

CS 4803 / 7643: Deep Learning

Topics:

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

Dhruv Batra
Georgia Tech

Administrativa

- HW1 Reminder
 - Due: 09/26, 11:55pm
- Fuller schedule + future reading posted (https://www.cc.gatech.edu/classes/AY2020/cs7643_fall/)
 - https://www.cc.gatech.edu/classes/AY2020/cs7643_fall/
 - Caveat: subject to change;
please don't make irreversible decisions based on this.

Recap from last time

Strategy: **Follow the slope**

$$\min_{\vec{w}} L(\vec{w}, D) \leftarrow$$
$$\underline{L(\vec{w})} = \frac{1}{N} \sum_i L_i(w)$$



Gradient Descent

```
# Vanilla Gradient Descent
```

```
while True:
```

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

backprop

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

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

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

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

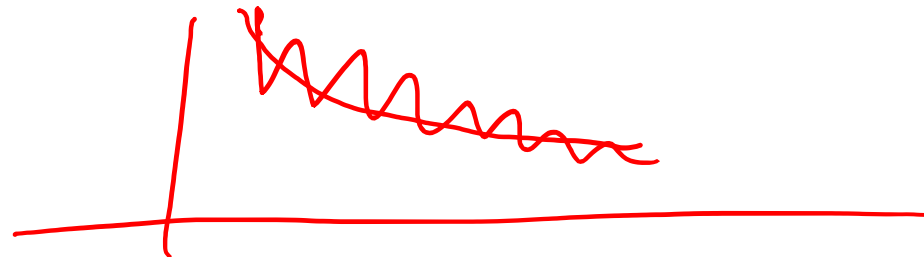
```
# Vanilla Minibatch Gradient Descent
```

```
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

Coding

```
f(x):
  v = x
  for i = 1 to 3
    v = 4*v*(1 - v)
  return v
```

or, in closed-form,

```
f(x):
  return 64*x*(1-x)*((1-2*x)^2)
  *(1-8*x+8*x*x)^2
```

Symbolic
Differentiation
of the Closed-form

```
f'(x):
  return 128*x*(1-x)*(-8+16*x)
  *((1-2*x)^2)*(1-8*x+8*x*x)
  + 64*(1-x)*((1-2*x)^2)*((1-
  -8*x+8*x*x)^2) - (64*x*(1-
  2*x)^2)*(1-8*x+8*x*x)^2 -
  256*x*(1-x)*(1-2*x)*(1-8*x
  + 8*x*x)^2
```

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

Automatic
Differentiation

Numerical
Differentiation

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

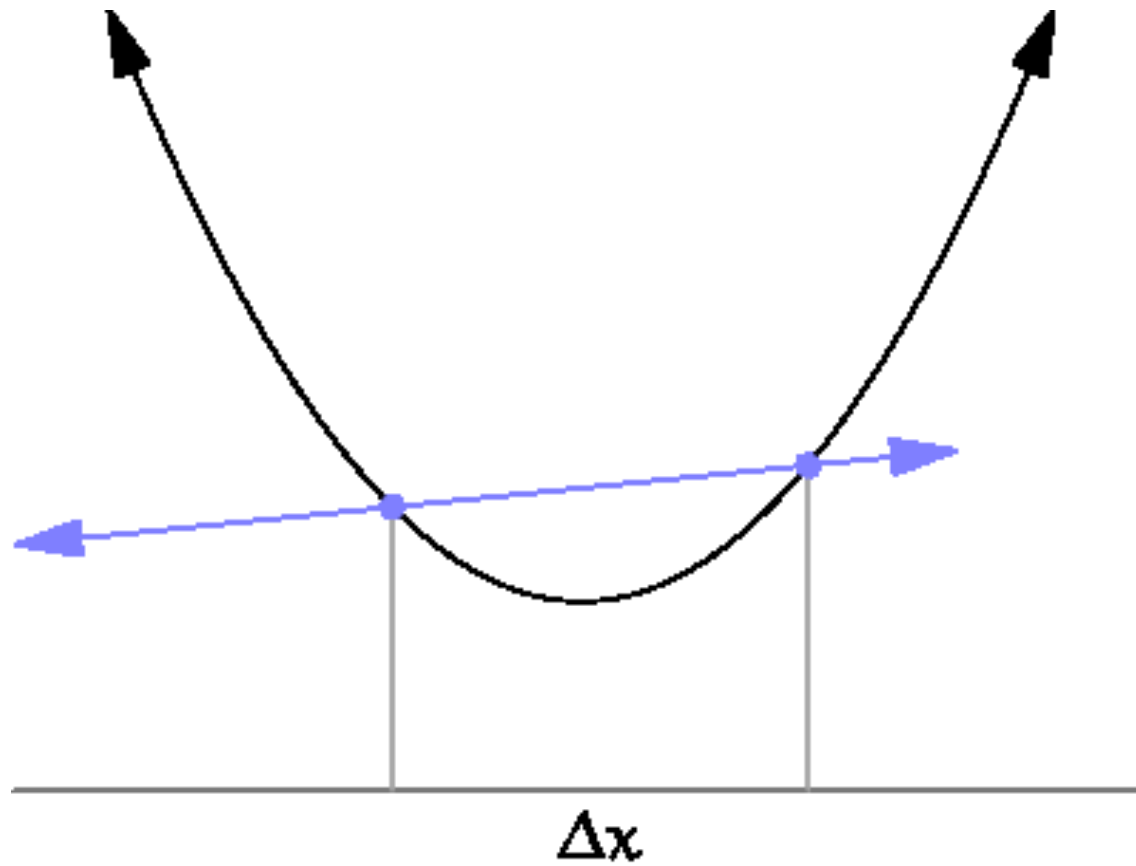
$f'(x_0) = f'(x_0)$
Exact

