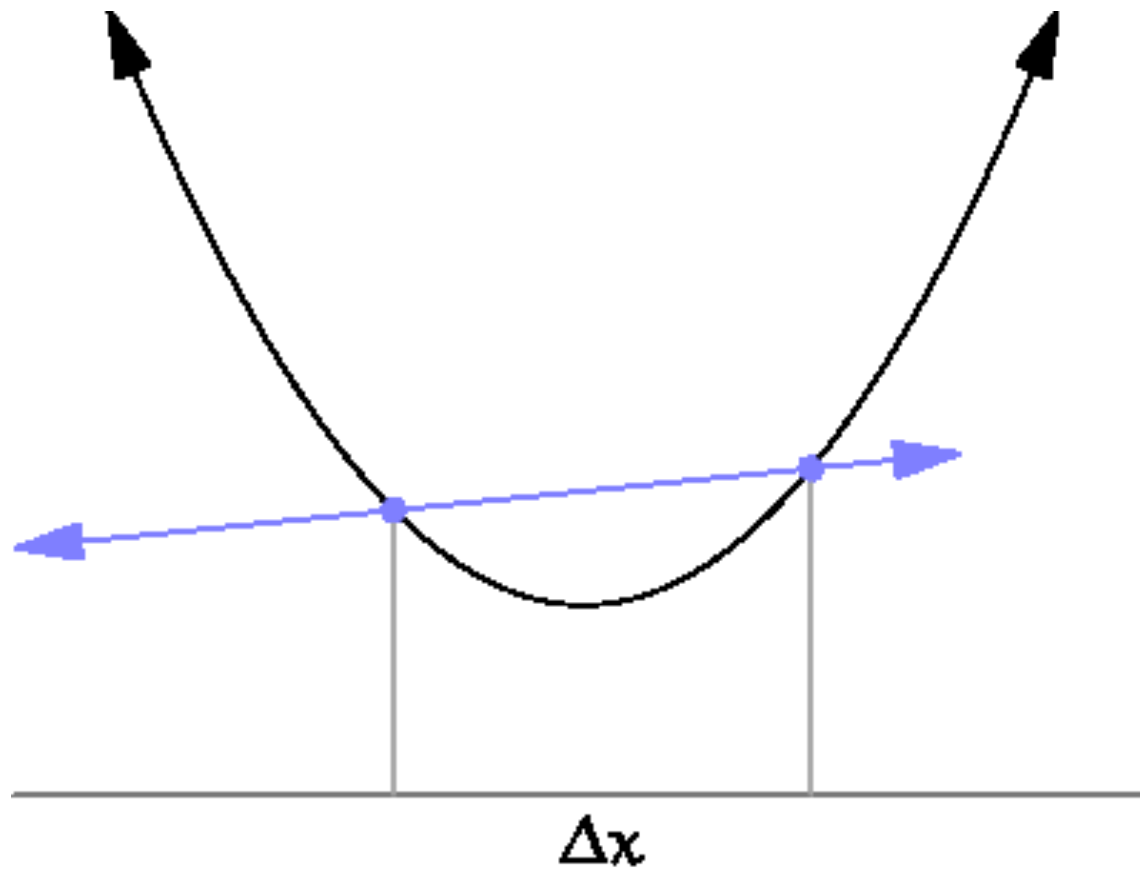
Exact

Numerical
Differentiation

```
f'(x):
h = 0.000001
(f(x+h) - f(x)) / h
f'(x_0) ≈ f'(x_0)
Approximate
```

How do we compute gradients?

- Analytic or “Manual” Differentiation
- Symbolic Differentiation
- 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

gradient dW:

[?,
?,
?,
?,
?,
?,
?,
?,
?,
?,...]

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:

[?,
?,
?,
?,
?,
?,
?,
?,
?,...]

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,
?,
?,

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

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

?,
?,...]

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 (second dim):

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

loss 1.25353

gradient dW:

[-2.5,
?,
?,
?,
?,
?,
?,
?,
?,...]

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 (second dim):

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

loss 1.25353

gradient dW:

[-2.5,
0.6,
?,
?,

$$(1.25353 - 1.25347)/0.0001 = 0.6$$

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

?,...]

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 (third dim):

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

loss 1.25347

gradient dW:

[-2.5,
0.6,
?,
?,
?,
?,
?,
?,
?,
?,...]

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 (third dim):

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

loss 1.25347

gradient dW:

[-2.5,
0.6,
0,
?,
?

$$(1.25347 - 1.25347)/0.0001 = 0$$

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

?,...]

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 :(

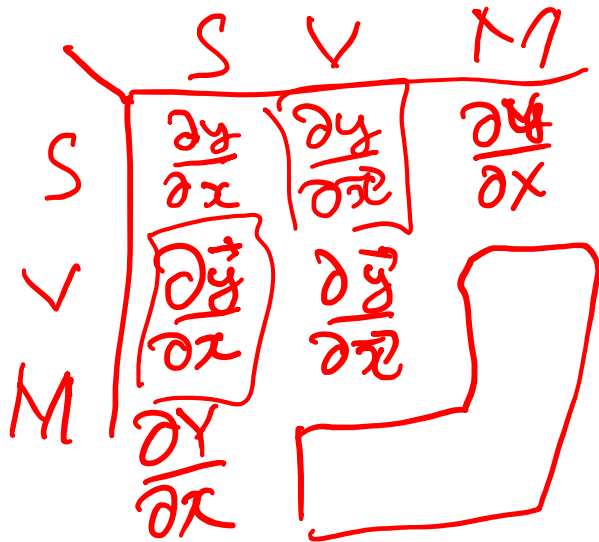
In practice: Derive analytic gradient, check your implementation with numerical gradient.

This is called a gradient check.

How do we compute gradients?

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

Matrix/Vector Derivatives Notation



$x, y \in \mathbb{R}^1$
 $\vec{x} \in \mathbb{R}^d$ $\vec{y} \in \mathbb{R}^c$
 $X, Y \in \mathbb{R}^{m \times n}$

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

$$\begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \vdots \\ \frac{\partial y_c}{\partial x} \end{bmatrix}$$

\downarrow num = dim 1

$$\left(\frac{\partial \vec{y}}{\partial \vec{x}} \right) =$$

$$\left[\frac{\partial y}{\partial x_1} \quad \frac{\partial y}{\partial x_2} \quad \dots \quad \frac{\partial y}{\partial x_d} \right]$$

\rightarrow den = dim 2

Matrix/Vector Derivatives Notation

$$\frac{\partial \vec{y}}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_c}{\partial x_1} & \dots & \frac{\partial y_c}{\partial x_d} \end{bmatrix} \quad c \times d$$

$\frac{\partial y_i}{\partial x_j}$

$$\frac{\partial (x^T w)}{\partial \vec{w}} = \left[\frac{\partial (x^T w)}{\partial w_1} \quad \dots \quad \frac{\partial (x^T w)}{\partial w_d} \right]$$

$\sum_{i=1}^c x_i w_i$

$$= x^T \quad [x_1 \quad \dots \quad x_d] = x^T$$

Vector Derivative Example

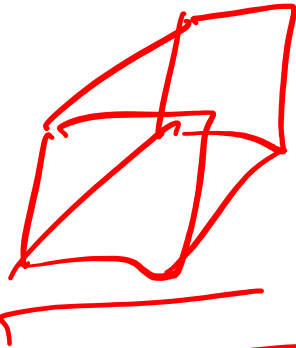
$$\frac{\partial (\underline{w}^T A \underline{x})}{\partial (\underline{w})} = \underline{w}^T A$$

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

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

$$\frac{\partial \underline{y}}{\partial \underline{x}} = A \left[\begin{array}{c} \frac{\partial y_i}{\partial x_j} a_{ij} \\ \vdots \\ \end{array} \right]$$

Extension to Tensors



$$x \in \mathbb{R}^{d_1 \times d_2 \times \dots \times d_n}$$

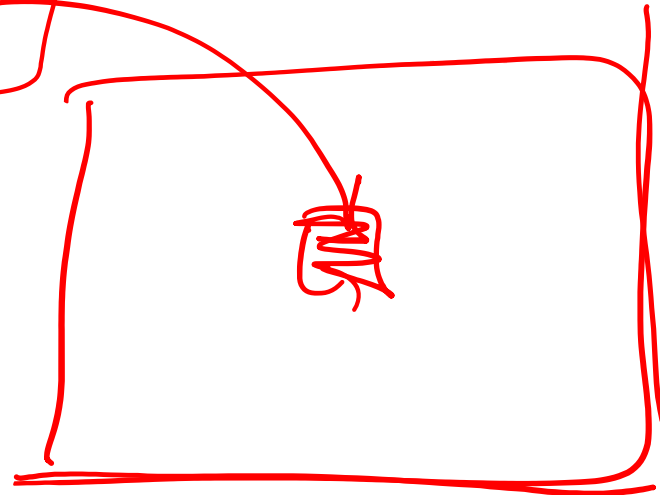
$$y \in \mathbb{R}^{c_1 \times c_2 \times \dots \times c_m}$$

$$y\text{-vec} = \underline{Y(\cdot)}$$

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

$$\frac{\partial (Y[i_1, i_2, \dots, i_m])}{\partial X[j_1, \dots, j_n]}$$

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



Chain Rule: Composite Functions

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

$$\begin{aligned} f(x) &= g_l(g_{l-1} \dots g_1(x)) \\ &= \underbrace{(g_l \circ g_{l-1} \dots \circ g_1)}_{\leftarrow} (x) \end{aligned}$$

Chain Rule: Scalar Case

$$\underbrace{x}_{\text{in}} \rightarrow z \rightarrow \underbrace{y}_{\text{out}} = \underbrace{f}_{\text{out}}(\underbrace{g}_{\text{in}}(x))$$

$$z = g(x)$$

$$y = f(z)$$

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

Chain Rule: Vector Case

$$\underline{\vec{x}} \in \mathbb{R}^d \rightarrow \underline{\vec{z}} \in \mathbb{R}^m \rightarrow \underline{y} \in \mathbb{R}^c$$

$$\vec{z} = g(\vec{x})$$

$$g: \mathbb{R}^d \rightarrow \mathbb{R}^m$$

$$\vec{y} = f(\vec{z})$$

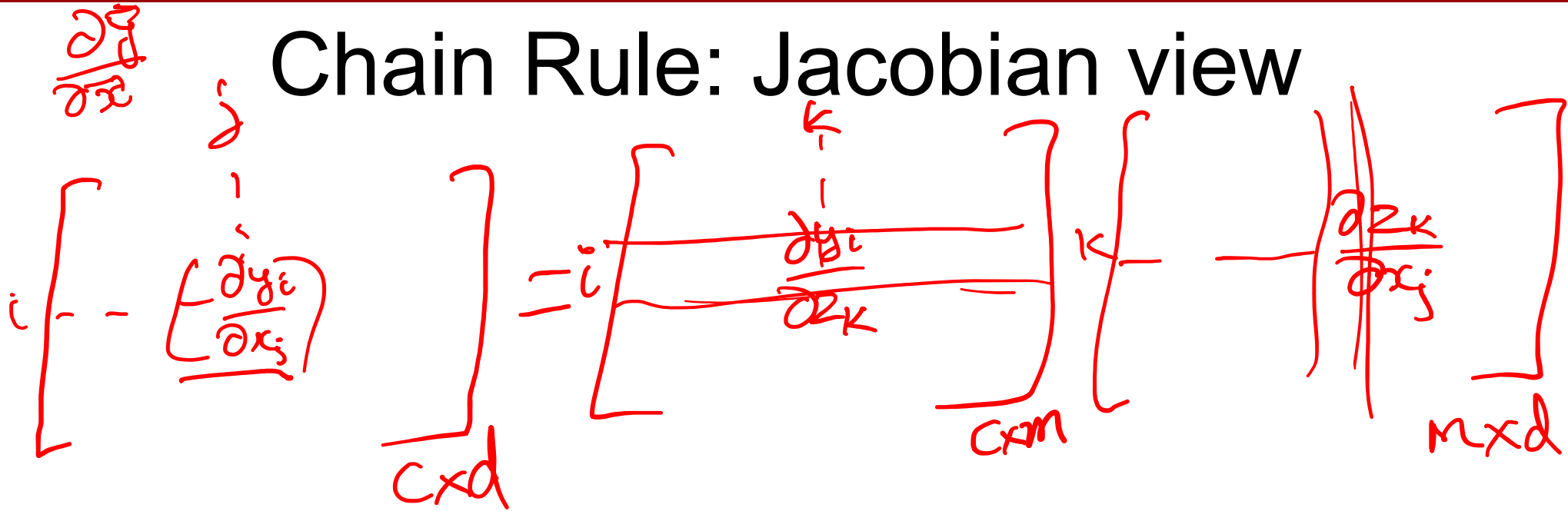
$$f: \mathbb{R}^m \rightarrow \mathbb{R}^c$$

$$\begin{array}{|c|} \hline \frac{\partial \vec{y}}{\partial \vec{x}} \\ \hline \frac{\partial \vec{z}}{\partial \vec{x}} \\ \hline \end{array} = \underline{J_{f \circ g}}$$

$$\begin{array}{|c|} \hline \frac{\partial \vec{y}}{\partial \vec{z}} \\ \hline \frac{\partial \vec{z}}{\partial \vec{x}} \\ \hline \end{array} \cdot \underline{J_f}$$

$$\begin{array}{|c|} \hline \frac{\partial \vec{z}}{\partial \vec{x}} \\ \hline \underline{J_g} \\ \hline \end{array}$$

Chain Rule: Jacobian view



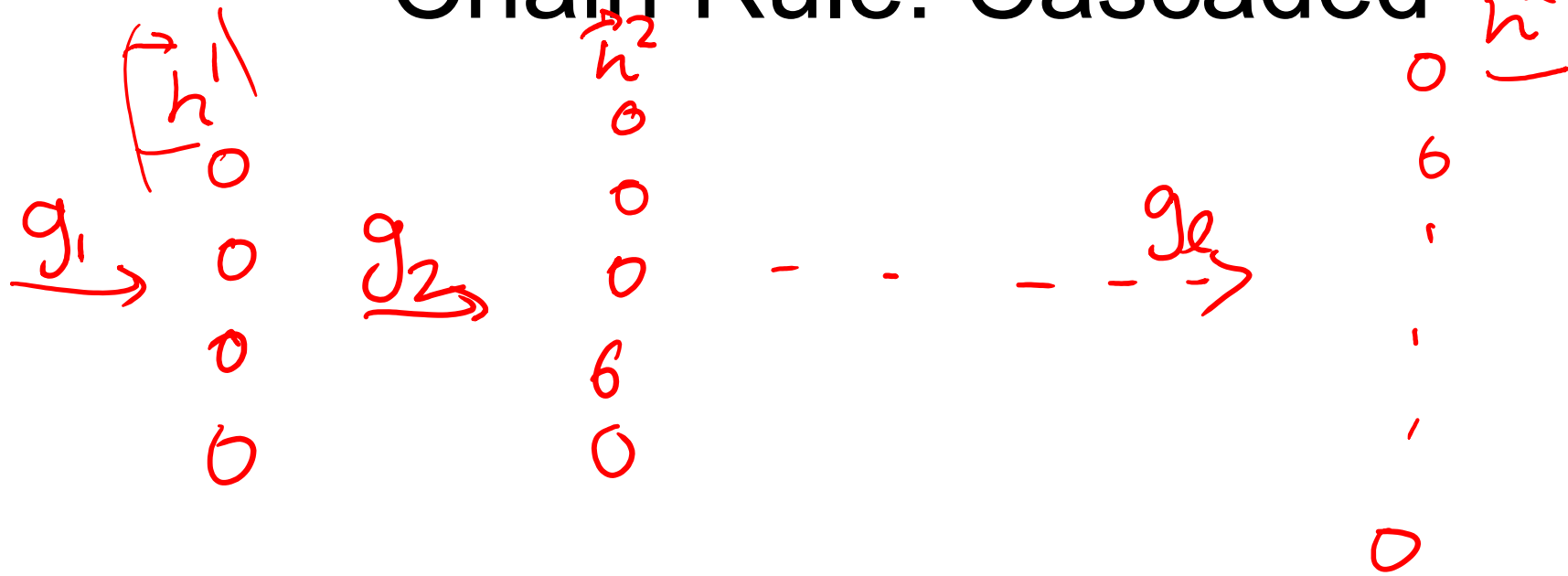
$$\frac{\partial y_i}{\partial x_j} = \sum_k \frac{\partial y_i}{\partial z_k} \frac{\partial z_k}{\partial x_j}$$

Chain Rule: Graphical view



$$\frac{\partial y_i}{\partial x_j} = \sum_{\text{paths}} \frac{\partial y_i}{\partial z_k} \frac{\partial z_k}{\partial x_j}$$

Chain Rule: Cascaded \vec{h}^e



$$\frac{\partial \vec{h}^e}{\partial \vec{h}^1} = \underline{J_{g_e}} \underline{J_{g_{e-1}}} \dots \underline{J_{g_2}}$$

$\xrightarrow{\alpha(d^2)}$
 $\xleftarrow{\alpha(d^3)}$

$$= \underline{\alpha(d^2)} \underline{J_{1 \times d}} \left[\dots \right]_{d \times d} \cdot \left[\dots \right]_{d \times d} \left[\dots \right]_{d \times d}$$

$\underbrace{\hspace{15em}}_{\alpha(d^3)}$

Chain Rule: How should we multiply?

Plan for Today

- (Finish) Computational Graphs }
 - Notation + example
- (Finish) Computing Gradients }
 - Forward mode vs Reverse mode AD
 - Patterns in backprop
 - Backprop in FC+ReLU NNs

Linear Classifier: Logistic Regression

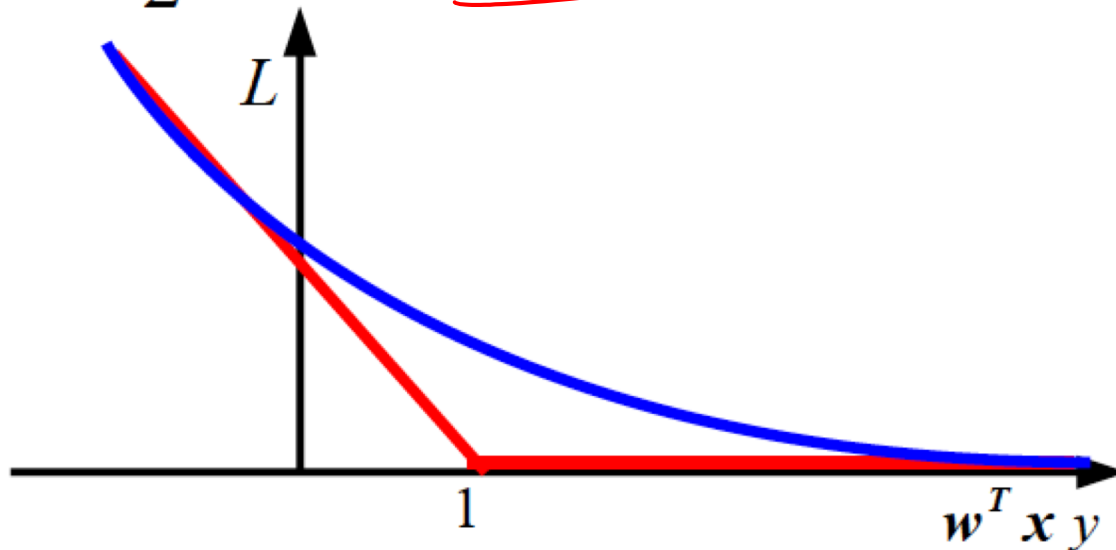
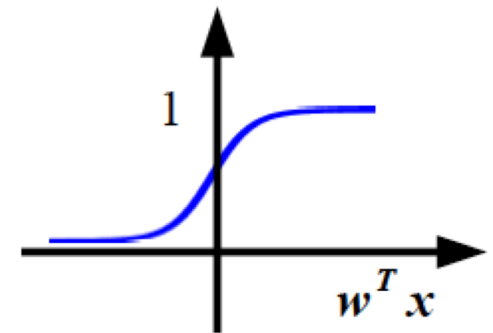
Input: $\mathbf{x} \in \mathbb{R}^D$

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

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

$$\frac{\partial L_i}{\partial \vec{w}} = \frac{-1}{1 + e^{-w^T x}} \cdot \frac{-1}{(1 + e^{-w^T x})^2} \cdot (-e^{-w^T x}) \cdot x^T$$

Convolutional network (AlexNet)

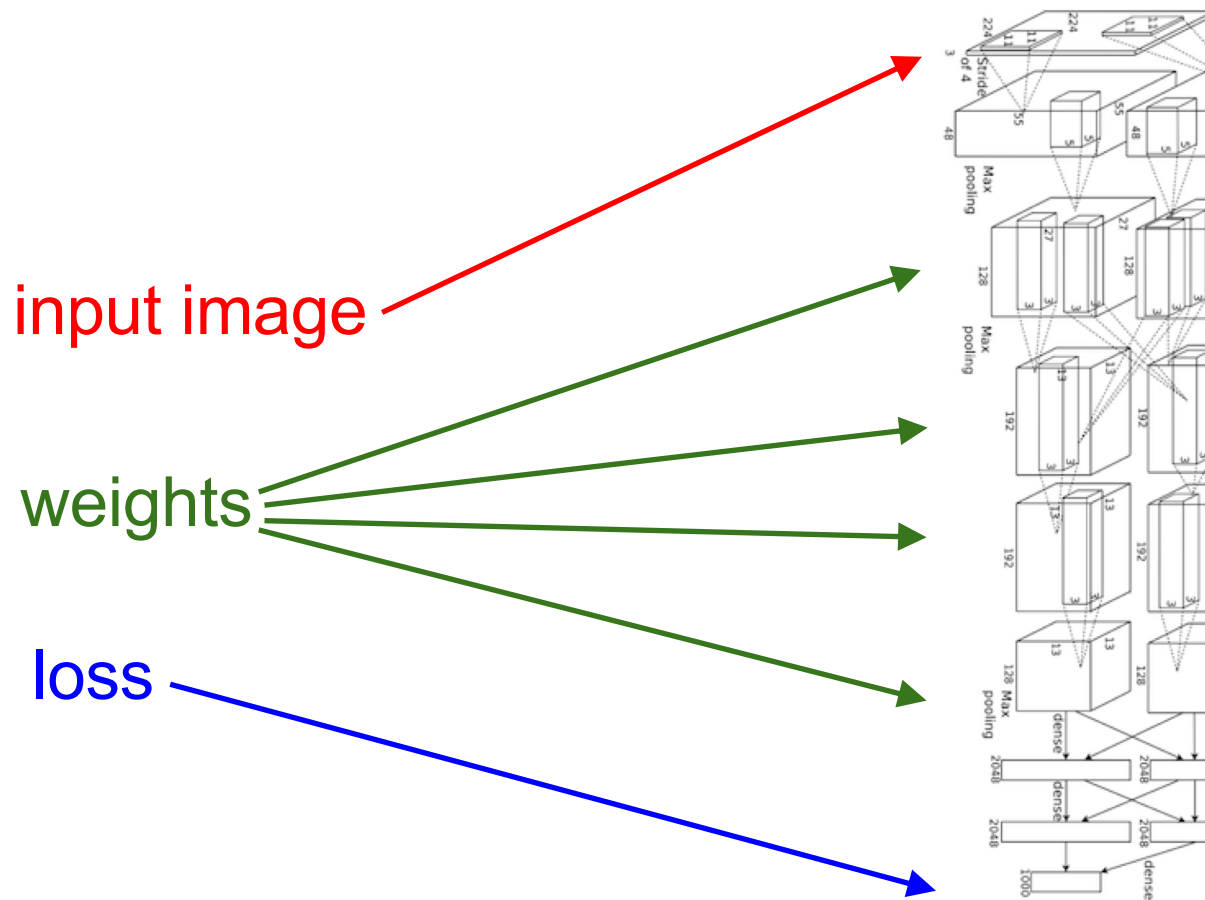


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Neural Turing Machine

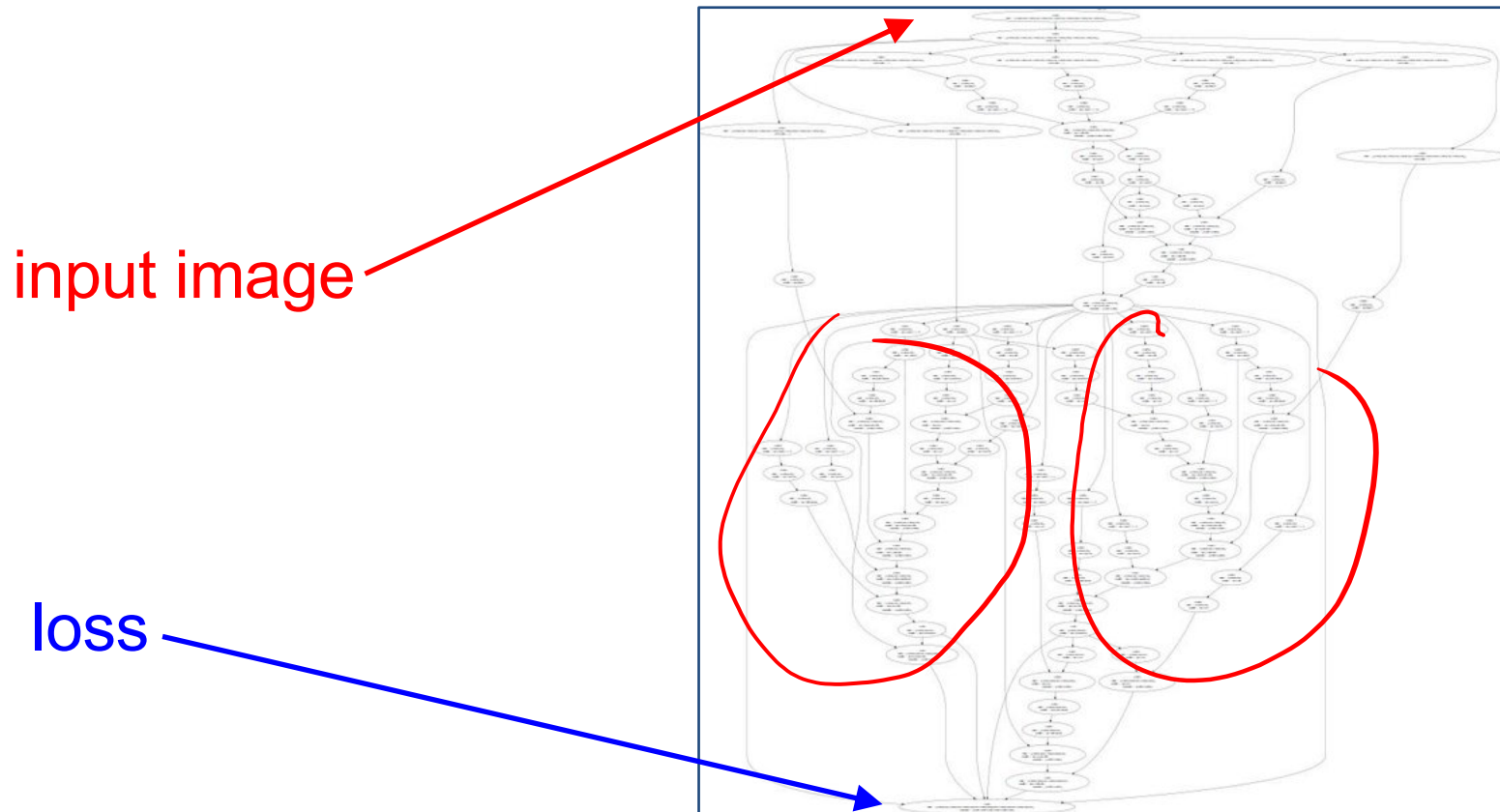
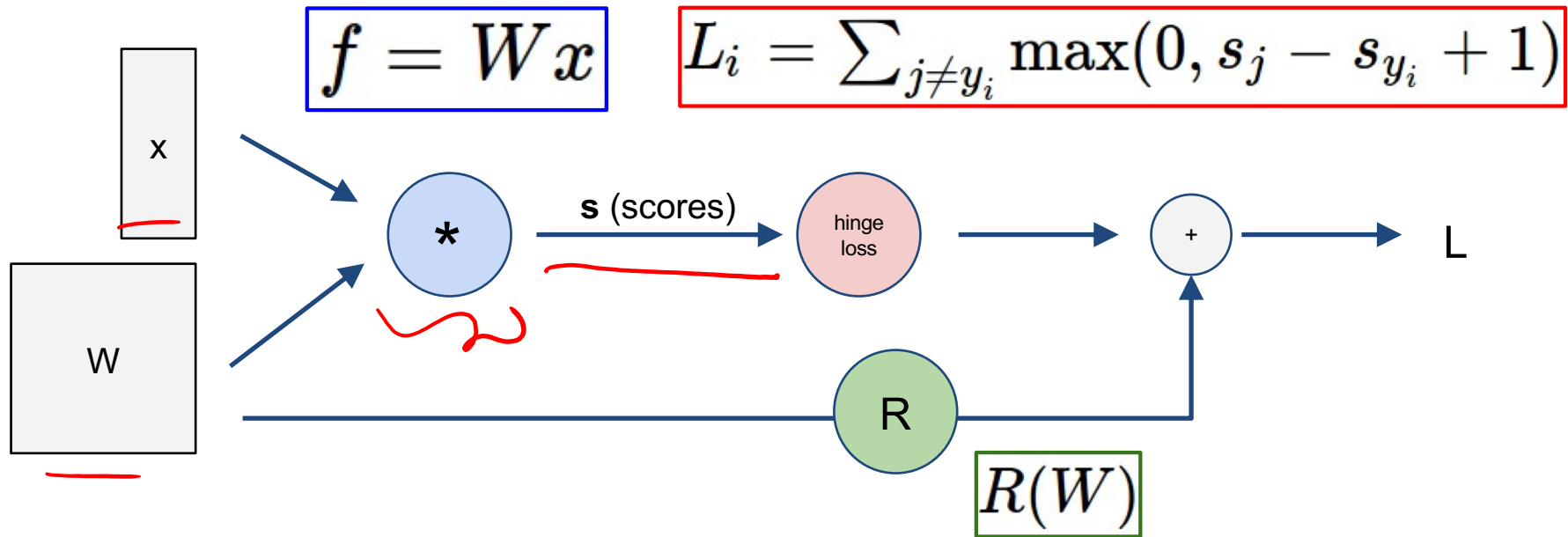


Figure reproduced with permission from a [Twitter post](#) by Andrej Karpathy.

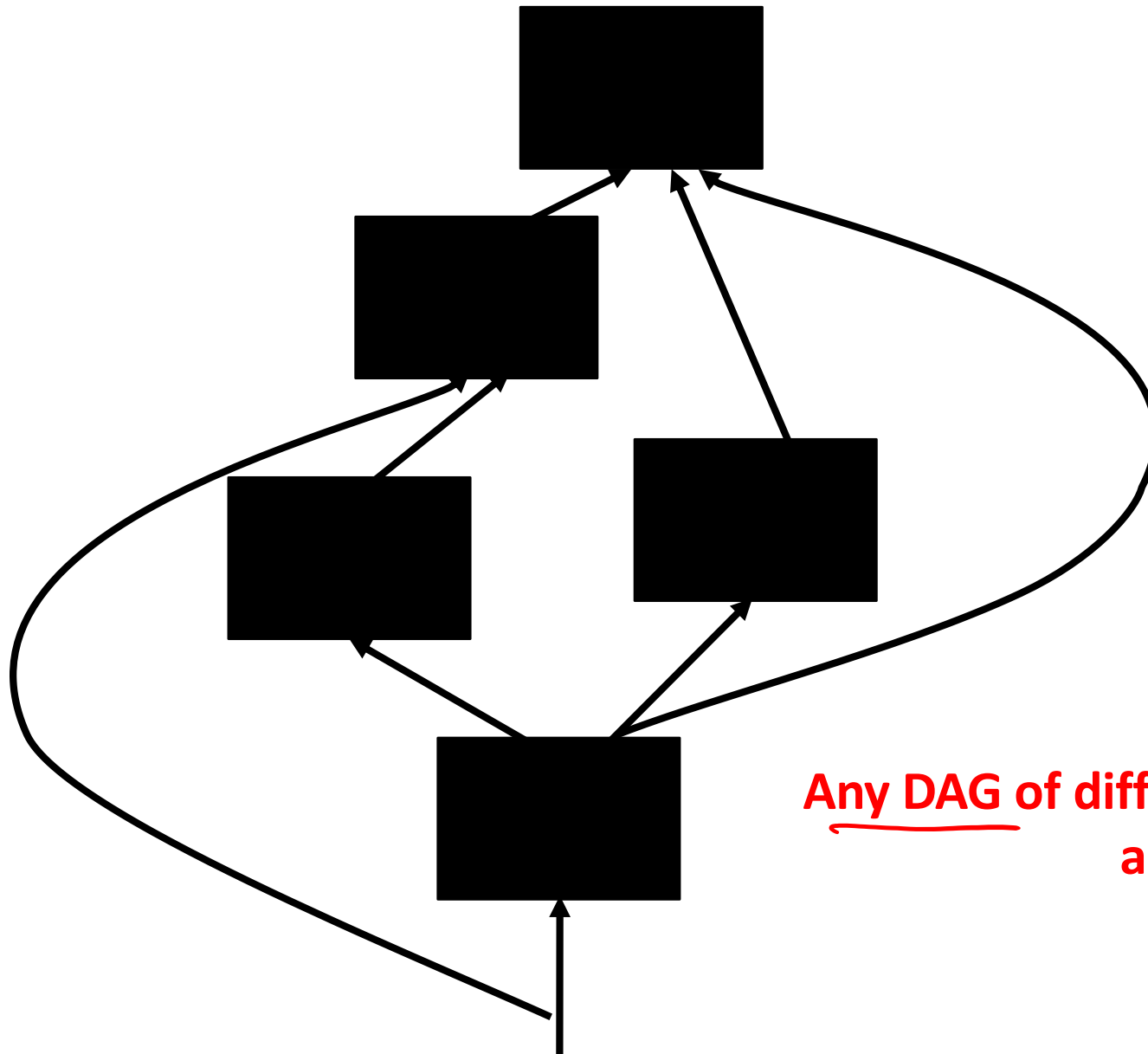
How do we compute gradients?

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

Computational Graph



Computational Graph



Any DAG of differentiable modules is allowed!

Directed Acyclic Graphs (DAGs)

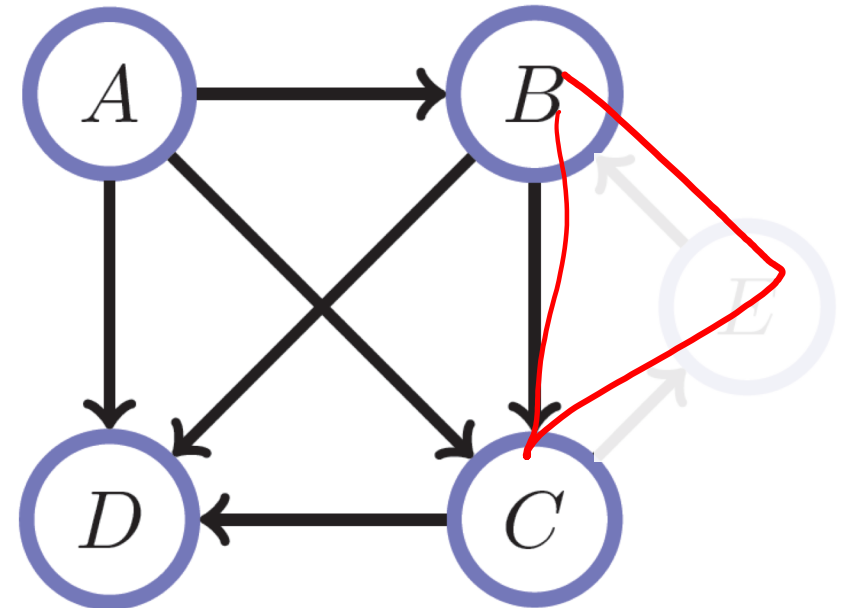
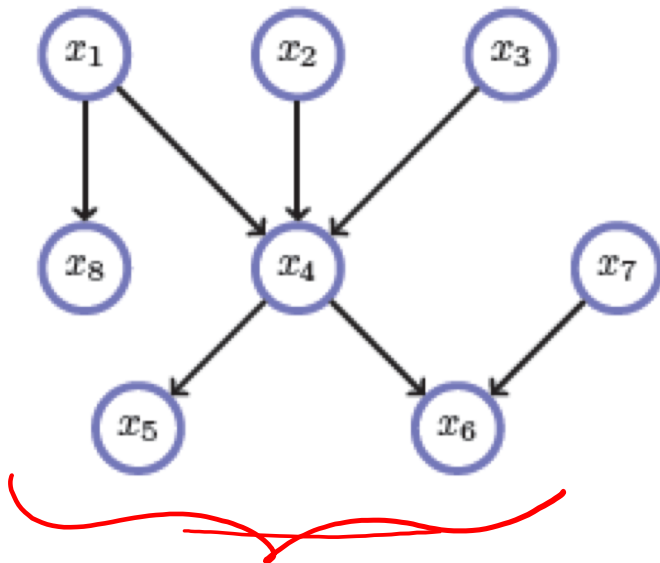
- Exactly what the name suggests

$$G = (V, E)$$
$$E = \{ (v_i, v_j) \mid v_i, v_j \in V \}$$

→ Directed edges

– No (directed) cycles

Underlying undirected cycles okay

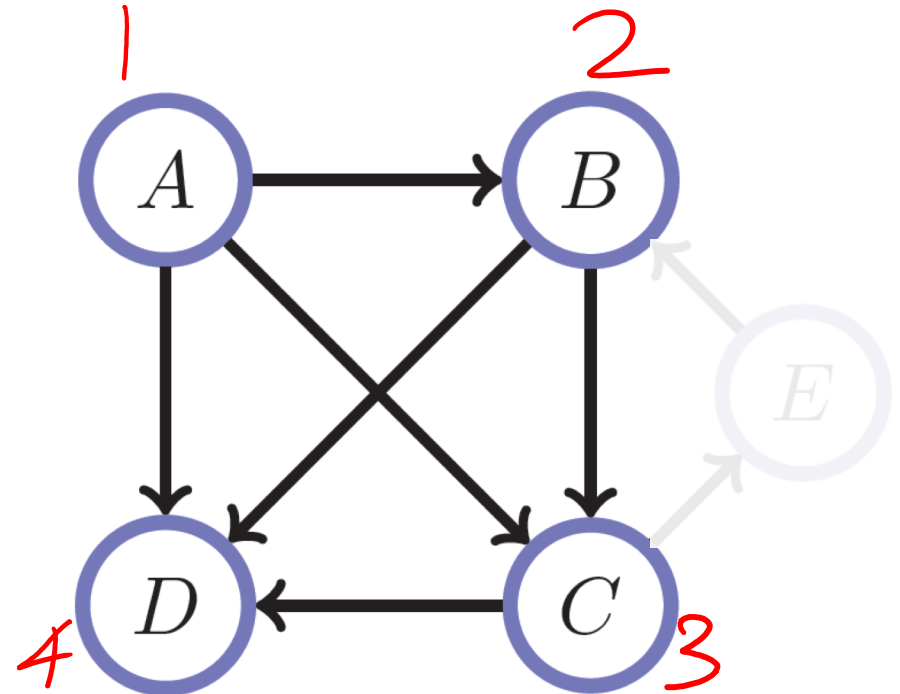
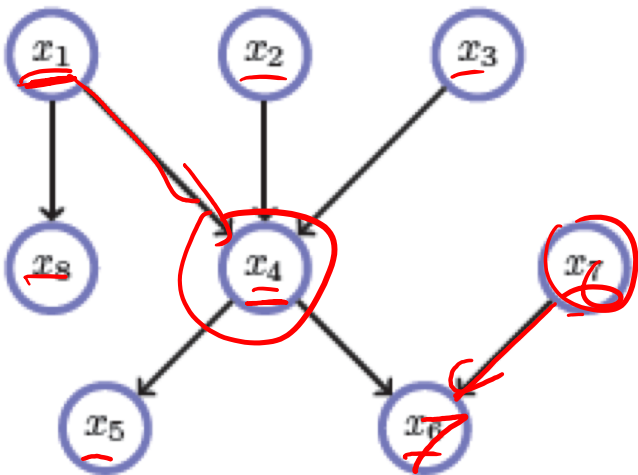