```
f'(x):
  h = 0.000001
  return (f(x+h) - f(x)) / h
```

$f'(x_0) \approx f'(x_0)$
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

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,
?,
?,

$$\frac{(1.25322 - 1.25347)}{0.0001} = -2.5$$

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

?,
?,...]

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,
0.6,
?,
?,

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

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

?,...]

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 :(

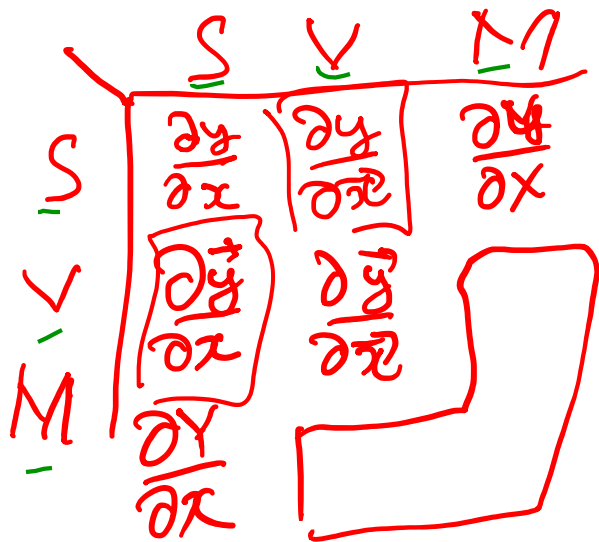
In practice: Derive analytic gradient, check your implementation with numerical gradient.

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



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

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

=

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

num = dim 1
row

den = dim 2

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

Vector/Matrix Derivatives Notation

$$\underbrace{\frac{\partial \vec{y}}{\partial \vec{x}}}_{c \times d} = \left[\begin{array}{c} \frac{\partial y_1}{\partial x_1} \quad \dots \quad \frac{\partial y_1}{\partial x_j} \\ \vdots \\ \frac{\partial y_i}{\partial x_1} \quad \dots \quad \frac{\partial y_i}{\partial x_j} \\ \vdots \\ \frac{\partial y_d}{\partial x_1} \quad \dots \quad \frac{\partial y_d}{\partial x_j} \end{array} \right]$$