Directed Acyclic Graphs (DAGs)

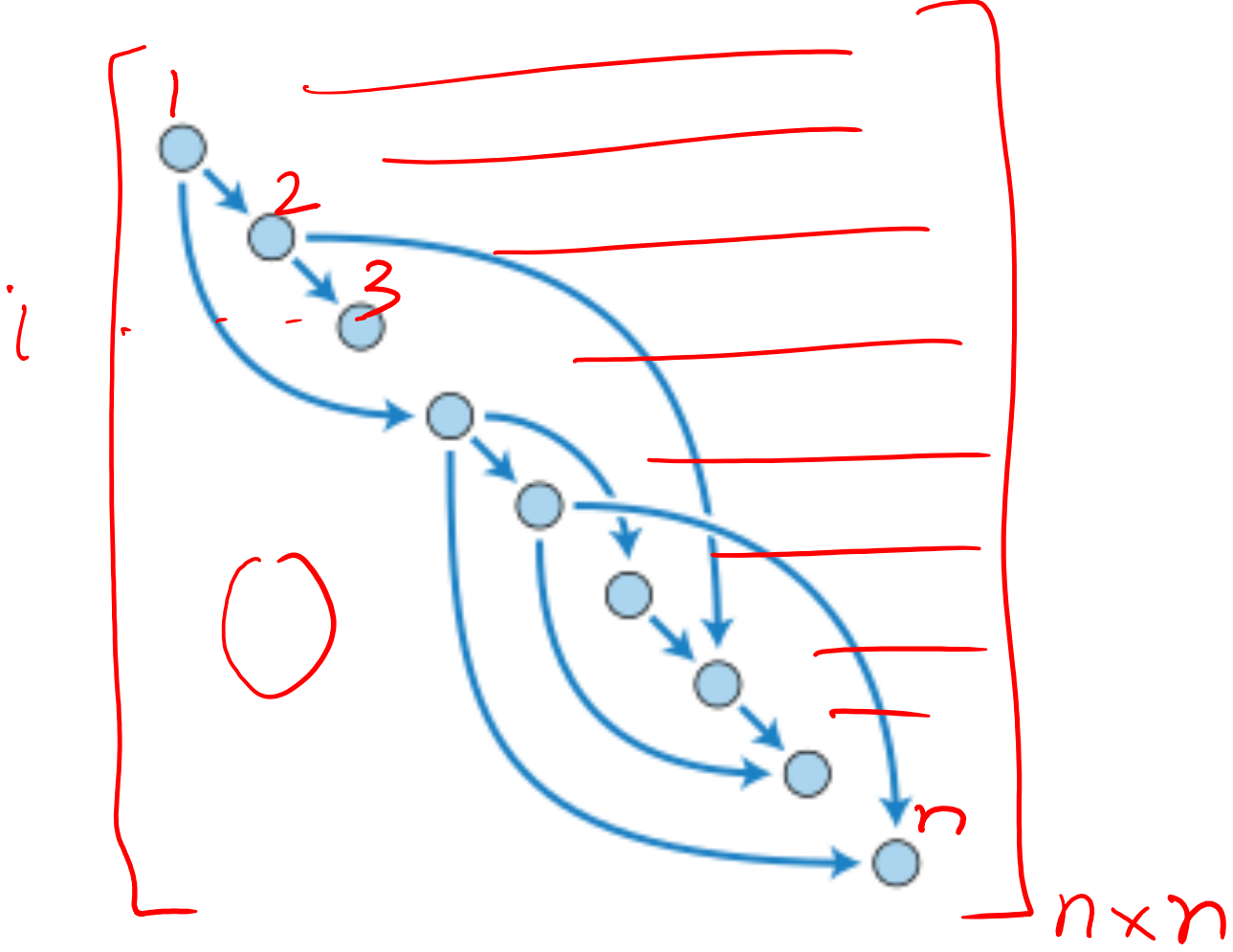
- Concept
 - Topological Ordering

$$\exists \sigma: V \rightarrow [n] = \{1, \dots, n\}$$

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



Directed Acyclic Graphs (DAGs)

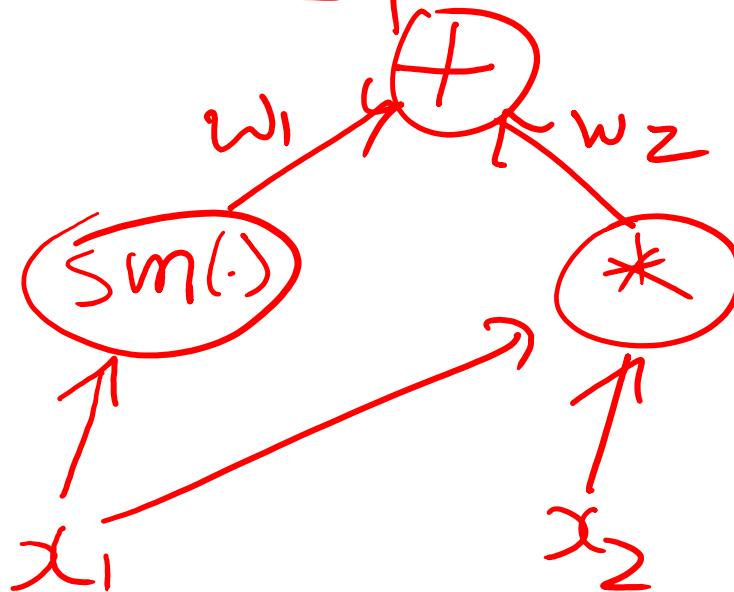


Computational Graphs

- Notation

$$f(x_1, x_2) = \underbrace{x_1 x_2}_{w_2} + \underbrace{\sin(x_1)}_{w_1}$$

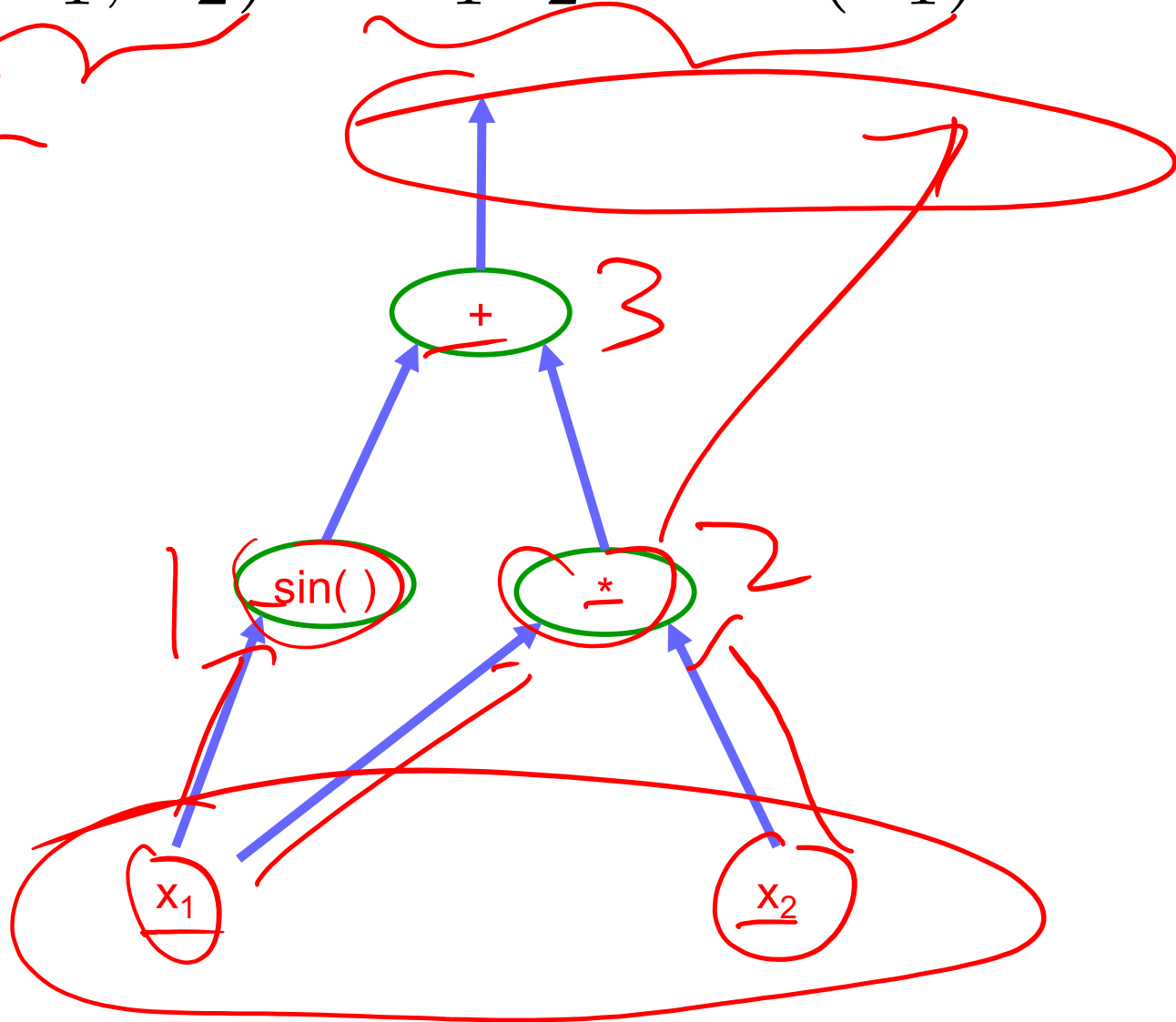
Handwritten red annotations:
 $f(x_1, x_2) = w_2 \uparrow \underbrace{x_1 x_2}_{w_2} + \underbrace{\sin(x_1)}_{w_1}$