The diagram illustrates the notation for vector/matrix derivatives. On the left, the derivative of a vector \vec{y} with respect to a vector \vec{x} is shown as $\frac{\partial \vec{y}}{\partial \vec{x}}$, which is a $c \times d$ matrix. The matrix is represented as a large bracketed structure containing rows of partial derivatives. The first row is $\frac{\partial y_1}{\partial x_1} \dots \frac{\partial y_1}{\partial x_j}$, and the i -th row is $\frac{\partial y_i}{\partial x_1} \dots \frac{\partial y_i}{\partial x_j}$. The d -th row is $\frac{\partial y_d}{\partial x_1} \dots \frac{\partial y_d}{\partial x_j}$. The matrix is labeled $c \times d$ at the bottom right. A green arrow points down from the matrix to the label $c \times d$.

Vector Derivative Example

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

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

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

$$\frac{\partial y}{\partial \vec{x}} = \left[\frac{\partial y}{\partial x_1} \quad \dots \quad \frac{\partial y}{\partial x_d} \right]$$

$$\frac{\partial (\sum w_i x_i)}{\partial x_1}$$

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

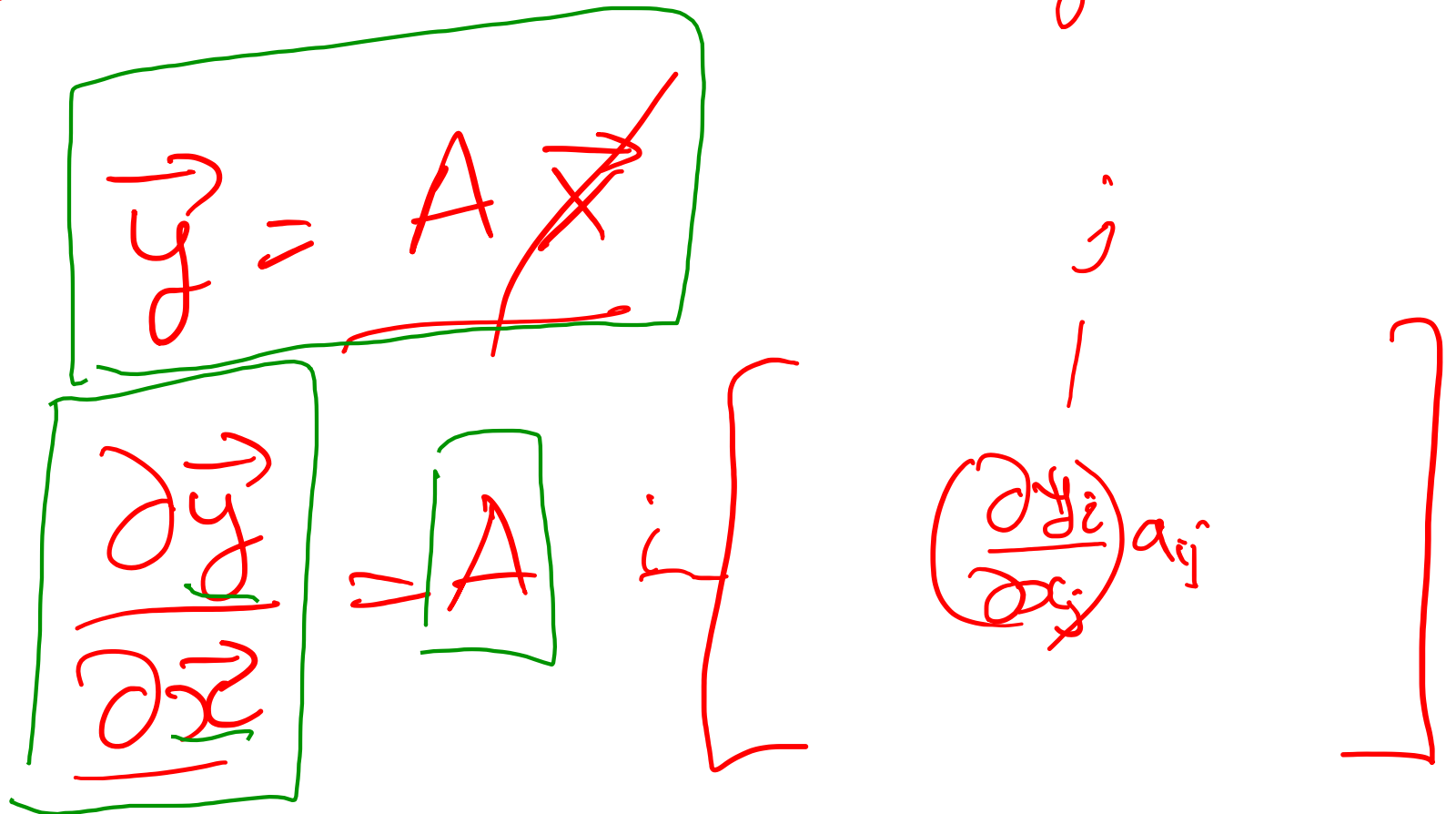
$$\left[w_1 \quad \dots \quad w_d \right]$$

\vec{w}^T

Vector Derivative Example

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

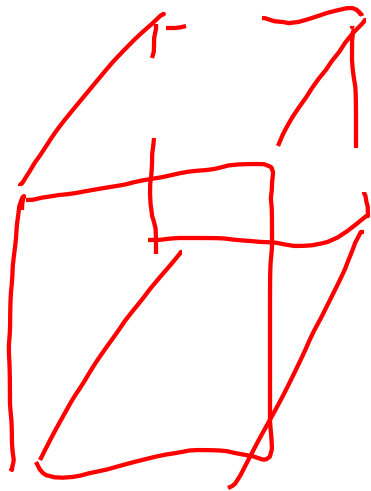
$y_i = \sum_j a_{ij} x_j$



Extension to Tensors

$$X \in \mathbb{R}^{d_1 \dots d_m}$$

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



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

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

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

$$\frac{\partial y\text{-vec}}{\partial x\text{-vec}} = \left[\right]$$

Plan for Today

- (Finish) Analytical Gradients
- Automatic Differentiation
 - Computational Graphs
 - Forward mode vs Reverse mode AD
 - Patterns in backprop

Chain Rule: Composite Functions

$$\left[\underline{L}(x) = \underline{f}(g(x)) = (f \circ g)(x) \right]$$

$$\begin{aligned} \underbrace{f(x)}_{L(w)} &= \underline{g}_l(\underline{g}_{l-1} \dots \underline{g}_1(x)) \\ &= (\underline{g}_l \circ \underline{g}_{l-1} \dots \circ \underline{g}_1)(x) \\ \frac{dL}{dx} & \end{aligned}$$

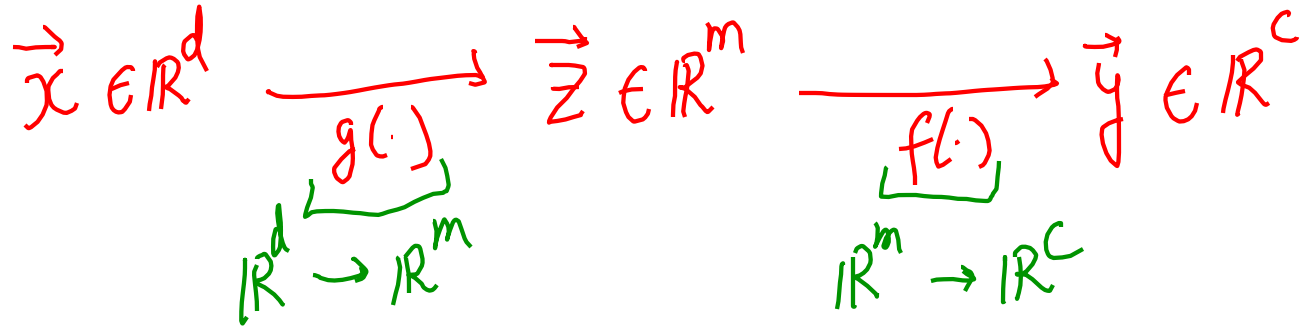
Chain Rule: Scalar Case

$$x \xrightarrow{g(\cdot)} z \xrightarrow{f(\cdot)} y \quad x, y, z \in \mathbb{R}^1$$
$$= f(\underbrace{g(x)}_z)$$

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

scalar prod.

Chain Rule: Vector Case



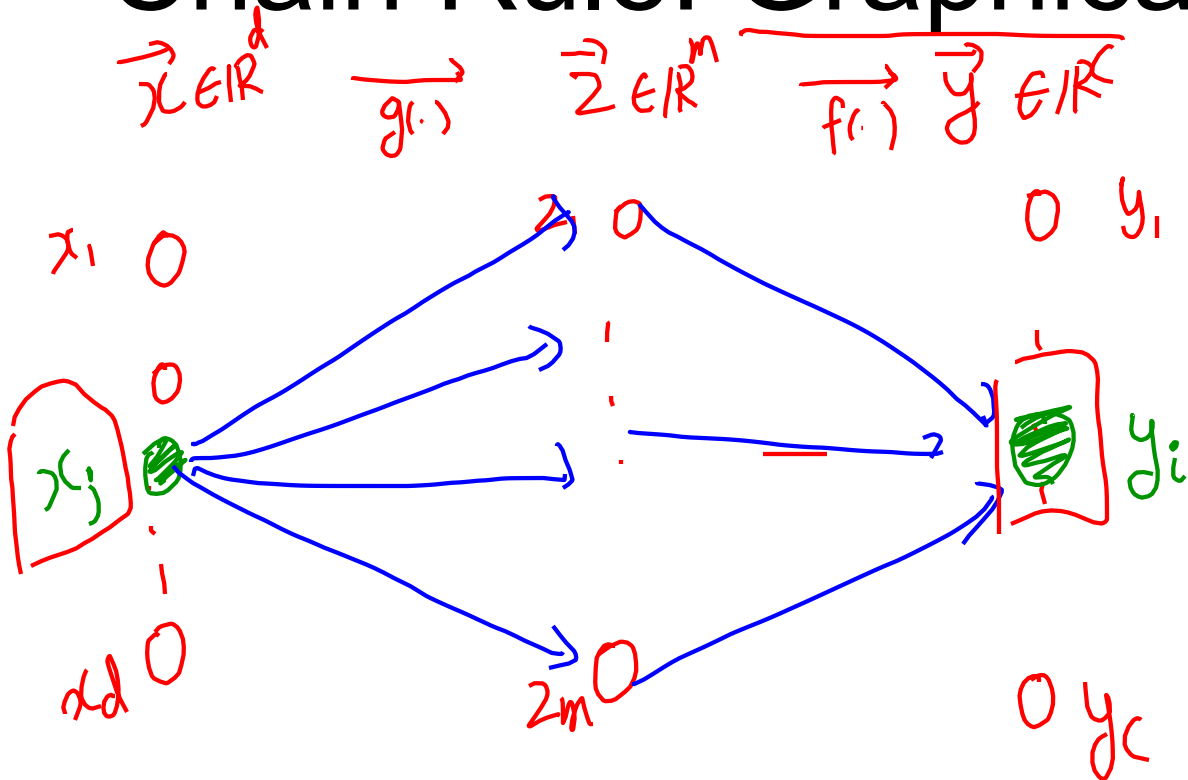