Example

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

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

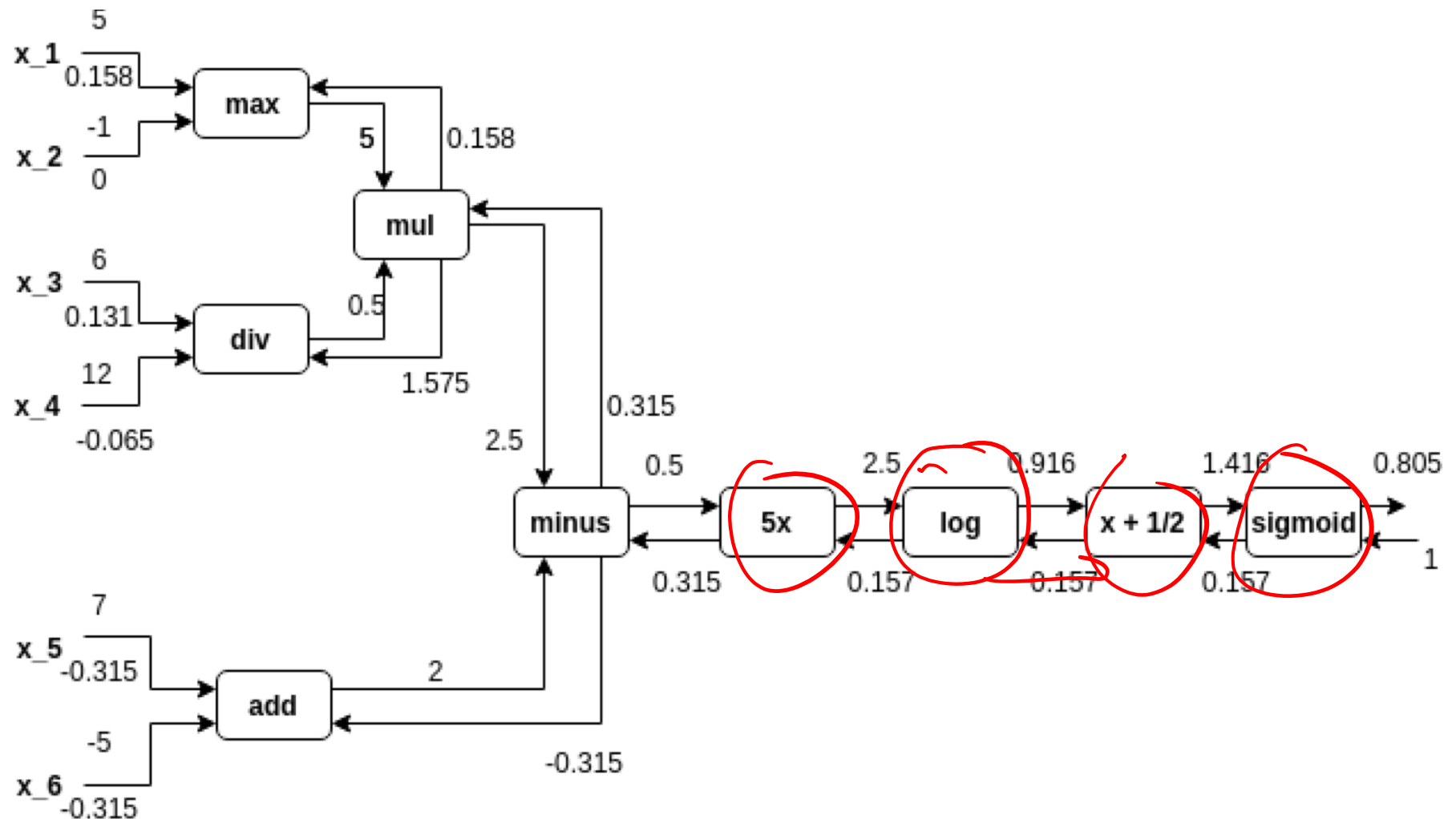


HW0

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

HW0 Submission by Samyak Datta

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

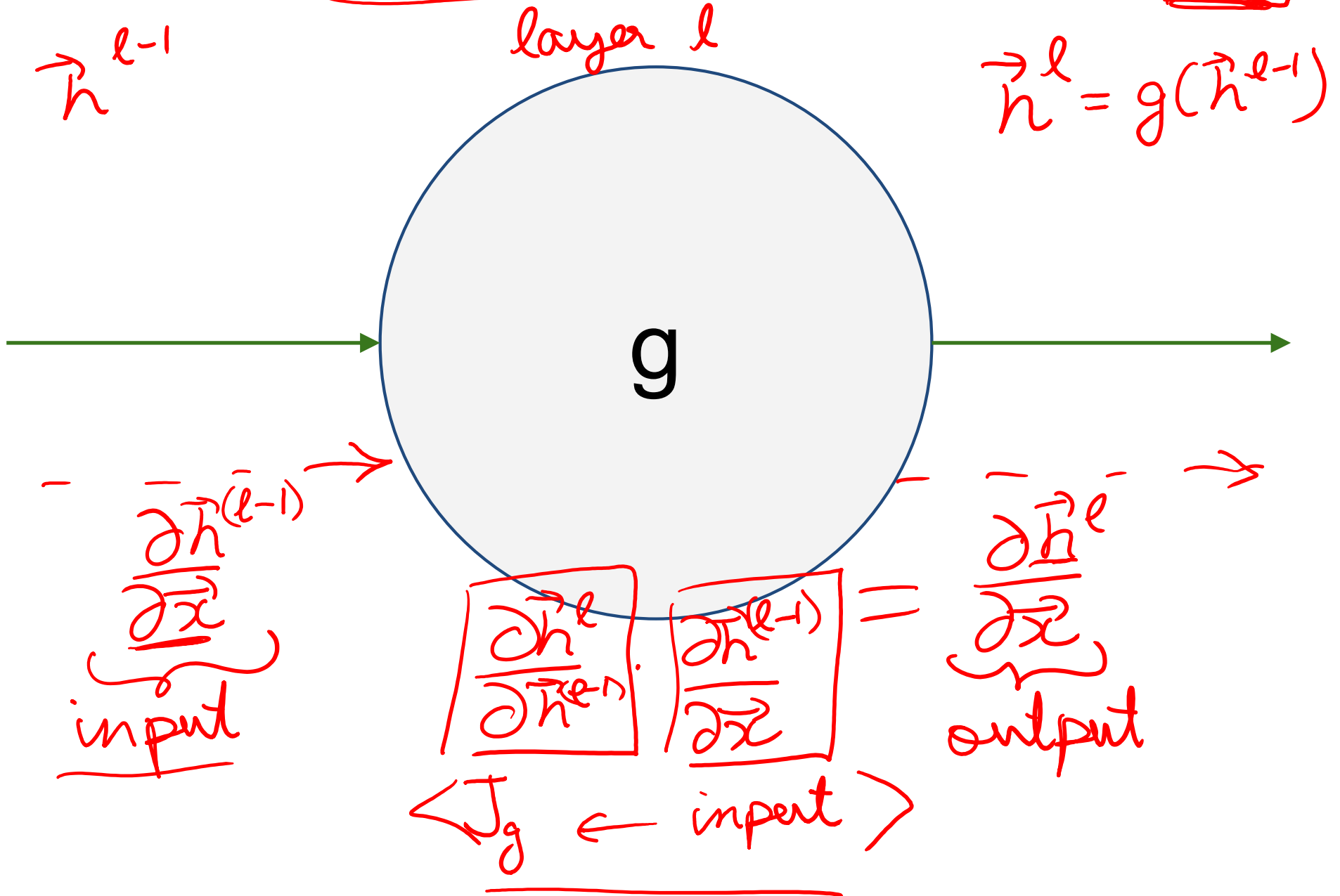


Forward mode vs Reverse Mode

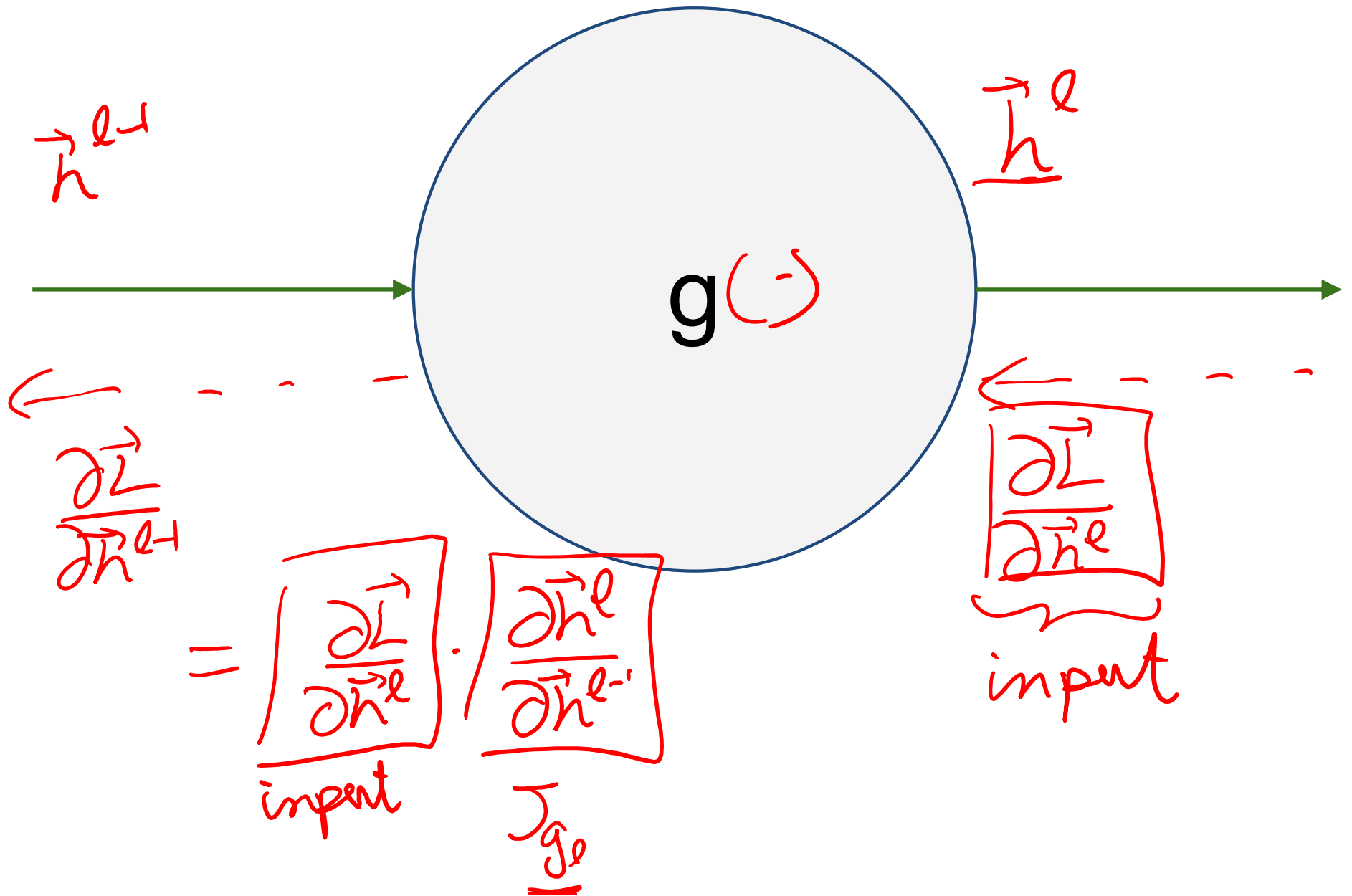
- Key Computations

Forward mode AD

$$\frac{\partial \mathcal{L}}{\partial \vec{x}}$$



Reverse mode AD



Example: Forward mode AD

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

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

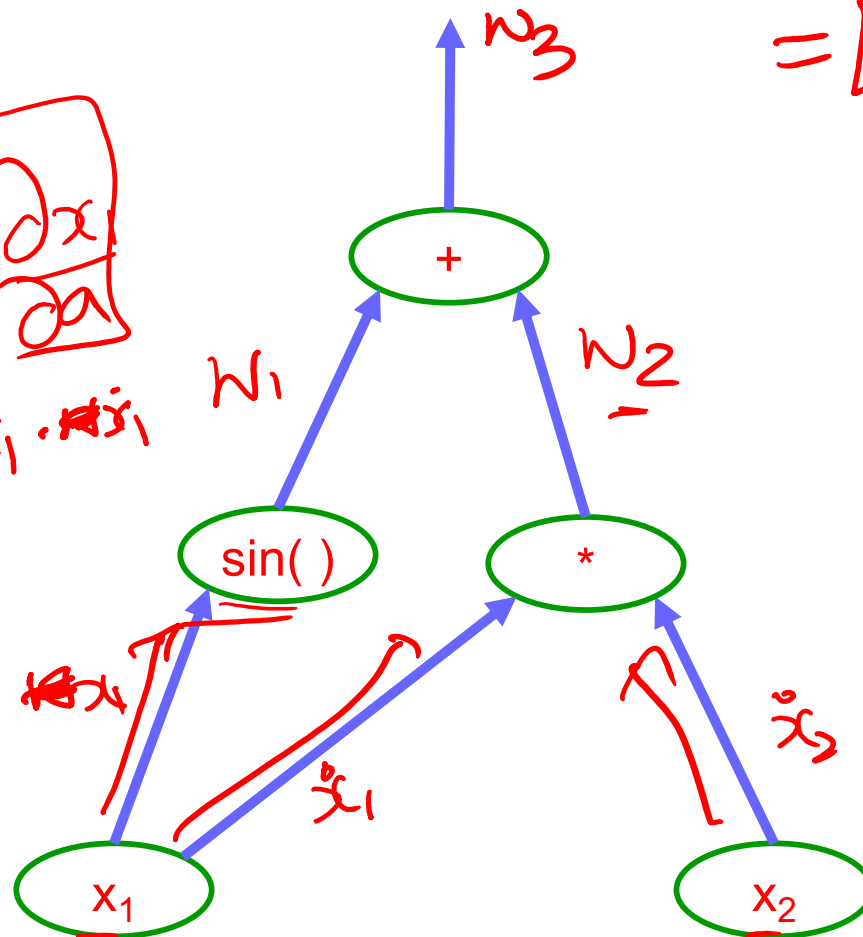
$$= \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \right]$$

$$\frac{\partial f}{\partial a} = \dot{f}$$

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

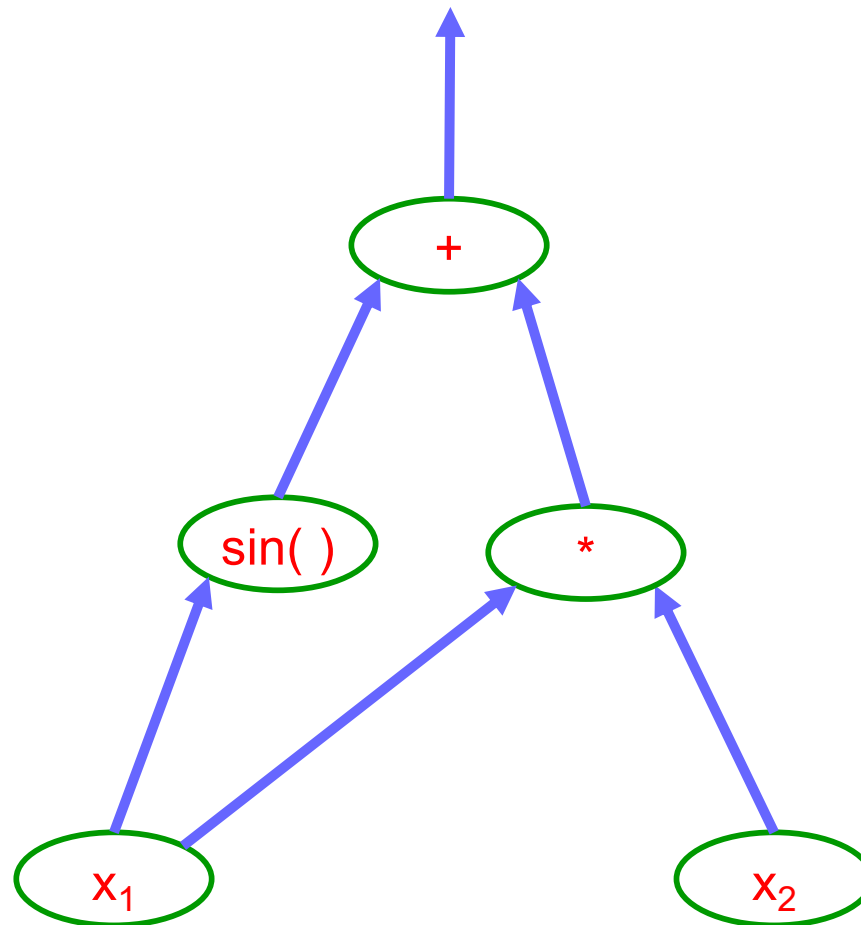
$$\dot{w}_1 = \frac{\partial w_1}{\partial a} = \frac{\partial w_1}{\partial x_1} \cdot \frac{\partial x_1}{\partial a} = \cos(x_1) \cdot \dot{x}_1$$

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



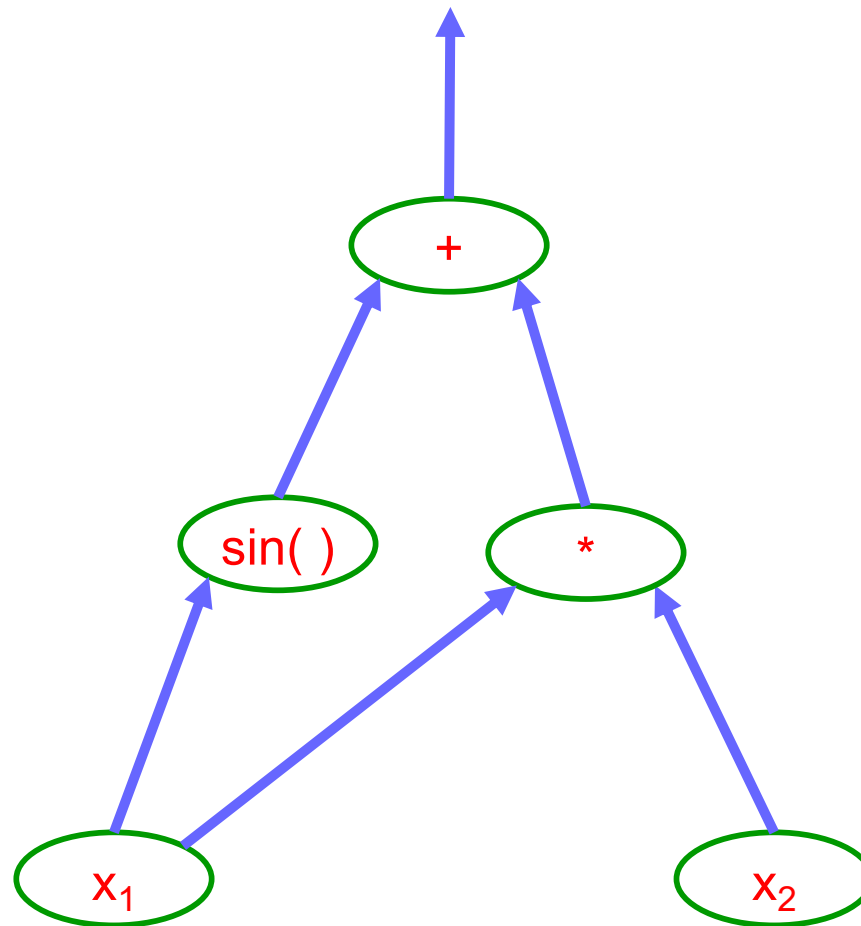
Example: Forward mode AD

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



Example: Forward mode AD

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



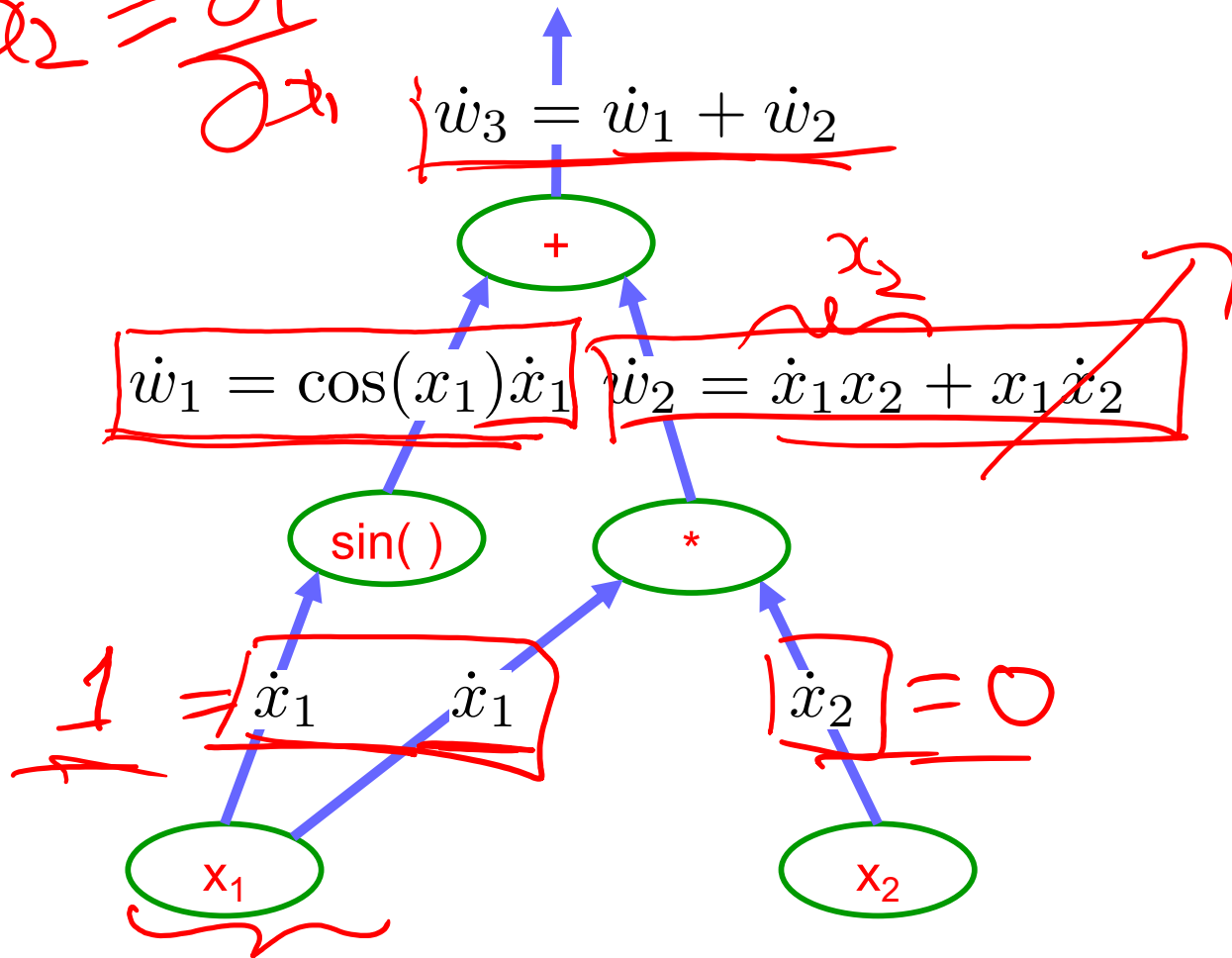
Example: Forward mode AD

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

$\frac{\partial f}{\partial a}$

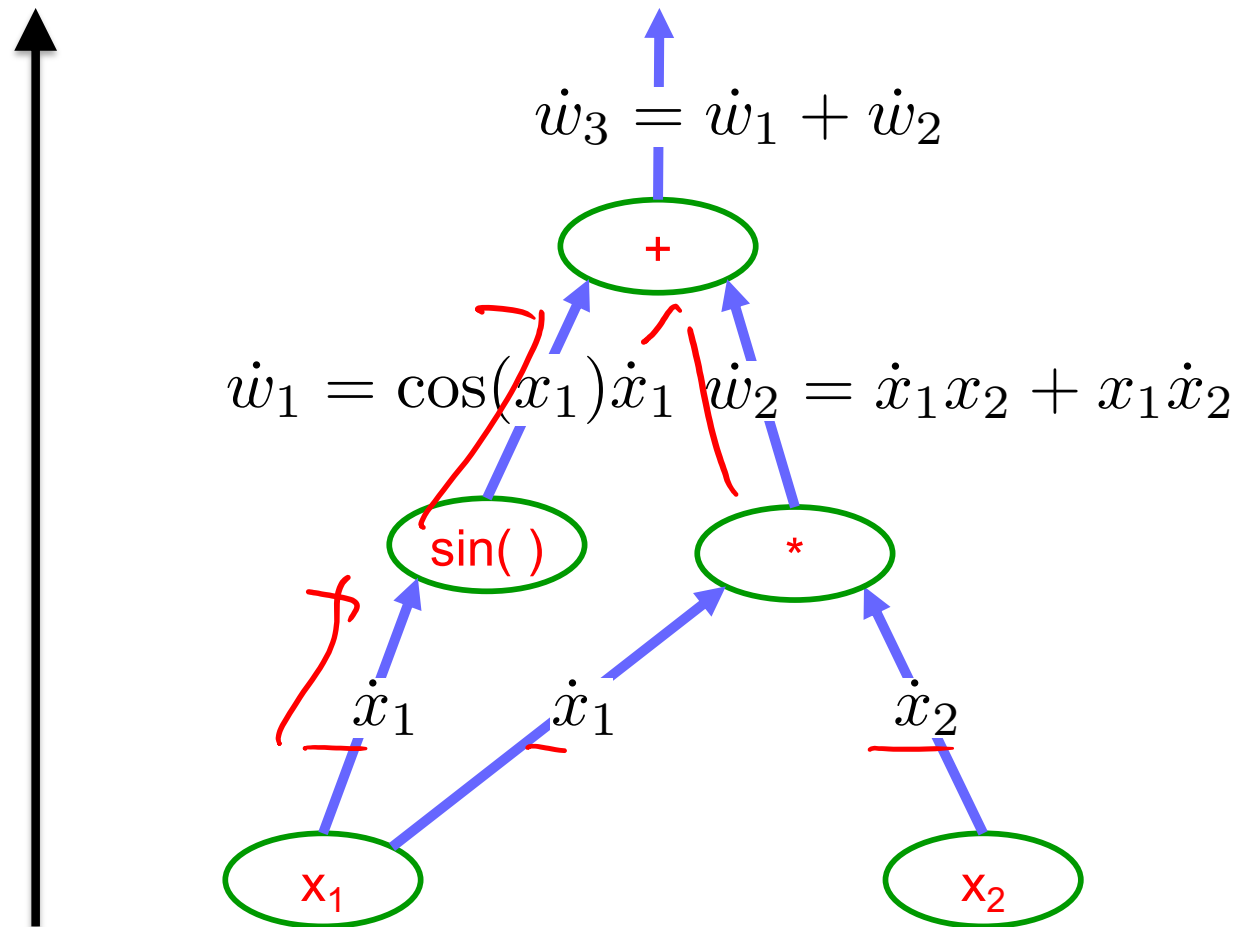
$\cos(x_1) + x_2 = \frac{\partial f}{\partial x_1}$

$\cos(x_1)$



Example: Forward mode AD

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



Example: Reverse mode AD

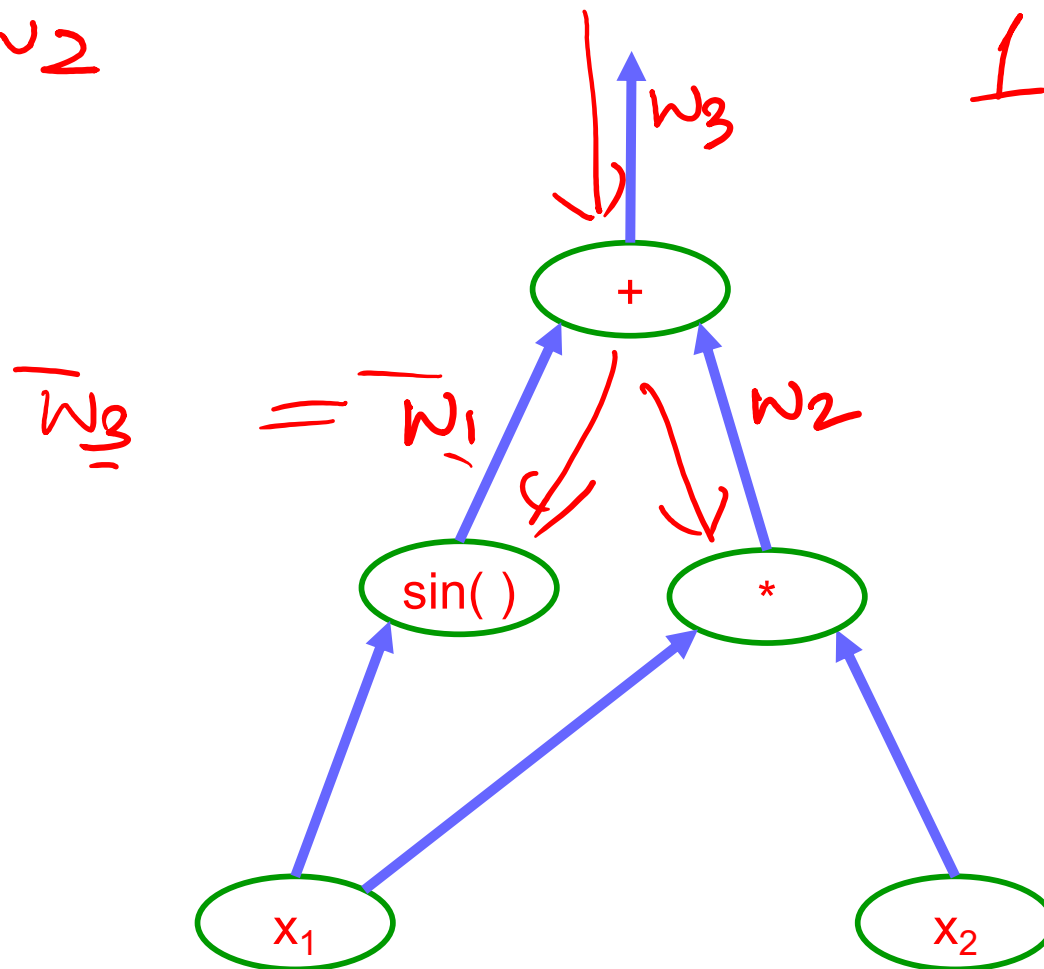
$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

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

$$w_3 = w_1 + w_2$$

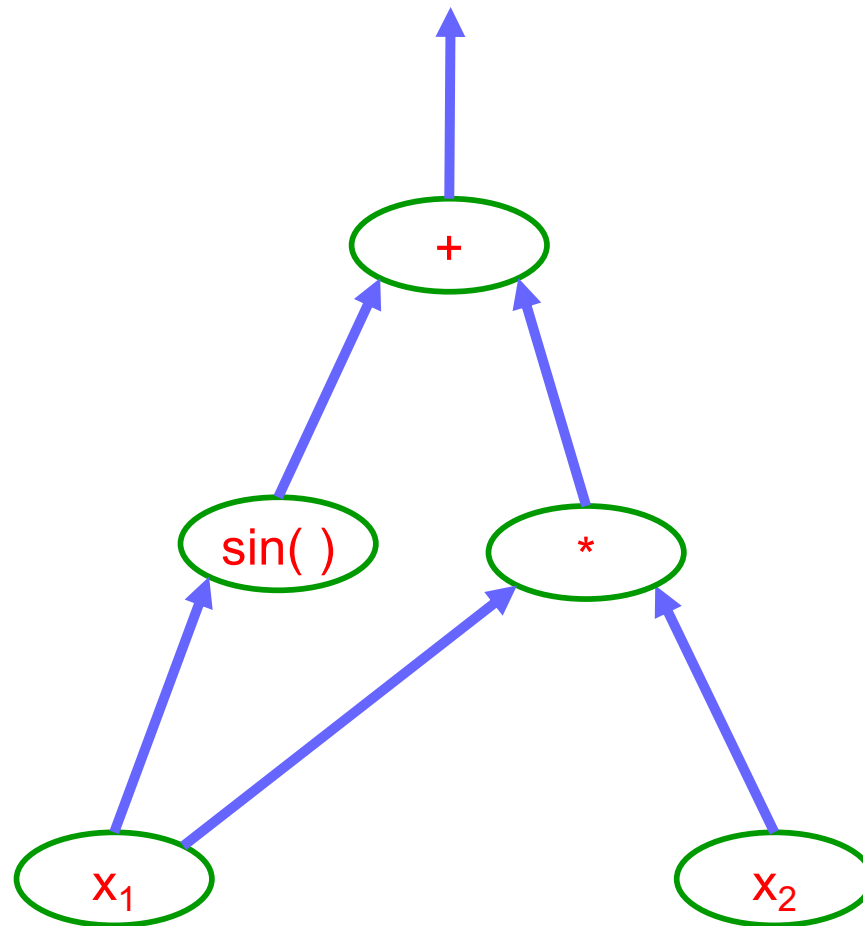
$$1 = \bar{w}_3 = \frac{\partial f}{\partial w_3}$$

$$\bar{x}_1 = \frac{\partial f}{\partial x_1}$$



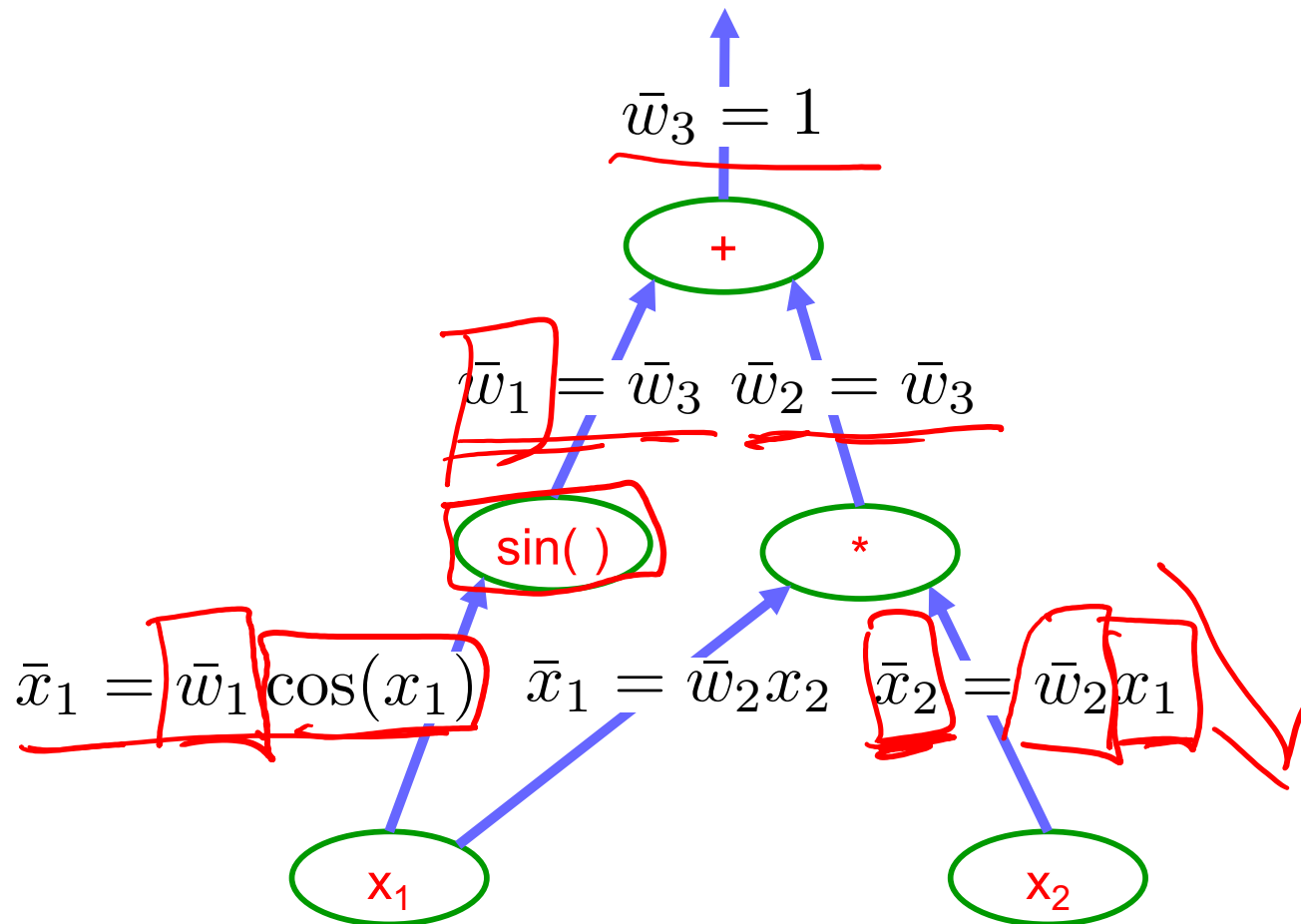
Example: Reverse mode AD

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



Example: Reverse mode AD

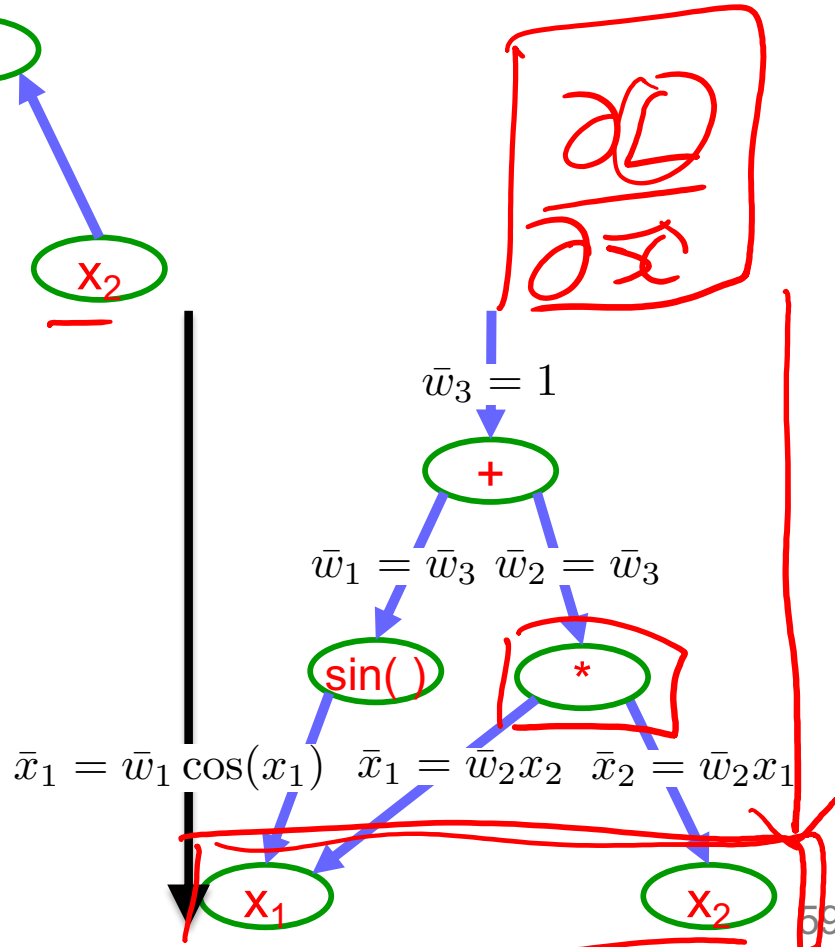
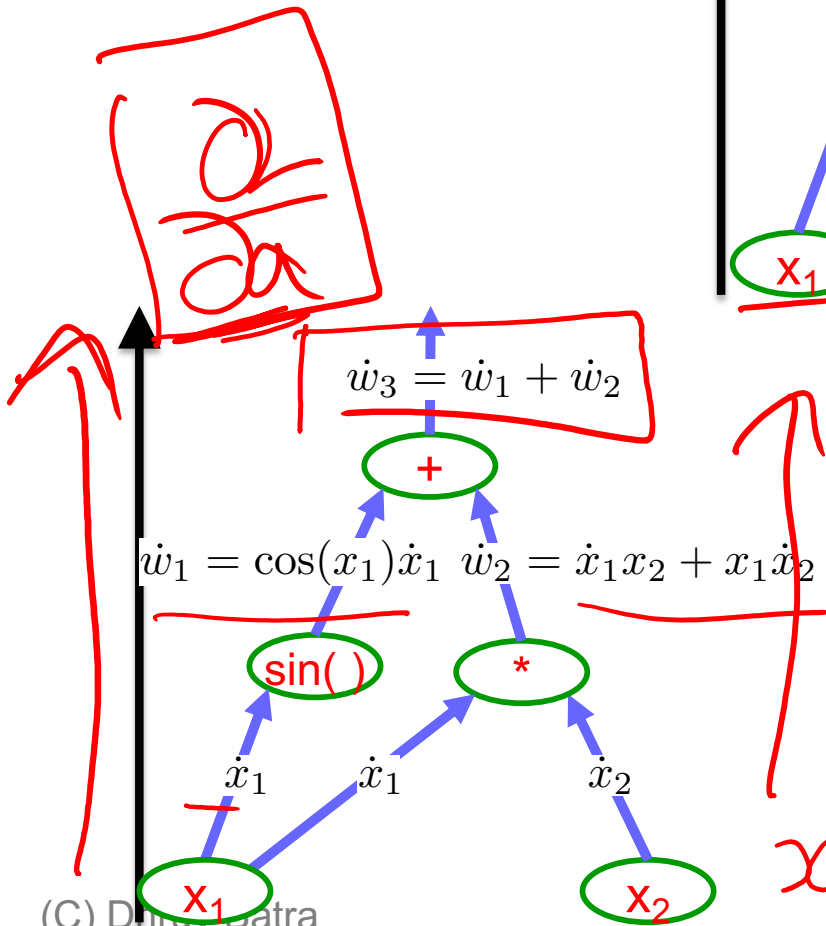
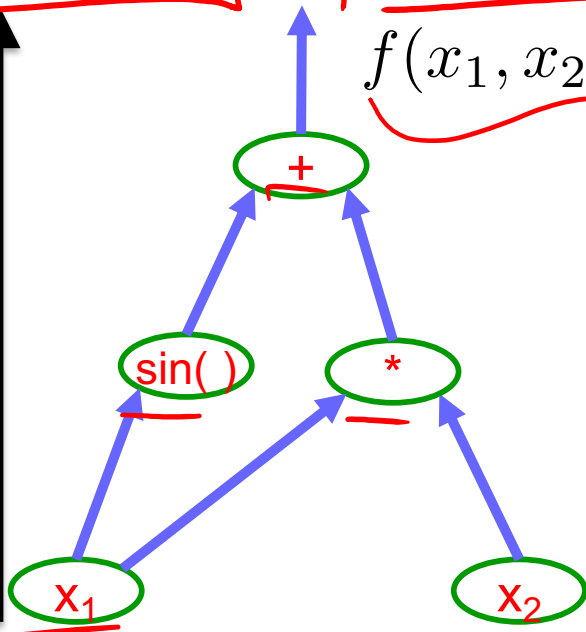
$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$