$$\underbrace{\begin{bmatrix} \frac{\partial \vec{y}}{\partial \vec{x}} \end{bmatrix}}_{J_{f \circ g}} = \begin{bmatrix} \frac{\partial \vec{y}}{\partial \vec{z}} \\ \cdot \\ \frac{\partial \vec{z}}{\partial \vec{x}} \end{bmatrix} \begin{bmatrix} \frac{\partial \vec{z}}{\partial \vec{x}} \\ \cdot \\ \frac{\partial \vec{z}}{\partial \vec{x}} \end{bmatrix}$$

J_f Matrix Mult J_g

Chain Rule: Jacobian view

$$\begin{aligned}
 & \left[\begin{array}{c} \frac{\partial y}{\partial \vec{x}} \\ \vdots \\ \frac{\partial y_i}{\partial x_j} \\ \vdots \end{array} \right]_{c \times d} = \left[\begin{array}{c} \frac{\partial \vec{y}}{\partial \vec{z}} \\ \vdots \\ \frac{\partial y_i}{\partial z_k} \\ \vdots \end{array} \right]_{c \times m} \left[\begin{array}{c} \frac{\partial z}{\partial \vec{x}} \\ \vdots \\ \frac{\partial z_k}{\partial x_j} \\ \vdots \end{array} \right]_{m \times d} \\
 & \left[\frac{\partial y_i}{\partial x_j} \right] = \sum_k \left[\frac{\partial y_i}{\partial z_k} \right] \left[\frac{\partial z_k}{\partial x_j} \right]
 \end{aligned}$$

Chain Rule: Graphical view



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

k is on path

Linear Classifier: Logistic Regression

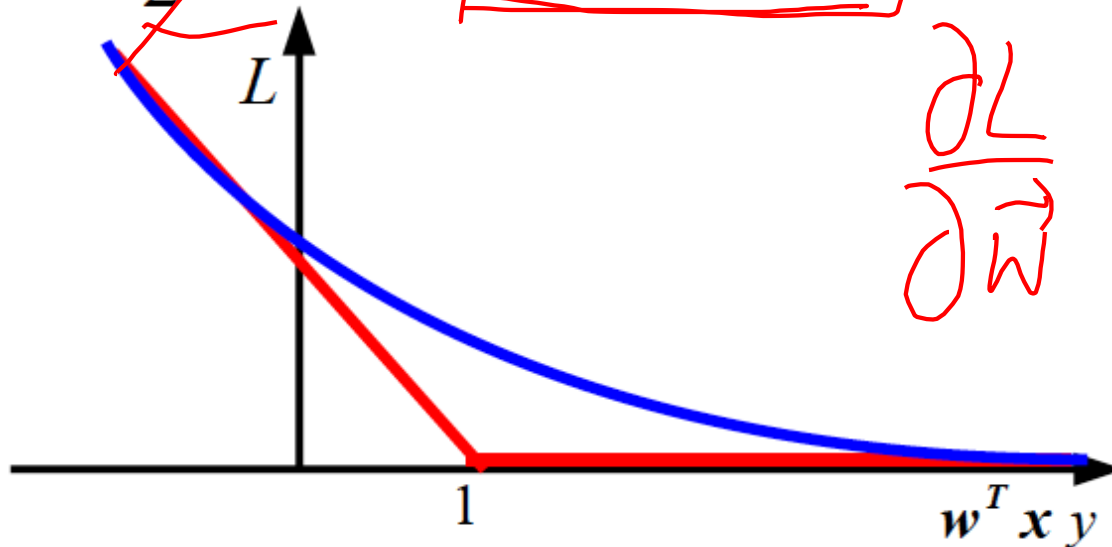
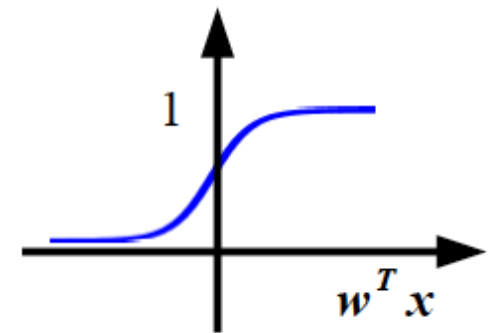
Input: $x \in R^D$

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

Parameters: $w \in R^D$

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

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



Log Loss

Logistic Regression Derivatives

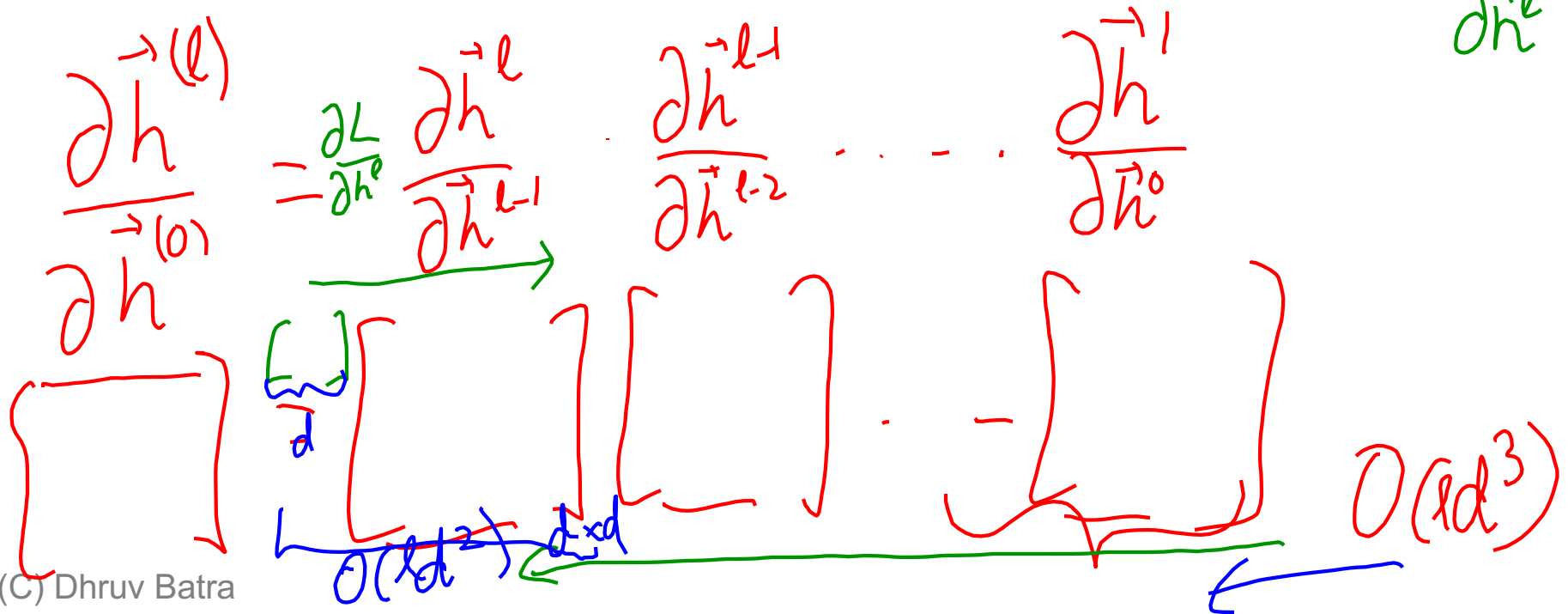