Forward Pass vs

Forward mode AD vs Reverse Mode AD

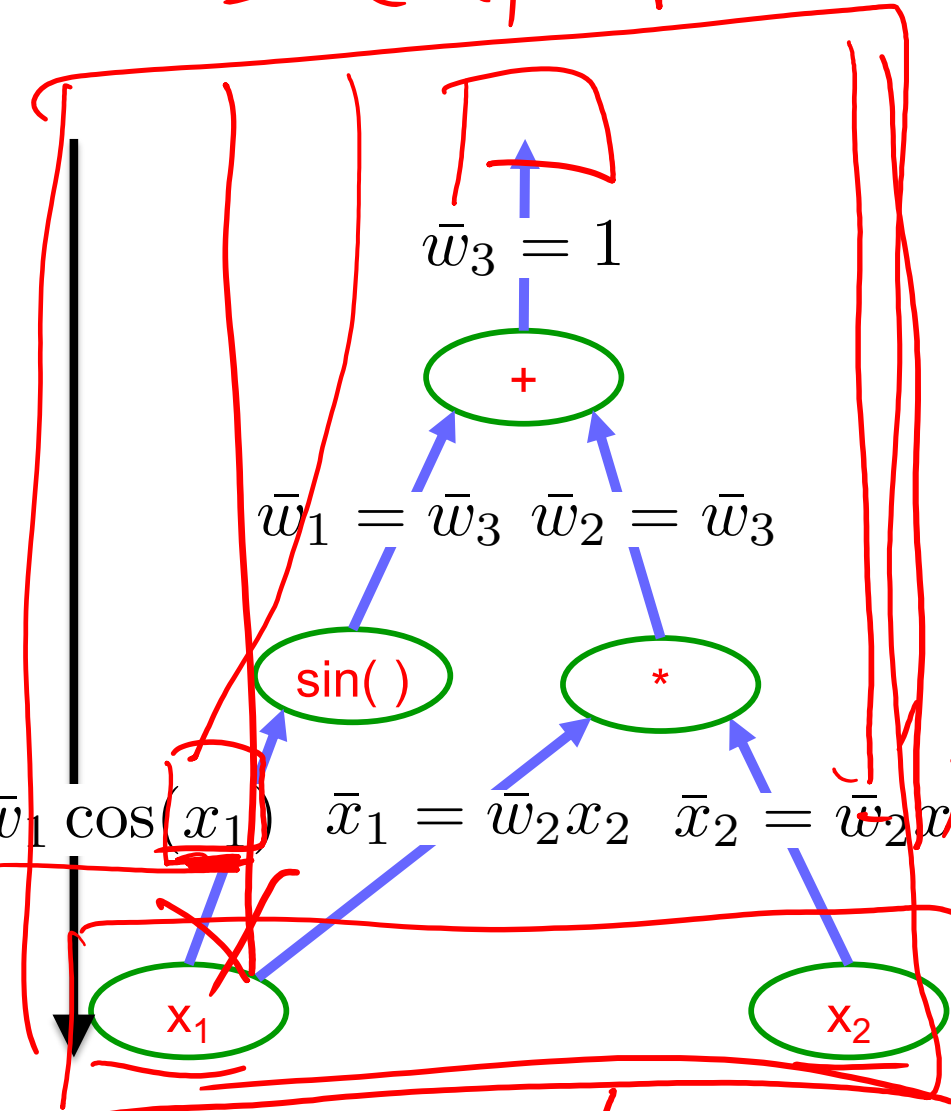
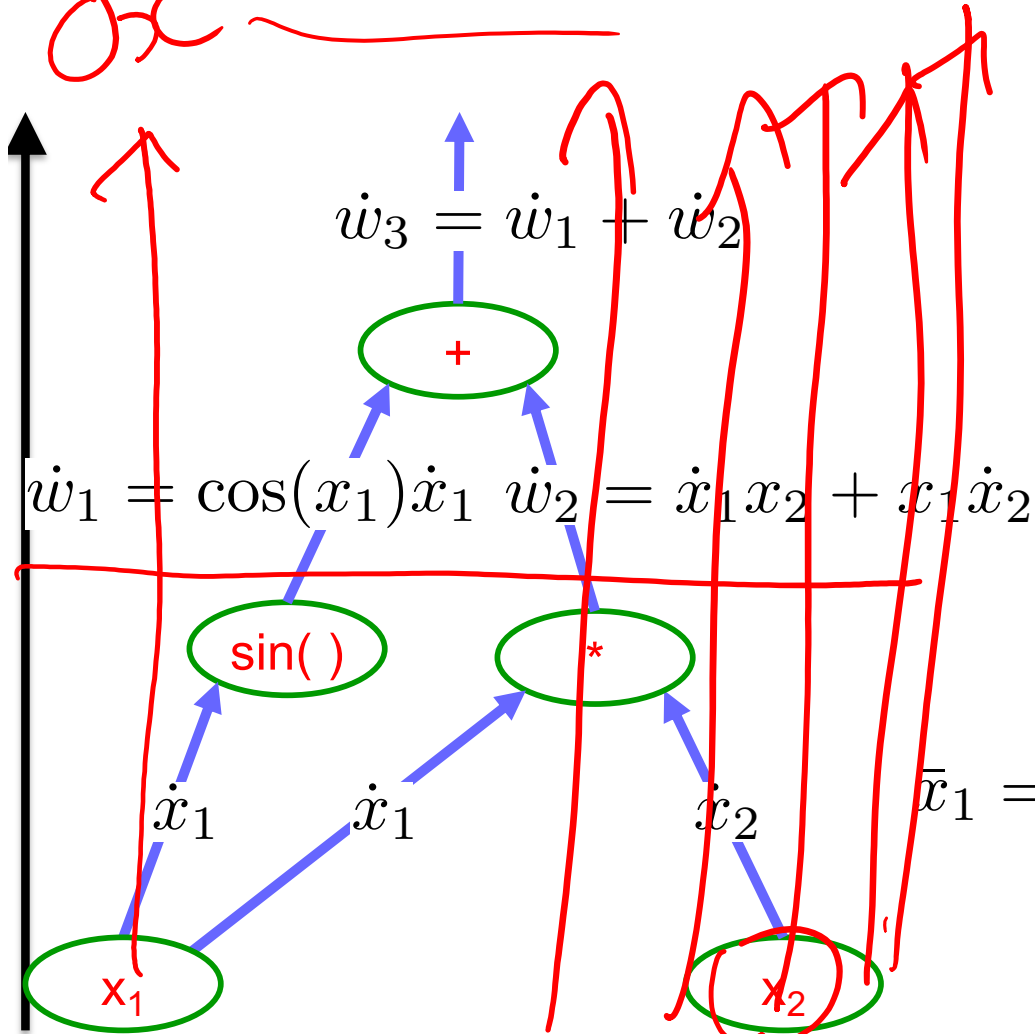
$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$



Forward mode vs Reverse Mode

What are the differences?

backprop



$\bar{x} \in \mathbb{R}^d$

$L \in \mathbb{R}^1$

Forward mode vs Reverse Mode

- What are the differences?

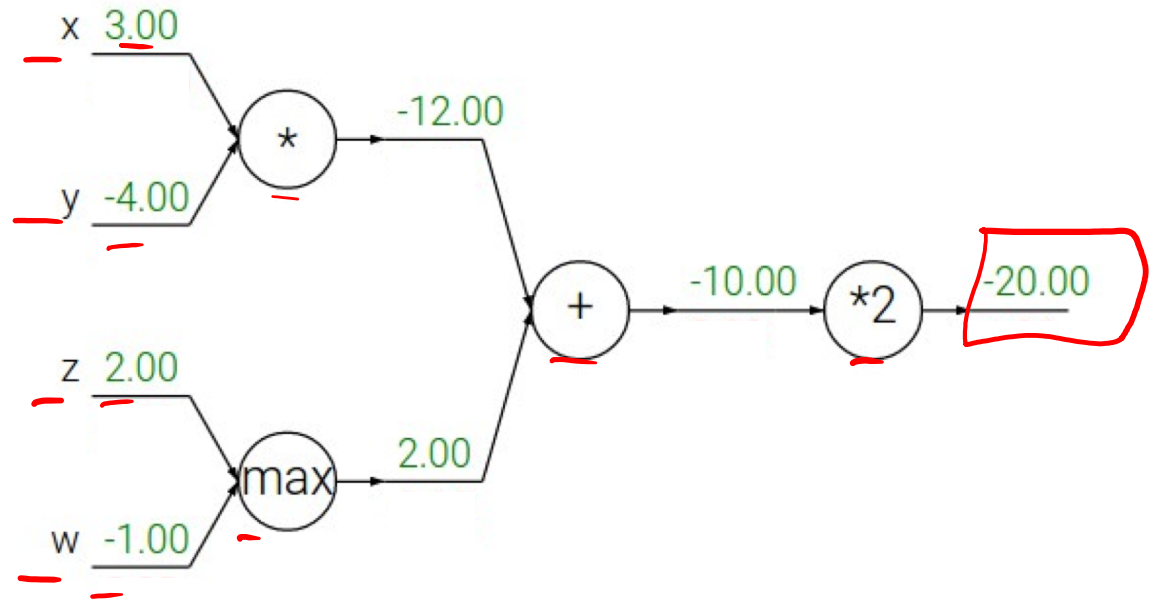
- Which one is faster to compute?
 - Forward or backward?

Forward mode vs Reverse Mode

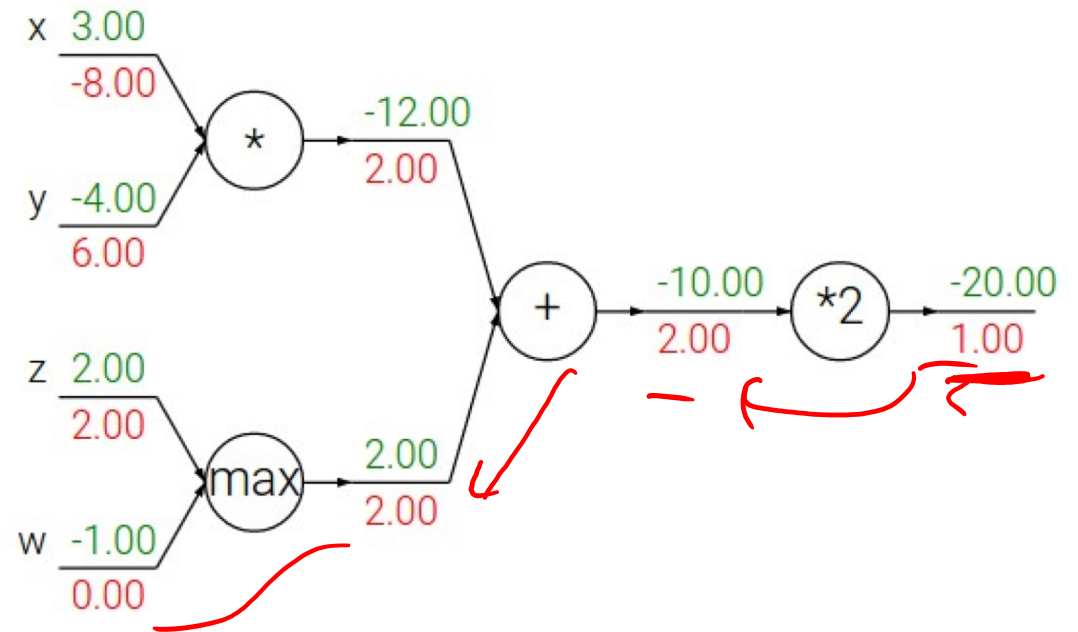
- What are the differences?
- Which one is faster to compute?
 - Forward or backward?
- Which one is more memory efficient (less storage)?
 - Forward or backward?

Patterns in backward flow

$$f(\cdot, \cdot) = 2(xy + \max\{z, w\})$$

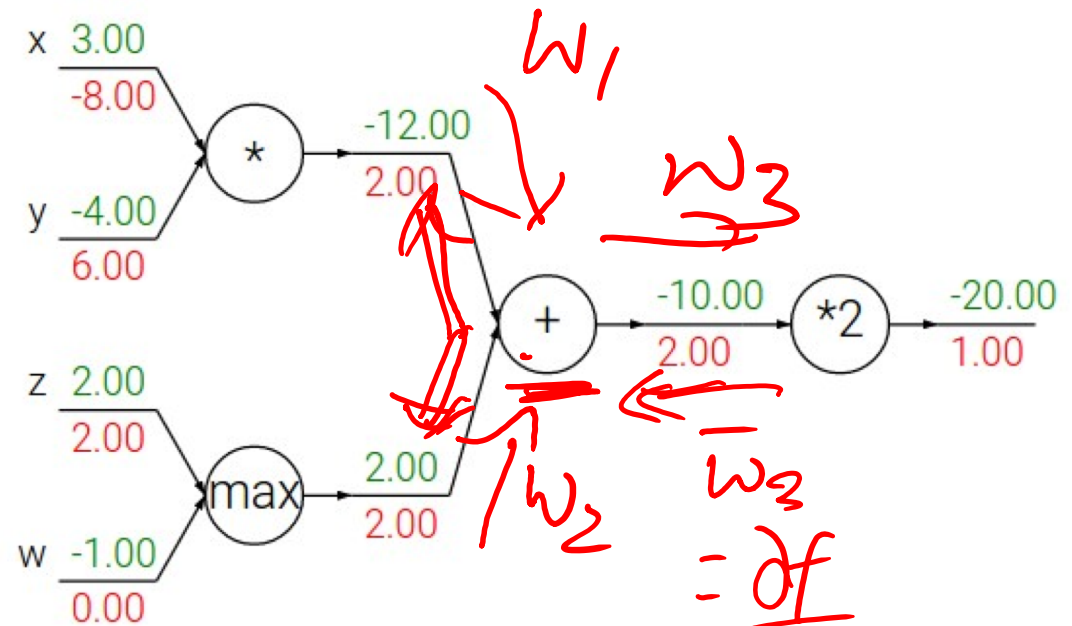


Patterns in backward flow



Patterns in backward flow

add gate gradient distributor



$$w_3 = w_1 + w_2$$

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

$$\boxed{\frac{\partial f}{\partial w_1}} = \boxed{\frac{\partial f}{w_3}} \cdot \boxed{\frac{\partial w_3}{\partial w_1}}$$

Patterns in backward flow

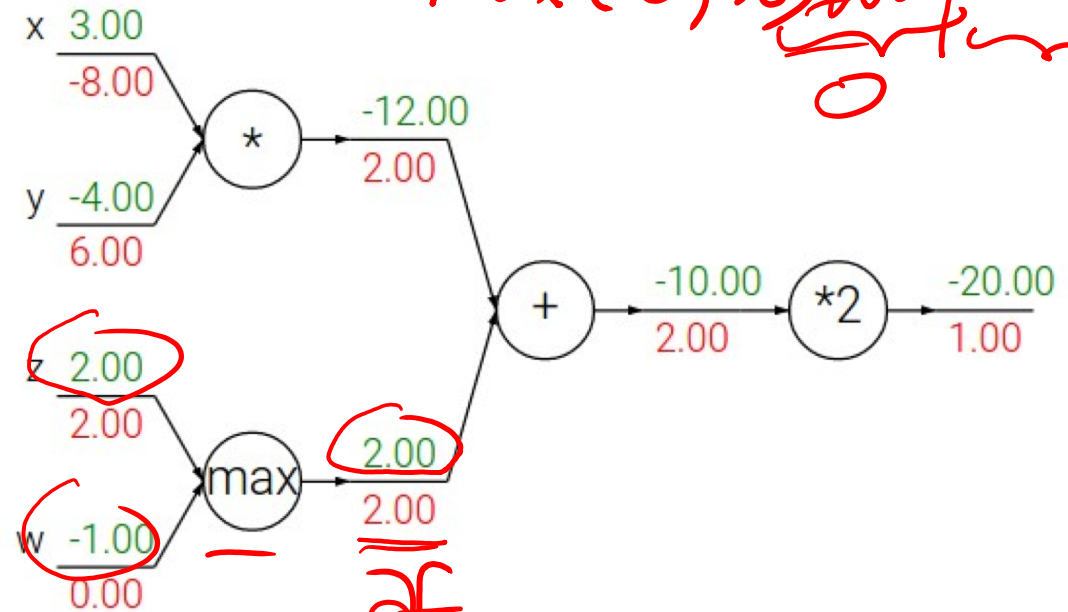
$$\underline{w}_2 = \max\{\underline{z}, \underline{w}\}$$

$$\max(0, x) \text{ derivative } \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$

add gate: gradient distributor

Q: What is a max gate?

$\frac{df}{dz}$
 $\frac{df}{dw}$

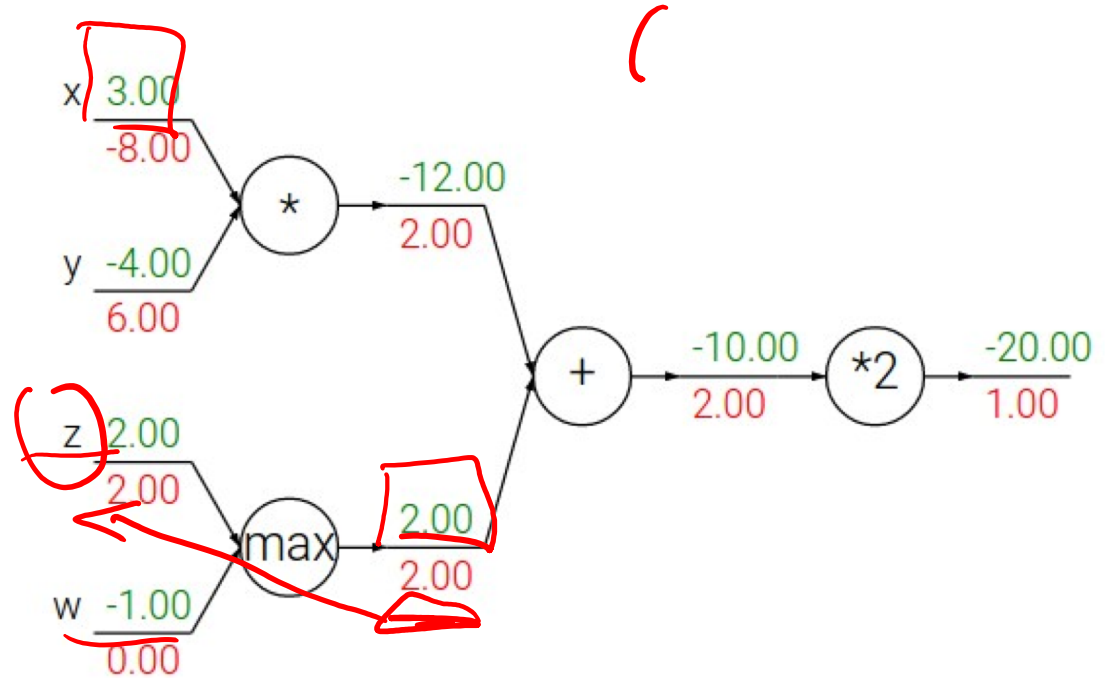


$\frac{df}{dw_2}$

Patterns in backward flow

add gate: gradient distributor

max gate: gradient router

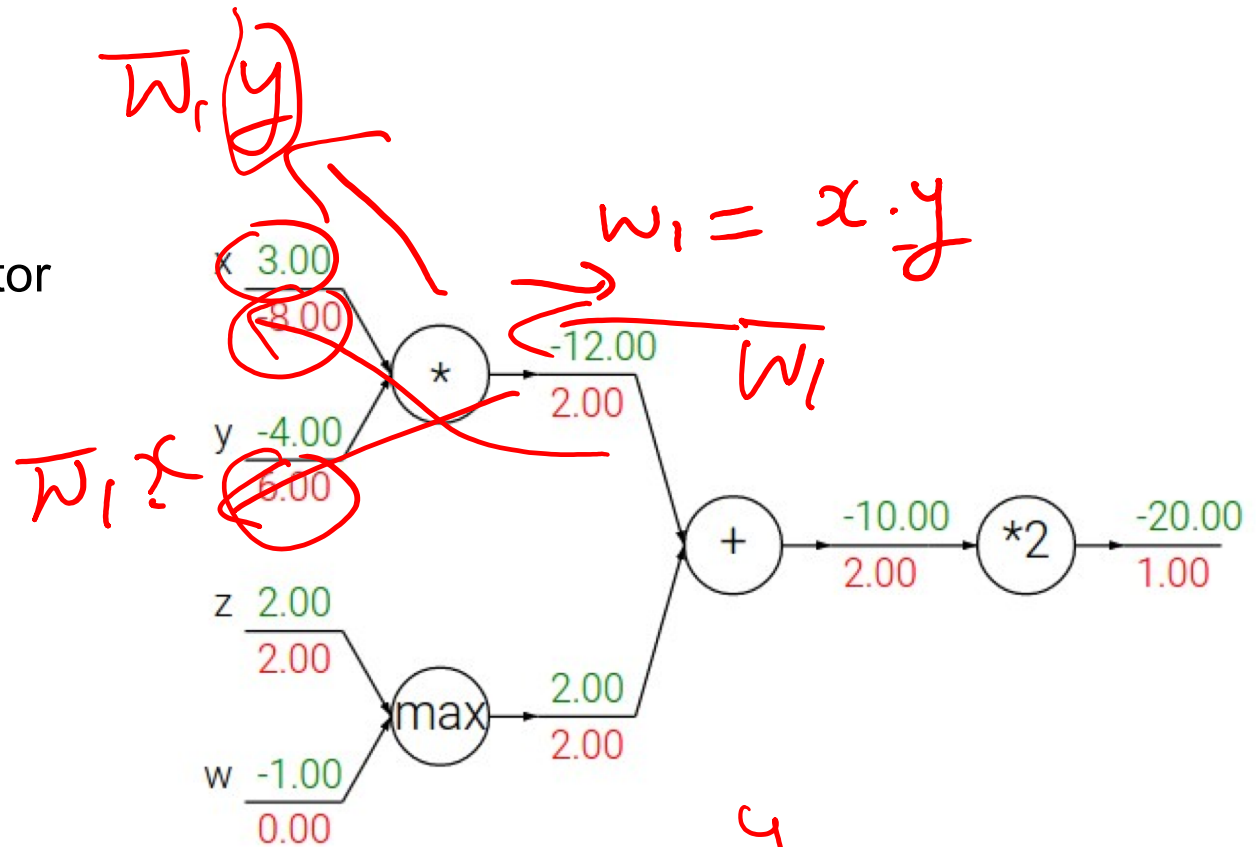


Patterns in backward flow

add gate: gradient distributor

max gate: gradient router

Q: What is a mul gate?