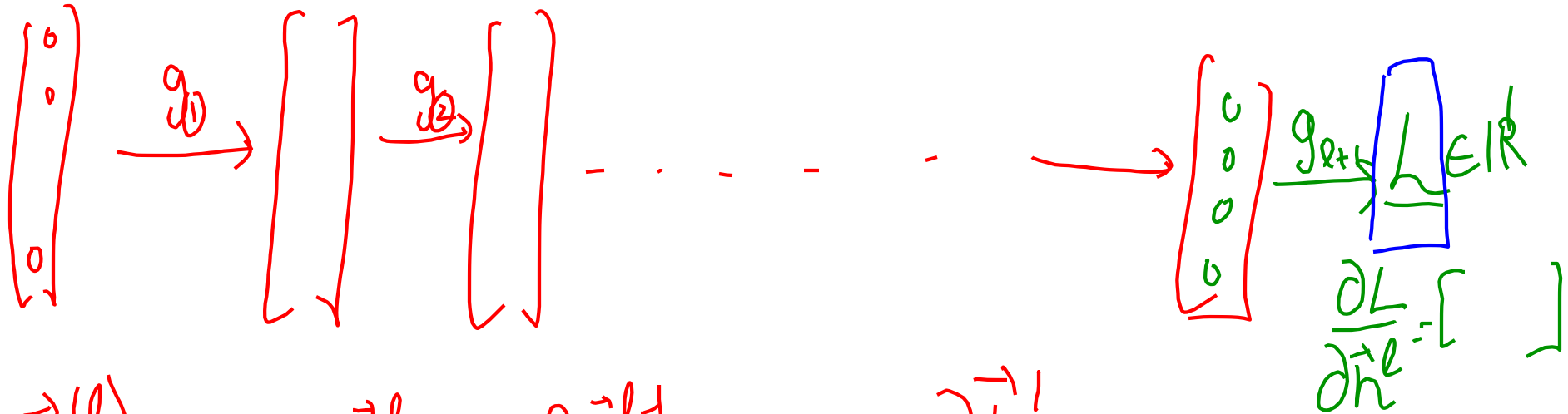
$$L_i = -\log\left(\frac{1}{1+e^{-w^T x}}\right) = \log\left(\underbrace{1+e^{-w^T x}}_p\right)$$

$$\frac{\partial L_i}{\partial \vec{w}} = \frac{\partial L_i}{\partial p} \cdot \frac{\partial p}{\partial u} \cdot \frac{\partial u}{\partial \vec{w}}$$

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

Logistic Regression Derivatives

Chain Rule: Cascaded



Chain Rule: How should we multiply?

Convolutional network (AlexNet)

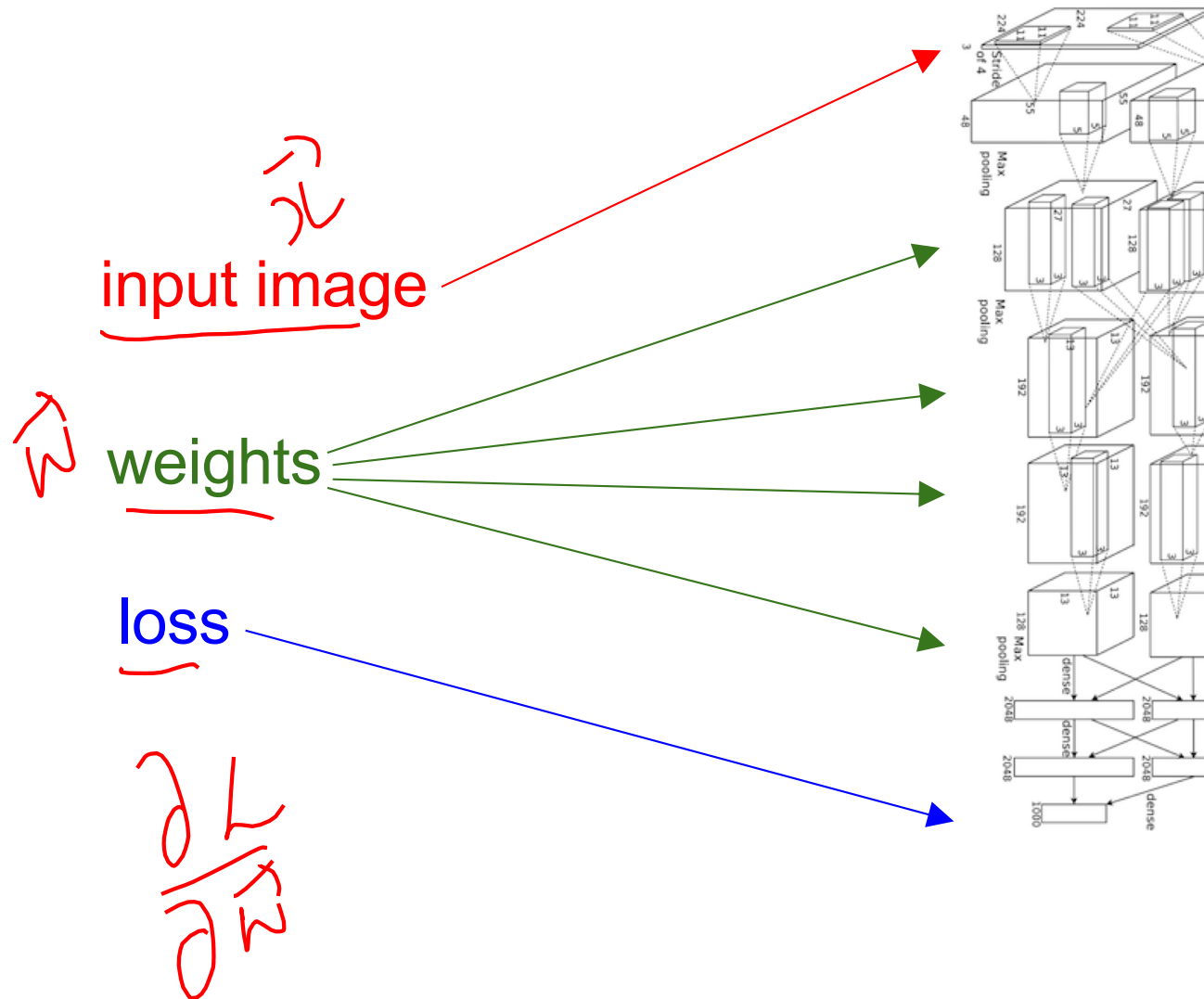


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Neural Turing Machine

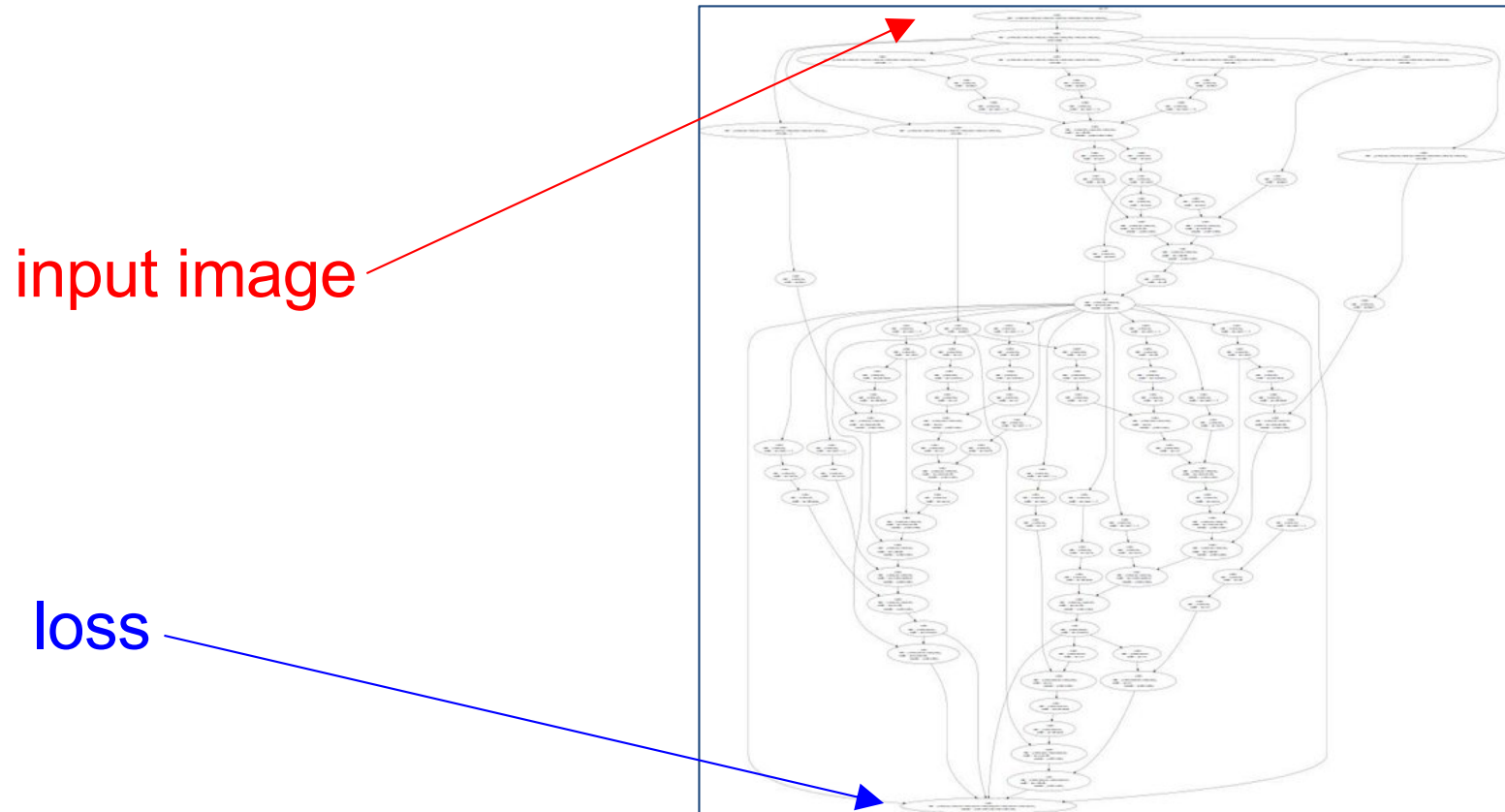
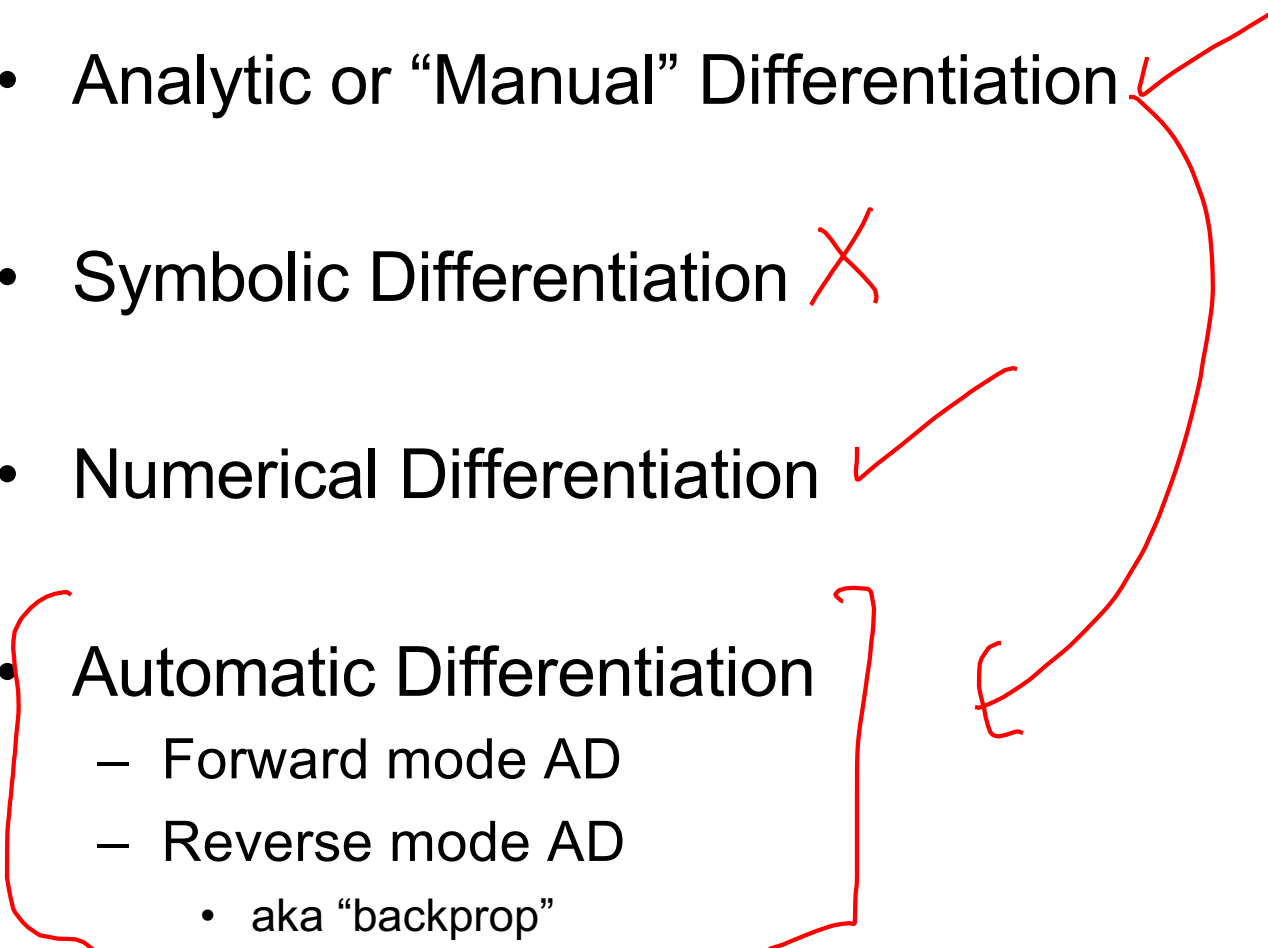


Figure reproduced with permission from a [Twitter post](#) by Andrej Karpathy.

How do we compute gradients?

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

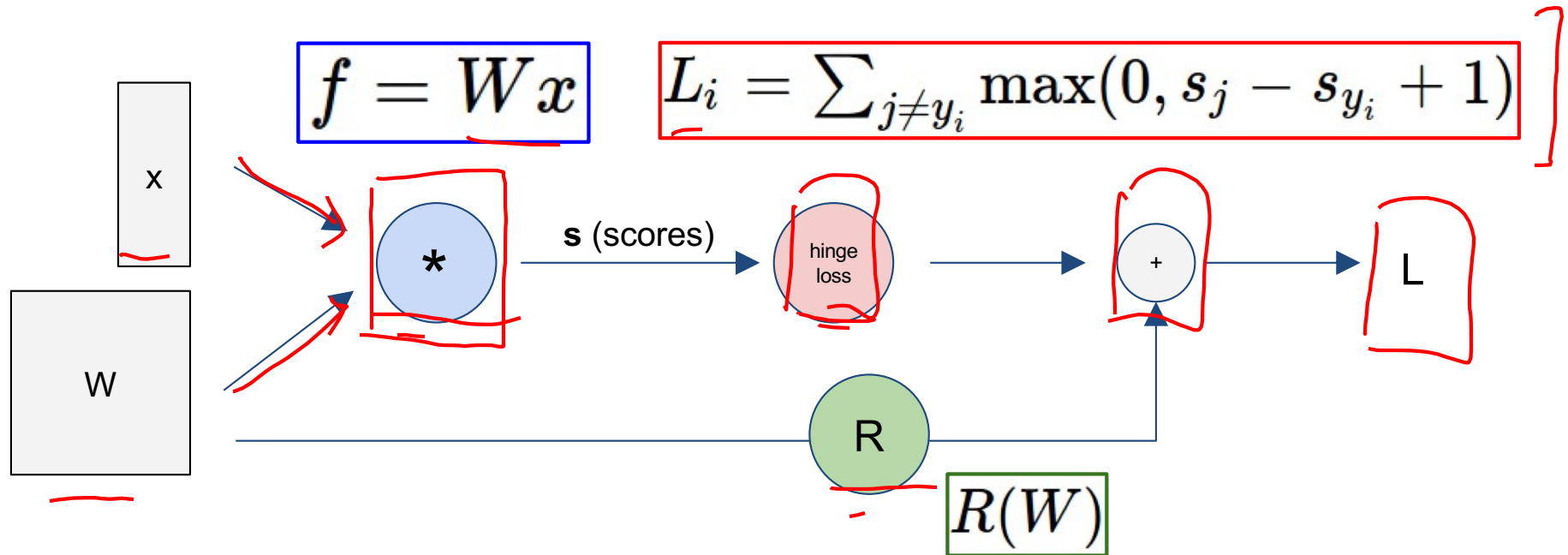
Plan for Today

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

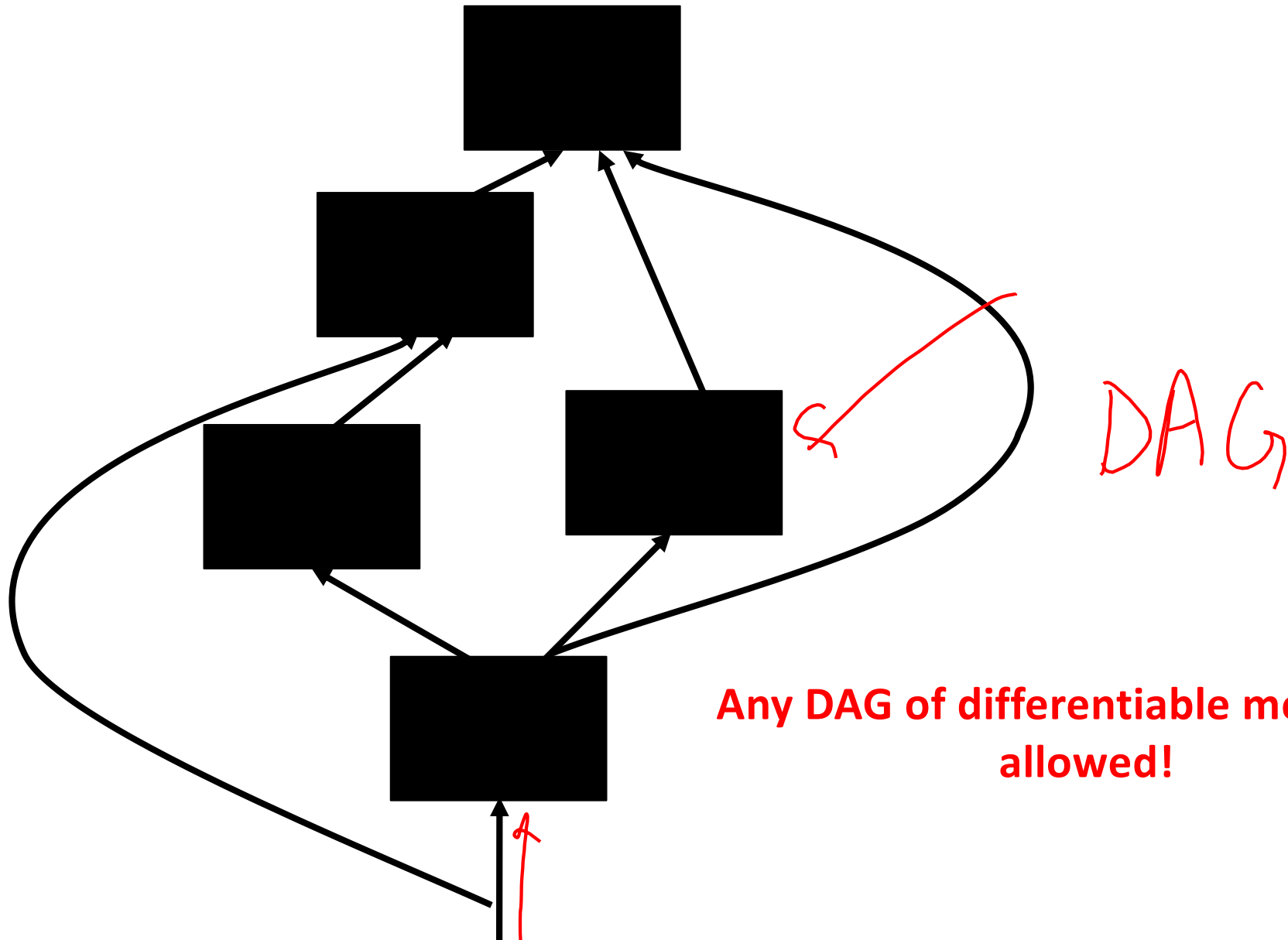
Deep Learning = Differentiable Programming

- 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



Computational Graph



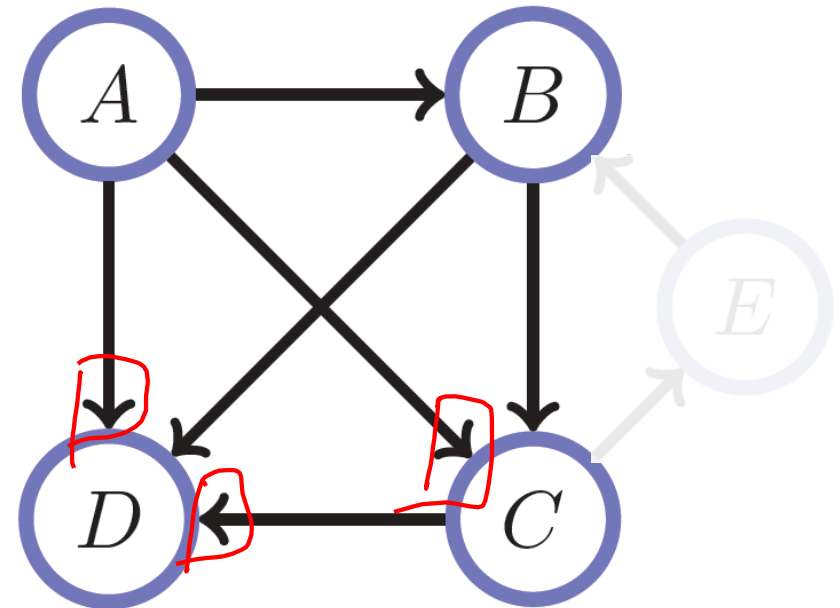
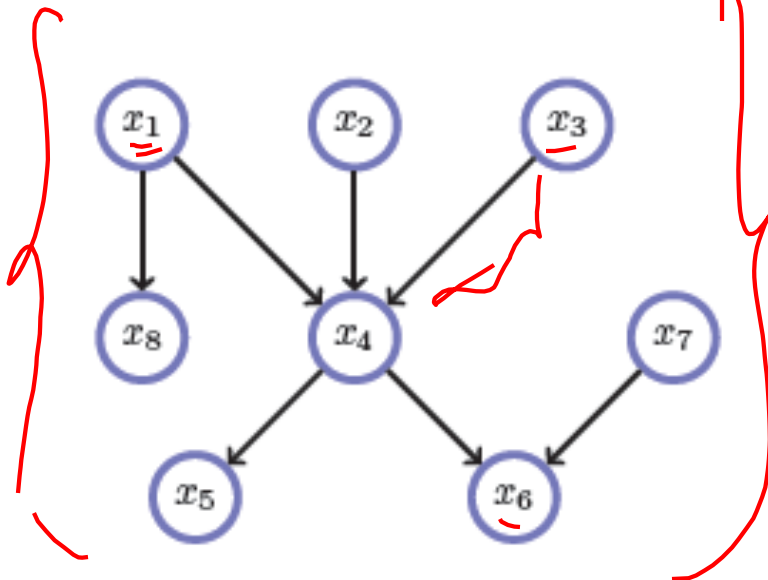
Directed Acyclic Graphs (DAGs)

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

$$G = (V, E)$$

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

Vertex
Vertices



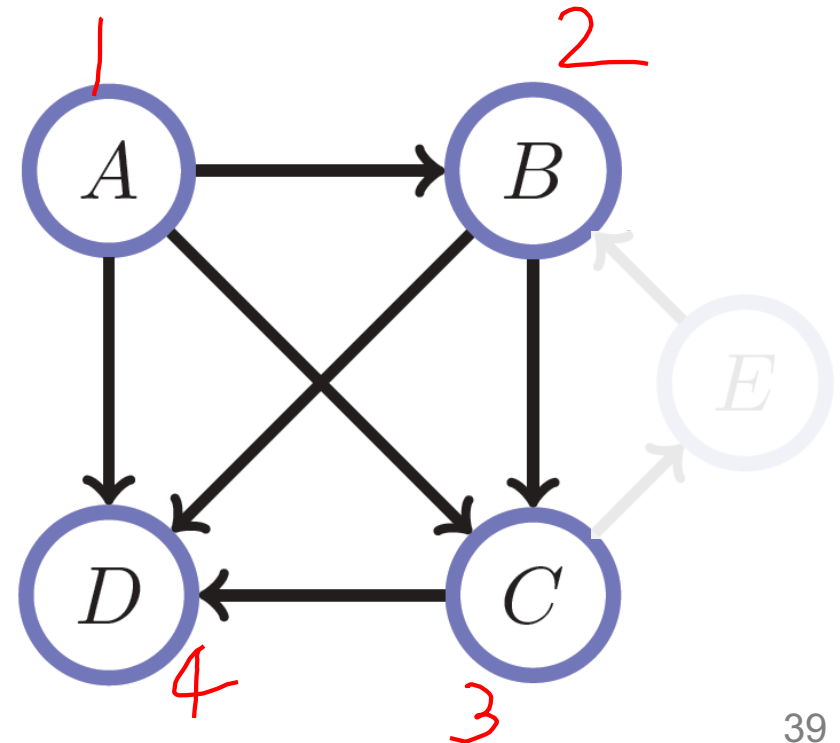
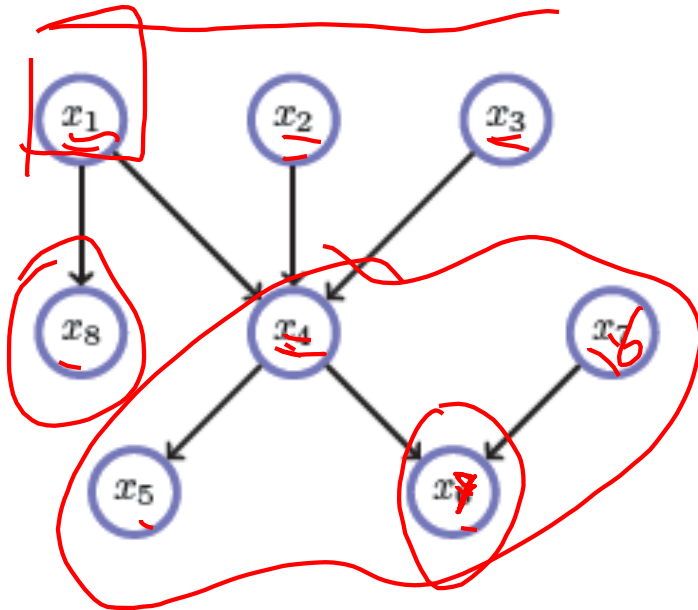
Directed Acyclic Graphs (DAGs)

- Concept
 - Topological Ordering

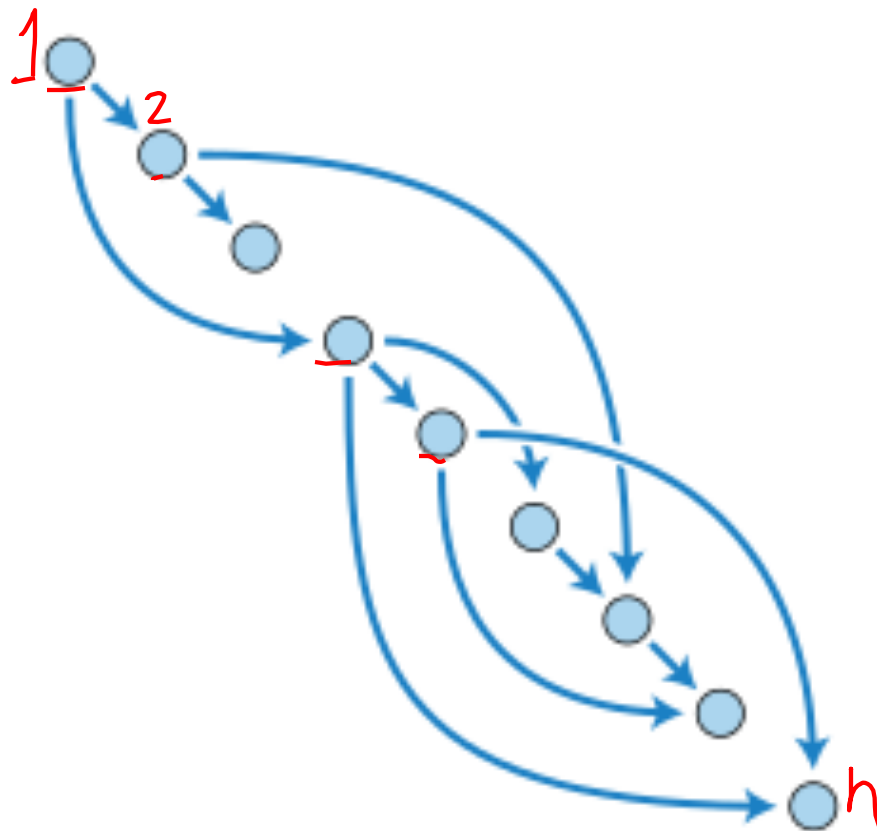
\exists bijection $\sigma: V \rightarrow \{1, \dots, n\}$

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

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



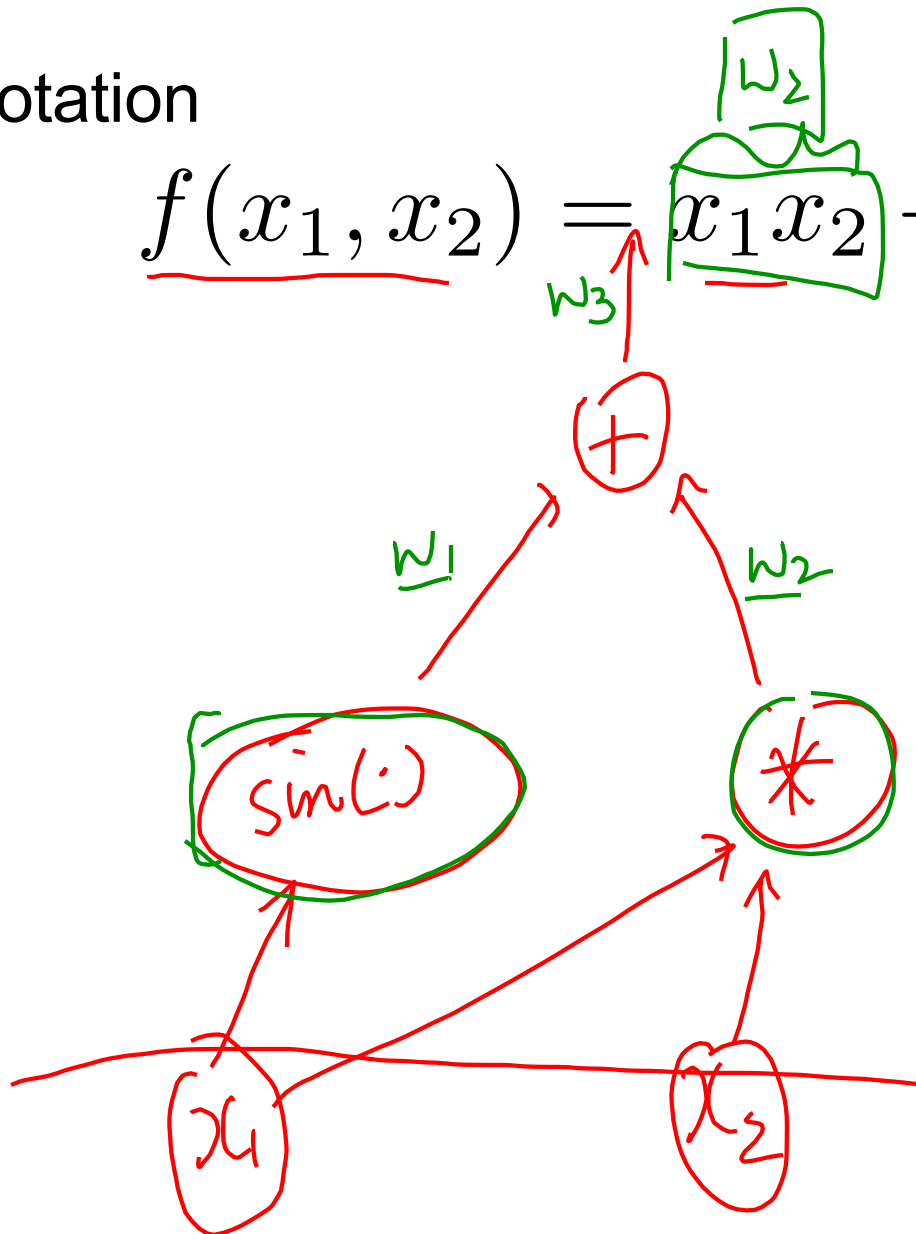
Directed Acyclic Graphs (DAGs)



Computational Graphs

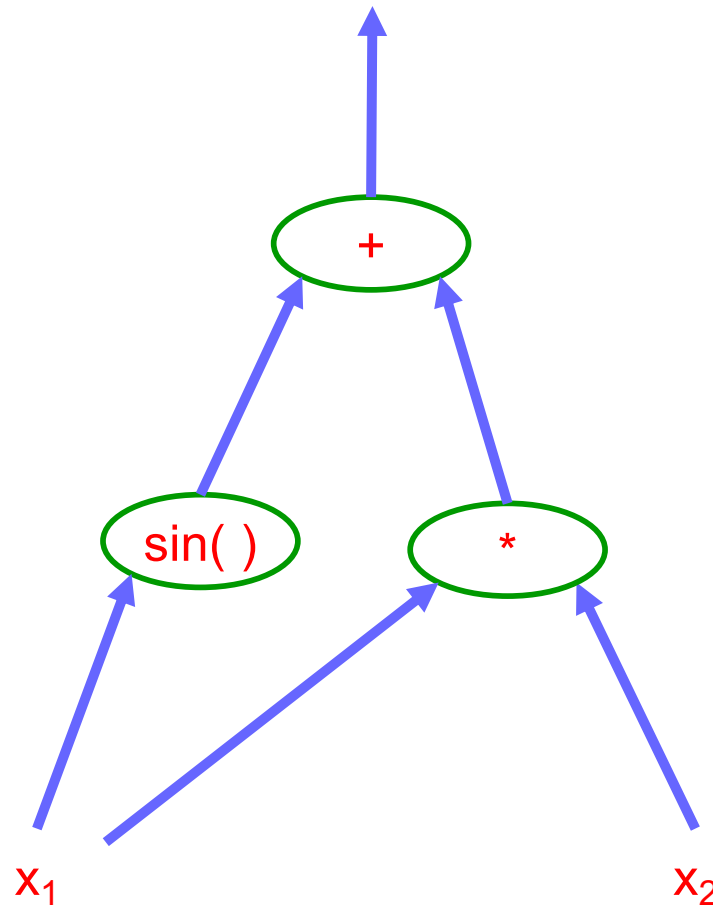
- Notation

$$\underline{f(x_1, x_2)} = \underline{x_1 x_2} + \underline{\sin(x_1)}$$



Example

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

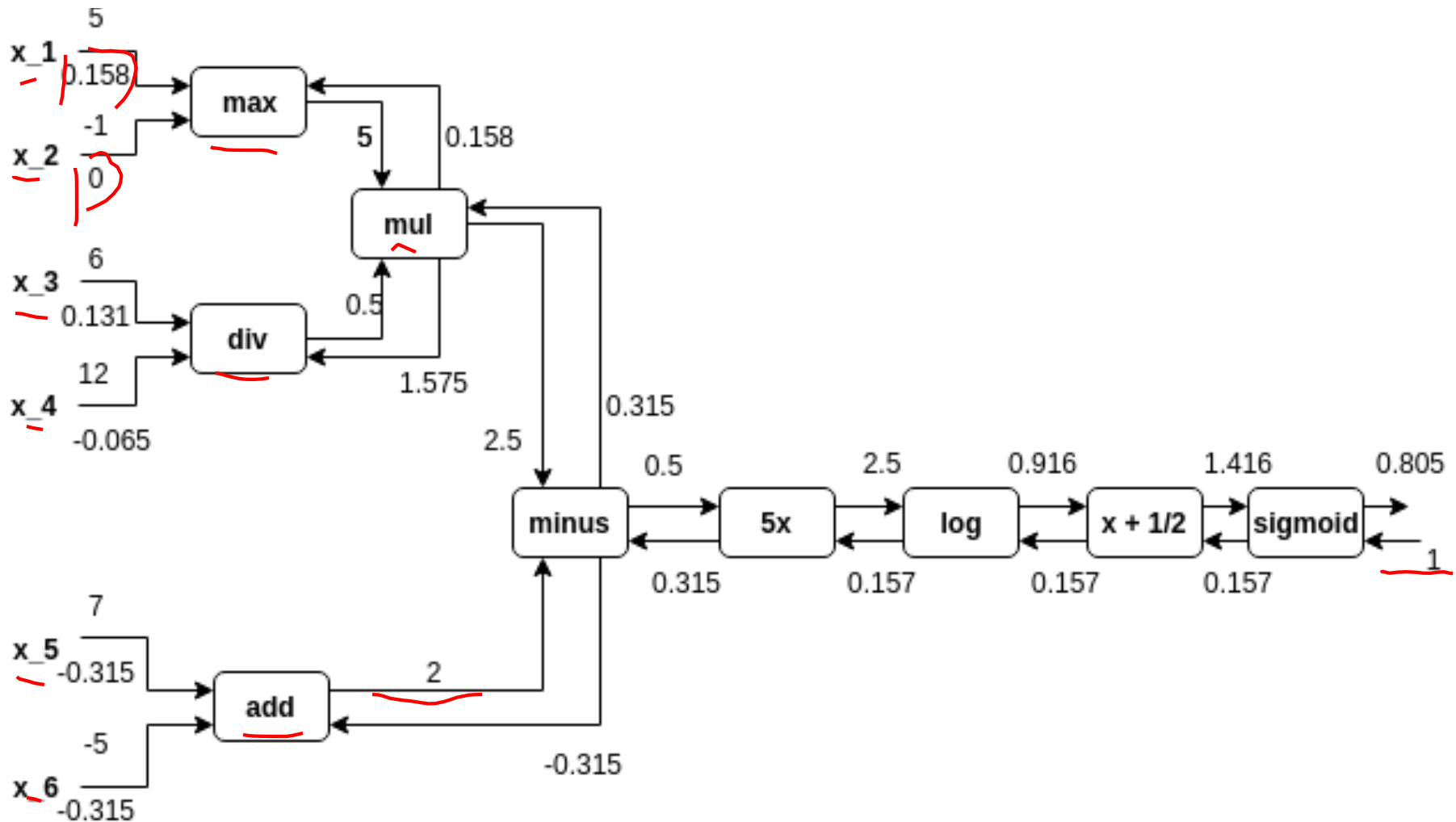


HW0

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

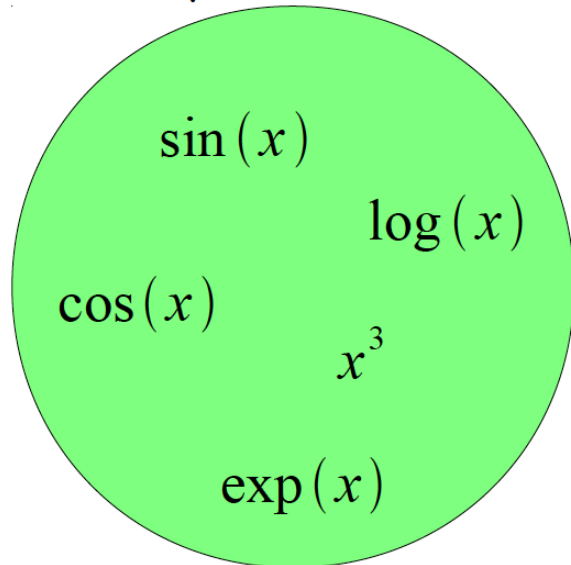
HW0 Submission by Samyak Datta

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



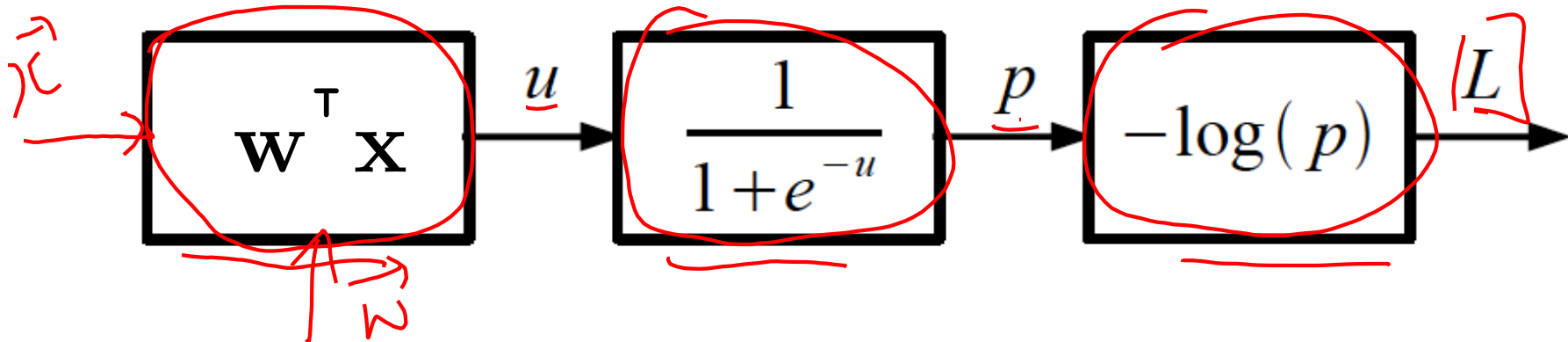
Logistic Regression as a Cascade

Given a library of simple functions



Compose into a
complicate function

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



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

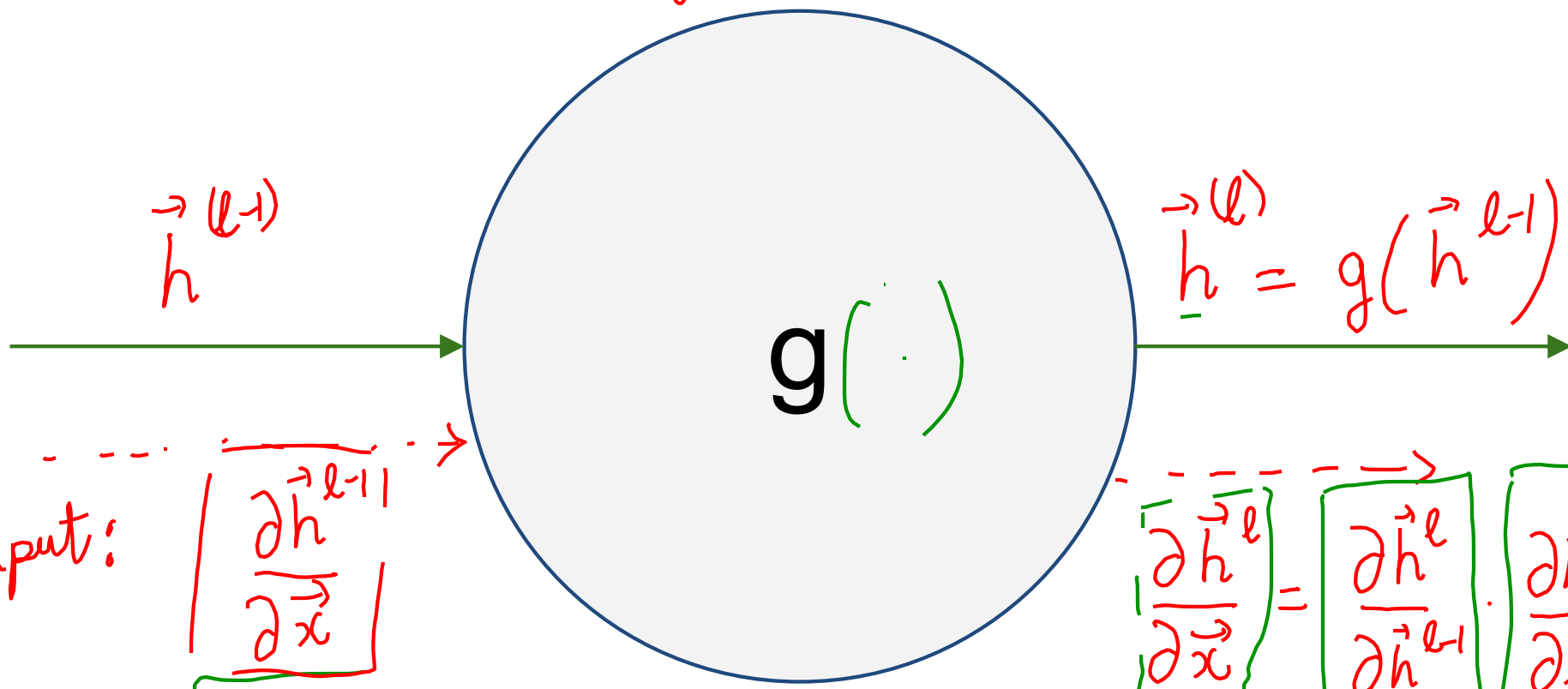
Forward mode vs Reverse Mode

- Key Computations

Forward mode AD

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

layer l



Input: $\begin{bmatrix} \frac{\partial \vec{h}^{(l-1)}}{\partial \vec{x}} \end{bmatrix}$