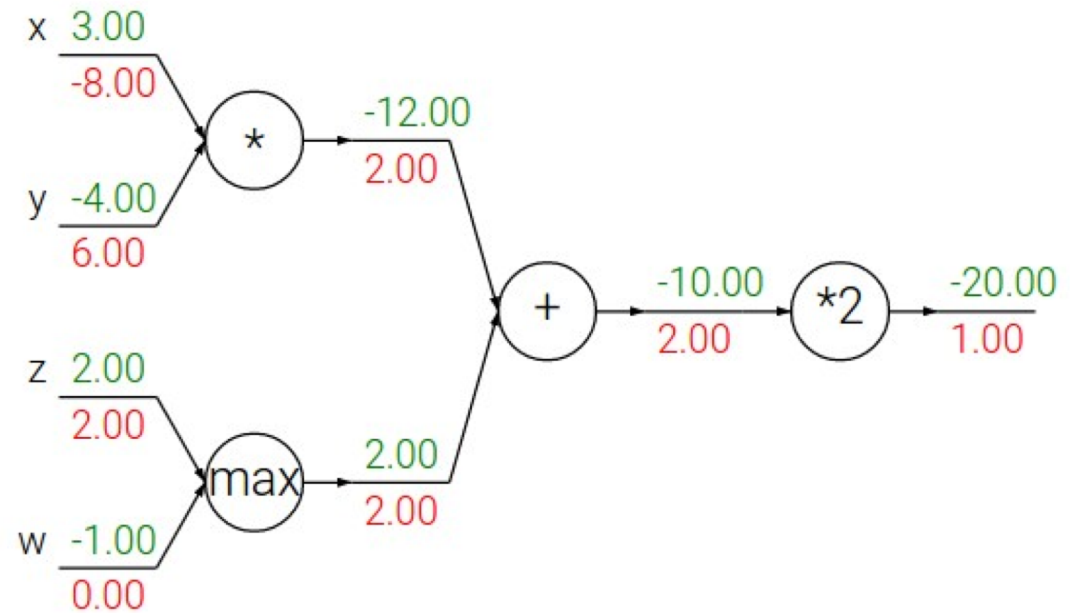
$$\frac{\partial f}{\partial x} = \left[\frac{\partial f}{\partial w_1} \right] \left[\frac{\partial w_1}{\partial x} \right]$$

Patterns in backward flow

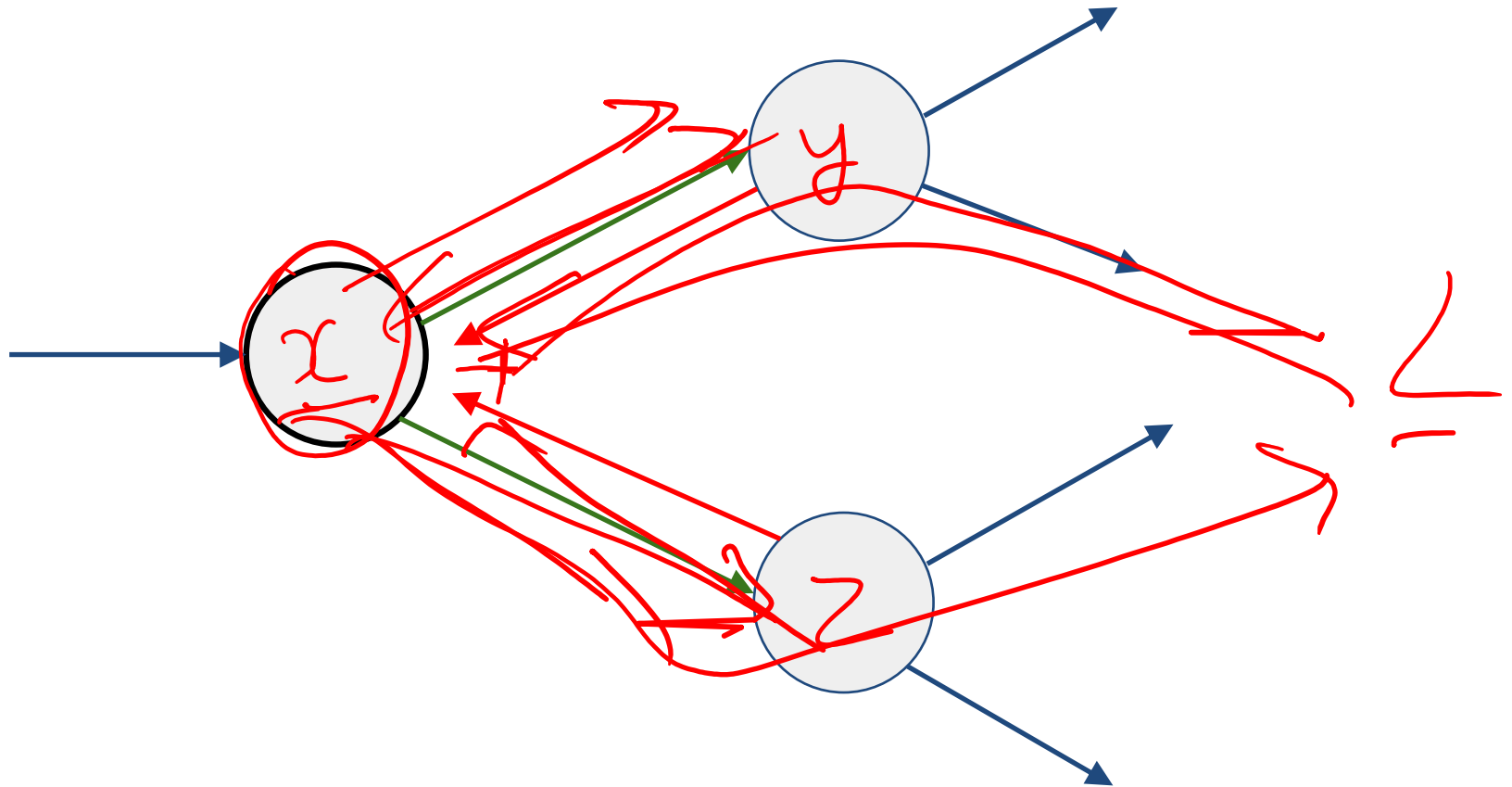
add gate: gradient distributor

max gate: gradient router

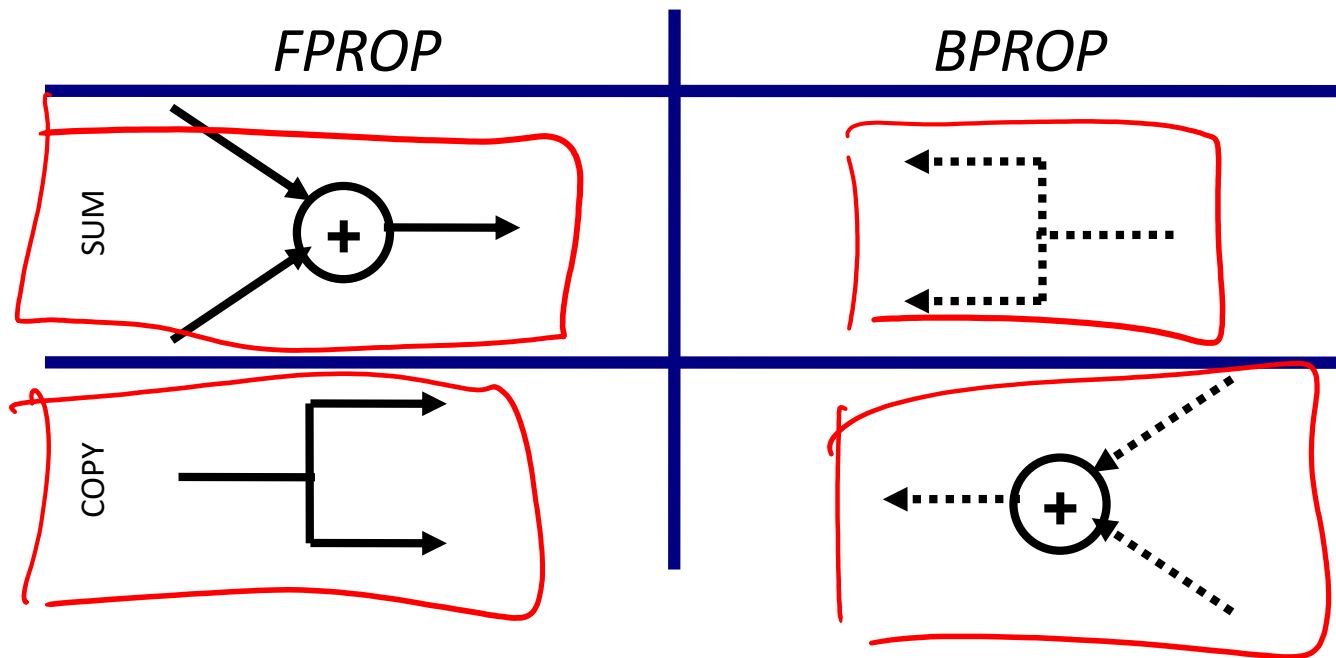
mul gate: gradient switcher



Gradients add at branches

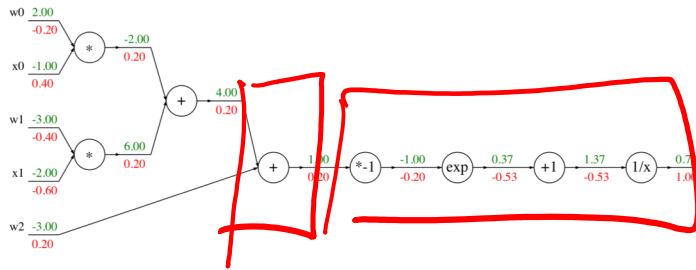


Duality in Fprop and Bprop



Modularized implementation: forward / backward API

Graph (or Net) object *(rough psuedo code)*



```
class ComputationalGraph(object):
```

```
#...
```

```
def forward(inputs):
```

```
# 1. [pass inputs to input gates...]
```

```
# 2. forward the computational graph:
```

```
for gate in self.graph.nodes_topologically_sorted():
```

```
    gate.forward()
```

```
return loss # the final gate in the graph outputs the loss
```

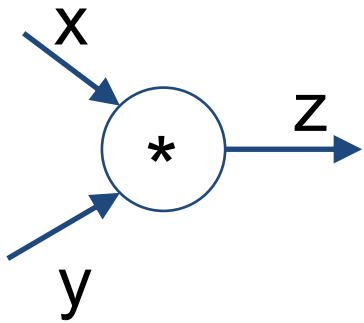
```
def backward():
```

```
for gate in reversed(self.graph.nodes_topologically_sorted()):
```

```
    gate.backward() # little piece of backprop (chain rule applied)
```

```
return inputs_gradients
```

Modularized implementation: forward / backward API



(x,y,z are scalars)

```
class MultiplyGate(object):
```

```
    def forward(x,y):
```

```
        z = x*y
```

```
        return z
```

```
    def backward(dz):
```

```
        # dx = ... #todo
```

```
        # dy = ... #todo
```

```
        return [dx, dy]
```

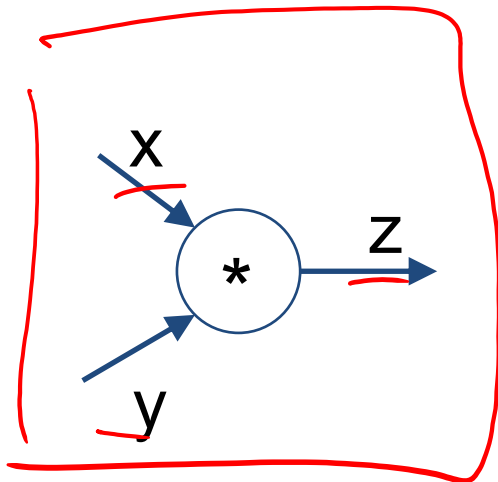
$$\frac{\partial L}{\partial z}$$

An arrow points from this box to the `dz` parameter in the `backward` method definition.

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

An arrow points from this box to the `dx` element in the `return` statement of the `backward` method.

Modularized implementation: forward / backward API



(x,y,z are scalars)

```
class MultiplyGate(object):  
    def forward(x,y):  
        z = x*y  
        self.x = x # must keep these around!  
        self.y = y  
        return z  
    def backward(dz):  
        dx = self.y * dz # [dz/dx * dL/dz]  
        dy = self.x * dz # [dz/dy * dL/dz]  
        return [dx, dy]
```


Example: Caffe layers

Branch: master | [caffe / src / caffe / layers /](#) | [Create new file](#) | [Upload files](#) | [Find file](#) | [History](#)

shelhamer committed on GitHub Merge pull request #4630 from BiGene/load_hdf6_fix Latest commit e687a71 21 days ago

absval_layer.cpp	dismantle layer headers	a year ago
absval_layer.cu	dismantle layer headers	a year ago
accuracy_layer.cpp	dismantle layer headers	a year ago
argmax_layer.cpp	dismantle layer headers	a year ago
base_conv_layer.cpp	enable dilated deconvolution	a year ago
base_data_layer.cpp	Using default from proto for prefetch	3 months ago
base_data_layer.cu	Switched multi-GPU to NCCL	3 months ago
batch_norm_layer.cpp	Add missing spaces besides equal signs in batch_norm_layer.cpp	4 months ago
batch_norm_layer.cu	dismantle layer headers	a year ago
batch_reindex_layer.cpp	dismantle layer headers	a year ago
batch_reindex_layer.cu	dismantle layer headers	a year ago
bias_layer.cpp	Remove incorrect cast of gemm int arg to Dtype in BiasLayer	a year ago
bias_layer.cu	Separation and generalization of ChannelwiseAffineLayer into BiasLayer	a year ago
bnl_layer.cpp	dismantle layer headers	a year ago
bnl_layer.cu	dismantle layer headers	a year ago
concat_layer.cpp	dismantle layer headers	a year ago
concat_layer.cu	dismantle layer headers	a year ago
contrastive_loss_layer.cpp	dismantle layer headers	a year ago
contrastive_loss_layer.cu	dismantle layer headers	a year ago
conv_layer.cpp	add support for 2D dilated convolution	a year ago
conv_layer.cu	dismantle layer headers	a year ago
crop_layer.cpp	remove redundant operations in Crop layer (#5138)	2 months ago
crop_layer.cu	remove redundant operations in Crop layer (#5138)	2 months ago
cuda_conv_layer.cpp	dismantle layer headers	a year ago
cuda_conv_layer.cu	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago

cuda_lcn_layer.cpp	dismantle layer headers	a year ago
cuda_lcn_layer.cu	dismantle layer headers	a year ago
cuda_lrn_layer.cpp	dismantle layer headers	a year ago
cuda_lrn_layer.cu	dismantle layer headers	a year ago
cuda_pooling_layer.cpp	dismantle layer headers	a year ago
cuda_pooling_layer.cu	dismantle layer headers	a year ago
cuda_relu_layer.cpp	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
cuda_relu_layer.cu	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
cuda_sigmoid_layer.cpp	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
cuda_sigmoid_layer.cu	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
cuda_softmax_layer.cpp	dismantle layer headers	a year ago
cuda_softmax_layer.cu	dismantle layer headers	a year ago
cuda_tanh_layer.cpp	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
cuda_tanh_layer.cu	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
data_layer.cpp	Switched multi-GPU to NCCL	3 months ago
deconv_layer.cpp	enable dilated deconvolution	a year ago
deconv_layer.cu	dismantle layer headers	a year ago
dropout_layer.cpp	supporting N-D Blobs in Dropout layer Reshape	a year ago
dropout_layer.cu	dismantle layer headers	a year ago
dummy_data_layer.cpp	dismantle layer headers	a year ago
eltwise_layer.cpp	dismantle layer headers	a year ago
eltwise_layer.cu	dismantle layer headers	a year ago
elu_layer.cpp	ELU layer with basic tests	a year ago
elu_layer.cu	ELU layer with basic tests	a year ago
embed_layer.cpp	dismantle layer headers	a year ago
embed_layer.cu	dismantle layer headers	a year ago
euclidean_loss_layer.cpp	dismantle layer headers	a year ago
euclidean_loss_layer.cu	dismantle layer headers	a year ago
exp_layer.cpp	Solving issue with exp layer with base e	a year ago
exp_layer.cu	dismantle layer headers	a year ago

[Caffe](#) is licensed under [BSD 2-Clause](#)

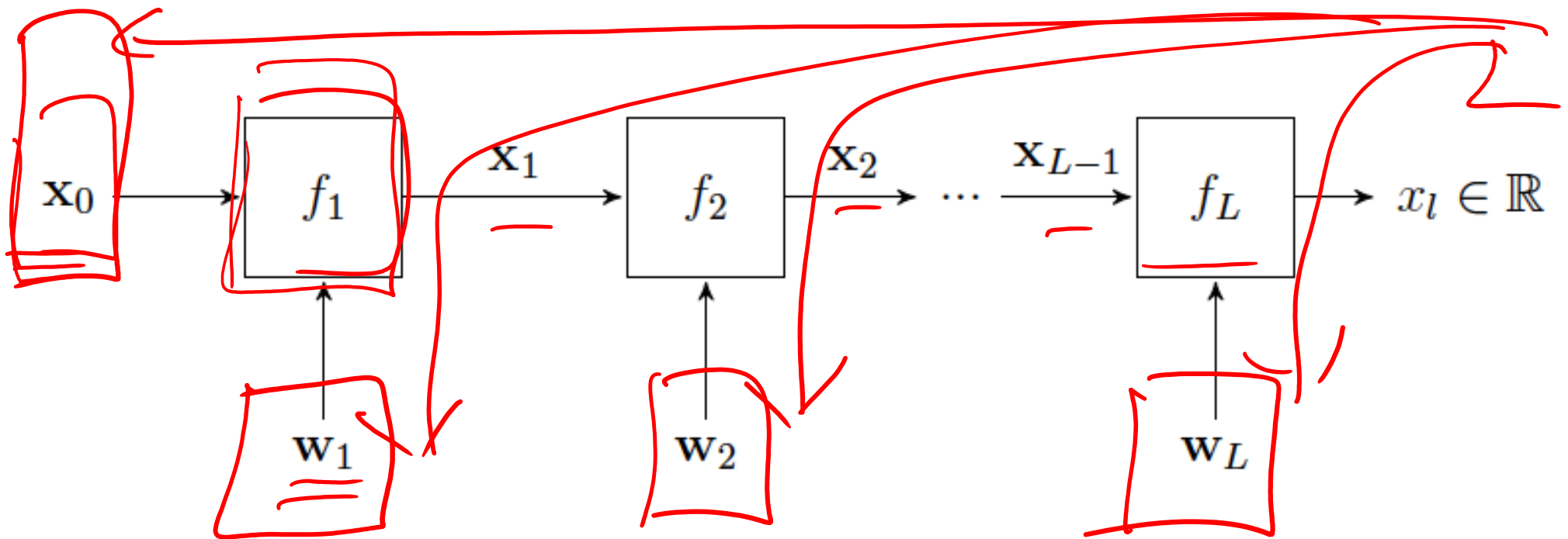
Caffe Sigmoid Layer

```
1 #include <cmath>
2 #include <vector>
3
4 #include "caffe/layers/sigmoid_layer.hpp"
5
6 namespace caffe {
7
8 template <typename Dtype>
9 inline Dtype sigmoid(Dtype x) {
10     return 1. / (1. + exp(-x));
11 }
12
13 template <typename Dtype>
14 void SigmoidLayer<Dtype>::Forward_cpu(const vector<Blob<Dtype>*>& bottom,
15     const vector<Blob<Dtype>*>& top) {
16     const Dtype* bottom_data = bottom[0]->cpu_data();
17     Dtype* top_data = top[0]->mutable_cpu_data();
18     const int count = bottom[0]->count();
19     for (int i = 0; i < count; ++i) {
20         top_data[i] = sigmoid(bottom_data[i]);
21     }
22 }
23
24 template <typename Dtype>
25 void SigmoidLayer<Dtype>::Backward_cpu(const vector<Blob<Dtype>*>& top,
26     const vector<bool>& propagate_down,
27     const vector<Blob<Dtype>*>& bottom) {
28     if (propagate_down[0]) {
29         const Dtype* top_data = top[0]->cpu_data();
30         const Dtype* top_diff = top[0]->cpu_diff();
31         Dtype* bottom_diff = bottom[0]->mutable_cpu_diff();
32         const int count = bottom[0]->count();
33         for (int i = 0; i < count; ++i) {
34             const Dtype sigmoid_x = top_data[i];
35             bottom_diff[i] = top_diff[i] * sigmoid_x * (1. - sigmoid_x);
36         }
37     }
38 }
39
40 #ifdef CPU_ONLY
41 STUB_GPU(SigmoidLayer);
42 #endif
43
44 INSTANTIATE_CLASS(SigmoidLayer);
45
46 } // namespace caffe
```

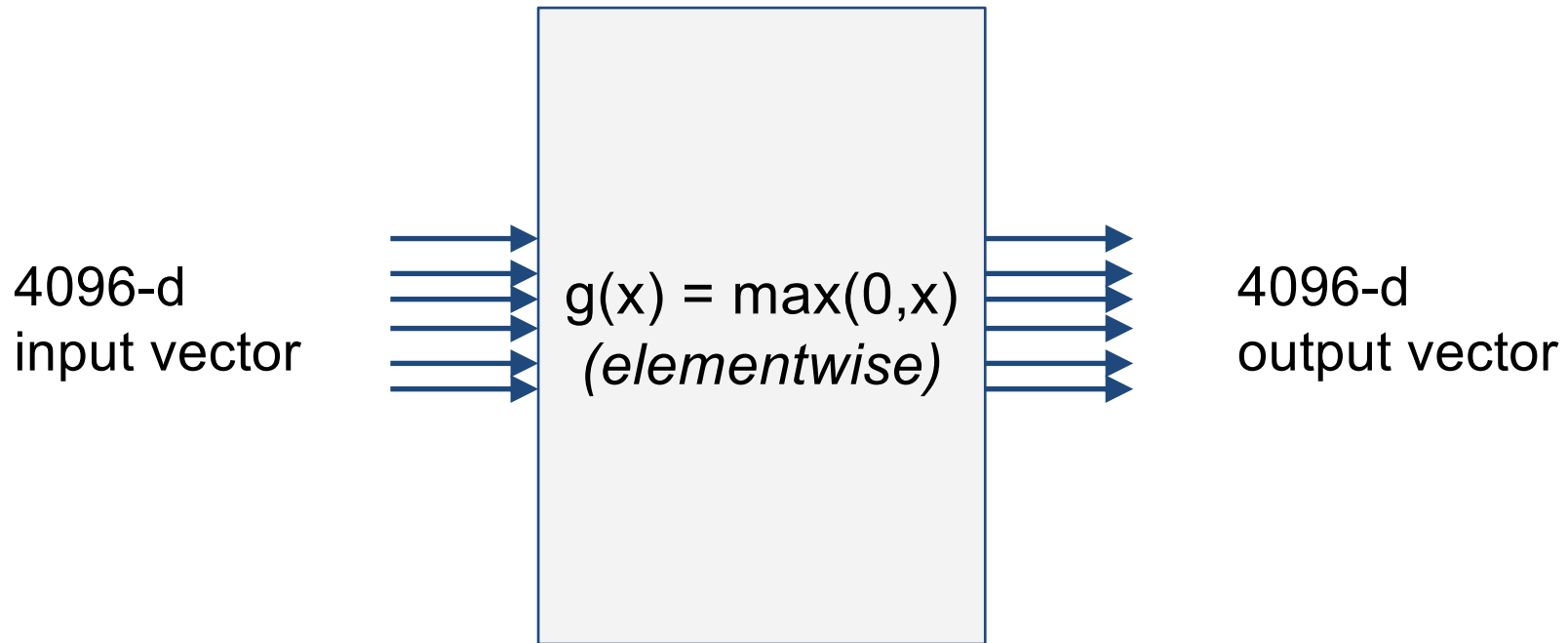
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$(1 - \sigma(x)) \sigma(x) * \text{top_diff} \text{ (chain rule)}$$

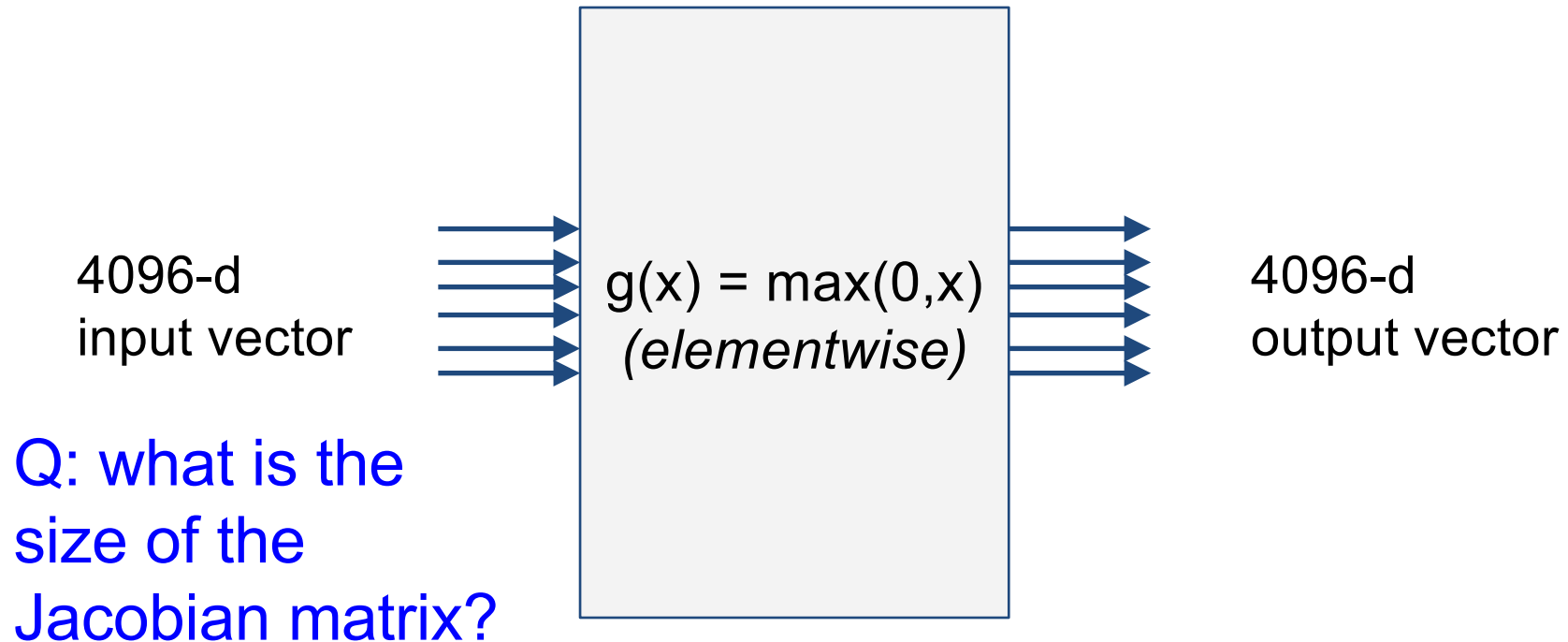
[Caffe](#) is licensed under [BSD 2-Clause](#)



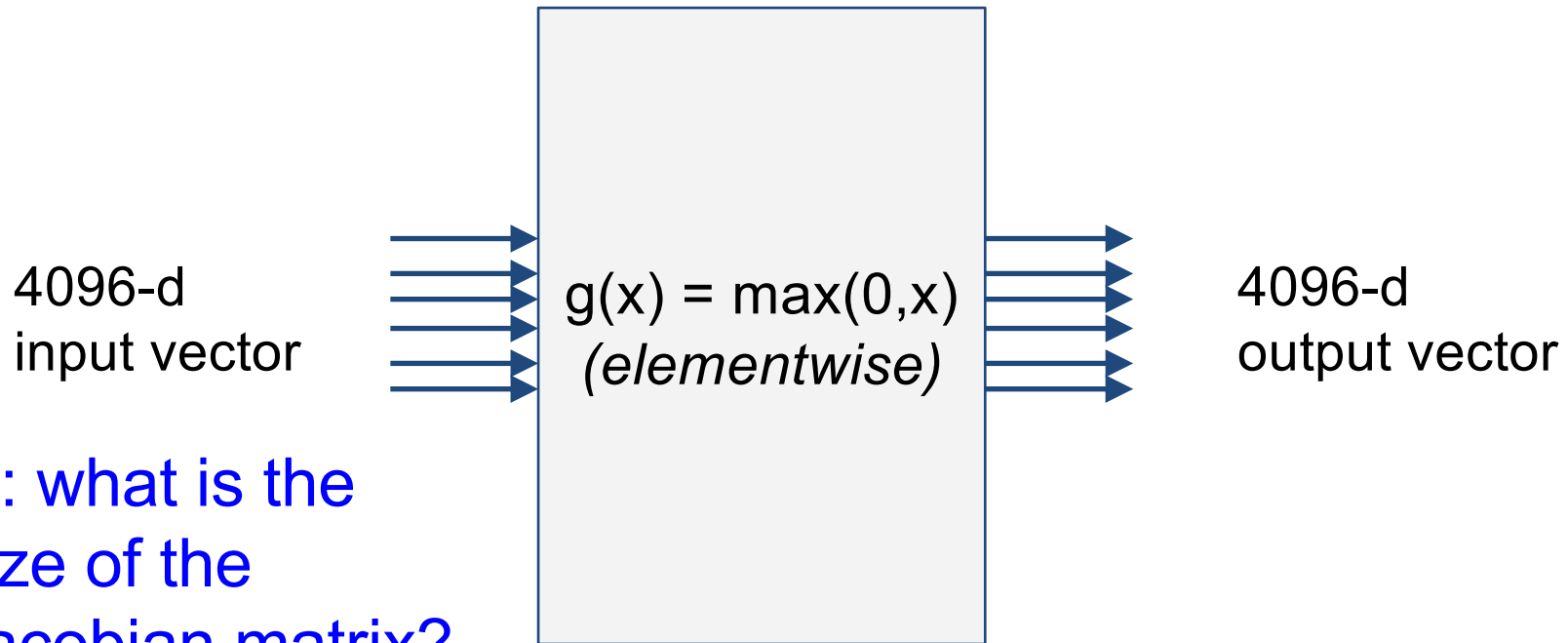
Jacobian of ReLU



Jacobian of ReLU

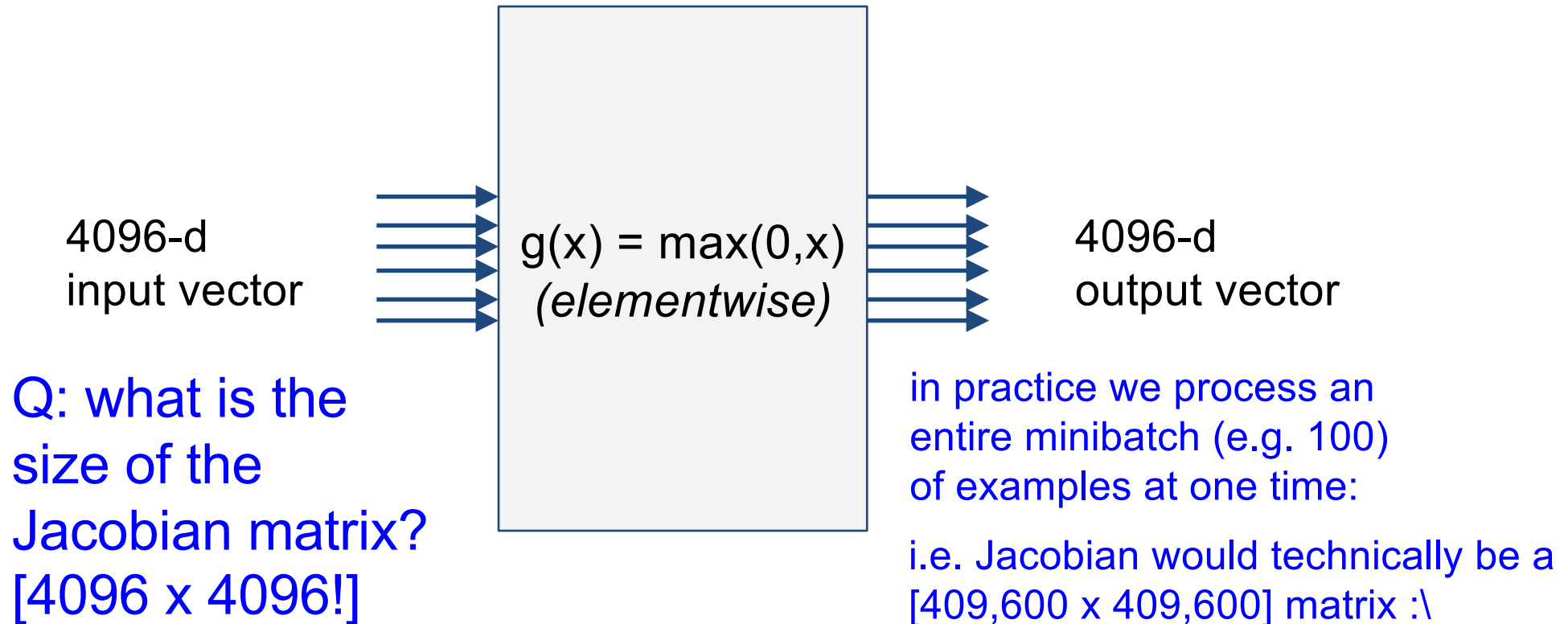


Jacobian of ReLU



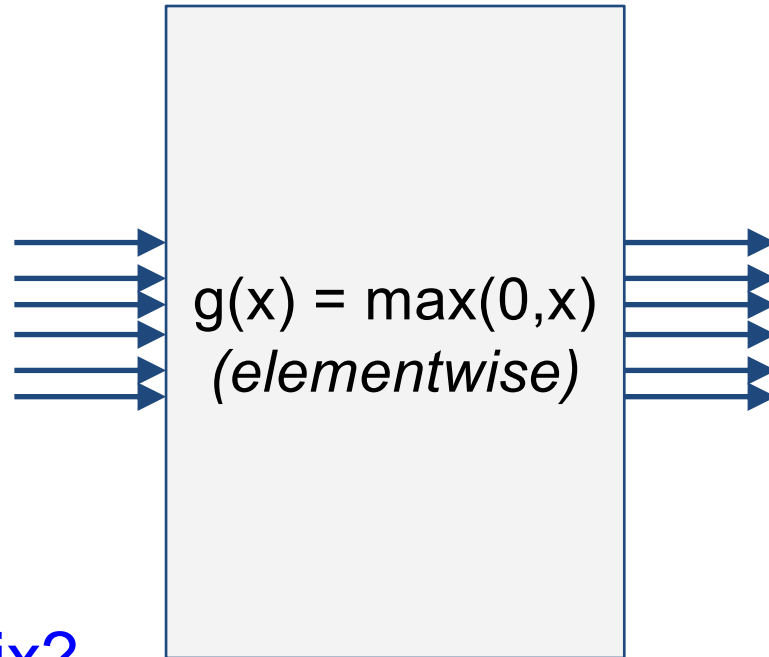
Q: what is the
size of the
Jacobian matrix?
[4096 x 4096!]

Jacobian of ReLU



Jacobian of ReLU

4096-d
input vector



4096-d
output vector

Q: what is the
size of the
Jacobian matrix?
[4096 x 4096!]

Q2: what does it
look like?

Jacobians of FC-Layer

Jacobians of FC-Layer

Jacobians of FC-Layer