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

Jacobian Input of g

$\vec{z} \rightarrow \vec{x} \rightarrow \vec{z}$

Reverse mode AD

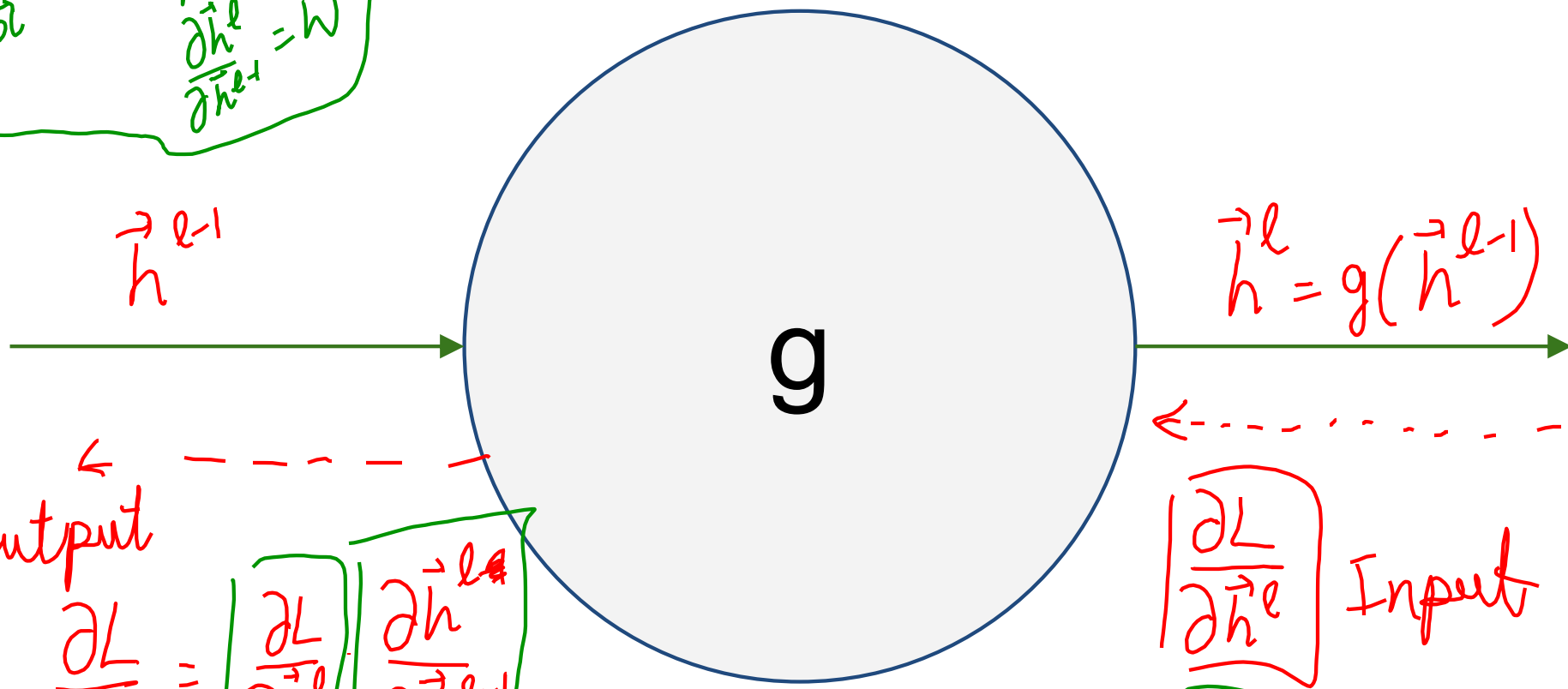
Goal: $\frac{\partial L}{\partial \vec{x}}$

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

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

$$\vec{h}^e = W\vec{h}^{e-1}$$

$$\frac{\partial \vec{h}^e}{\partial \vec{h}^{e-1}} = W$$



Output

$$\frac{\partial L}{\partial \vec{h}^{e-1}} = \left[\frac{\partial L}{\partial \vec{h}^e} \right] \left[\frac{\partial \vec{h}^e}{\partial \vec{h}^{e-1}} \right]$$

Input

Jacobian of g

$$\left[\frac{\partial L}{\partial \vec{h}^e} \right] \text{ Input}$$

Example: Forward mode AD

Goal:

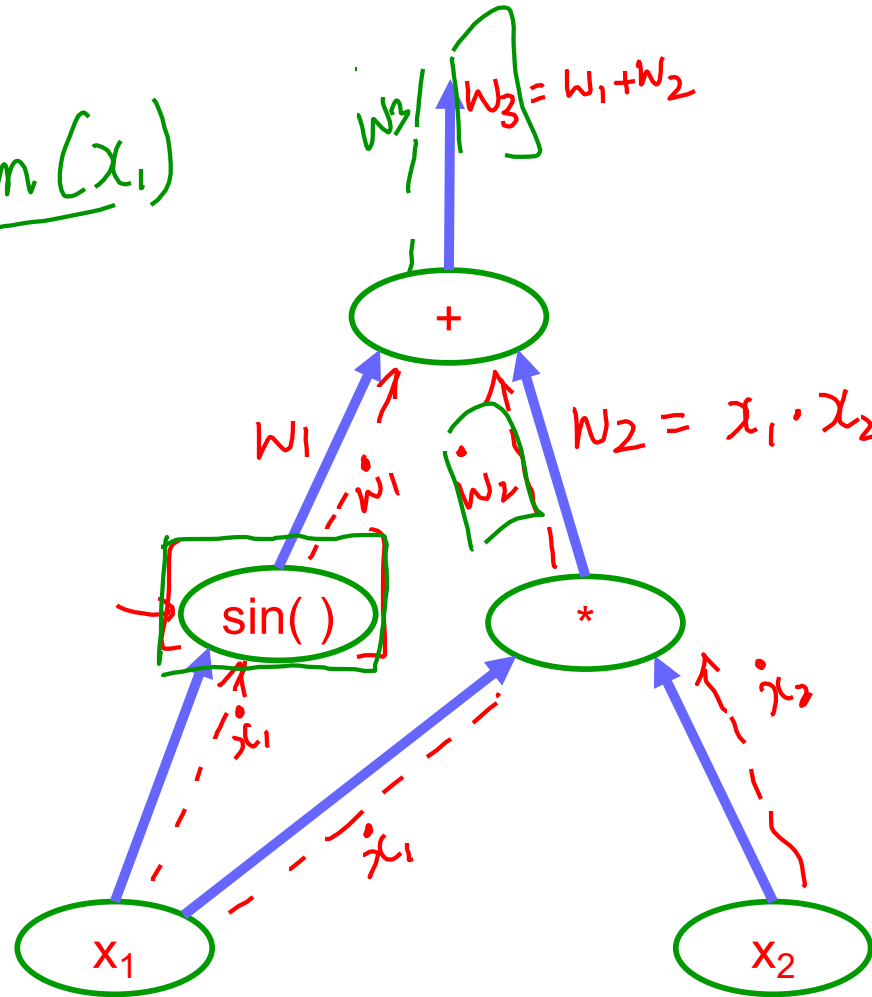
$$\frac{df}{dx} = \left[\frac{df}{dx_1} \quad \frac{df}{dx_2} \right]$$

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

$\frac{\partial x_1}{\partial a} = \dot{x}_1$
 $\frac{\partial x_2}{\partial a}$
 $\frac{\partial w_1}{\partial a}$
 $\frac{\partial w_2}{\partial a}$
 $\frac{\partial w_3}{\partial a}$

$w_1 = \sin(x_1)$

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



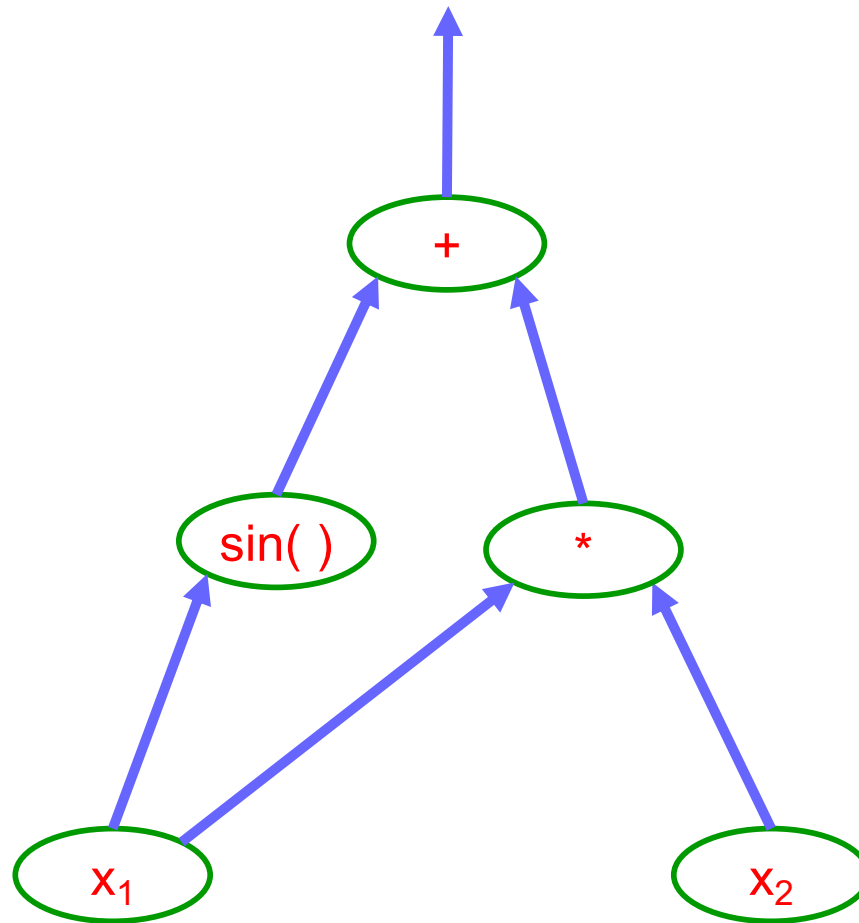
$$\frac{\partial w_2}{\partial a} = x_1 \frac{\partial x_2}{\partial a} + x_2 \frac{\partial x_1}{\partial a}$$

$$= x_1 \dot{x}_2 + x_2 \dot{x}_1$$

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

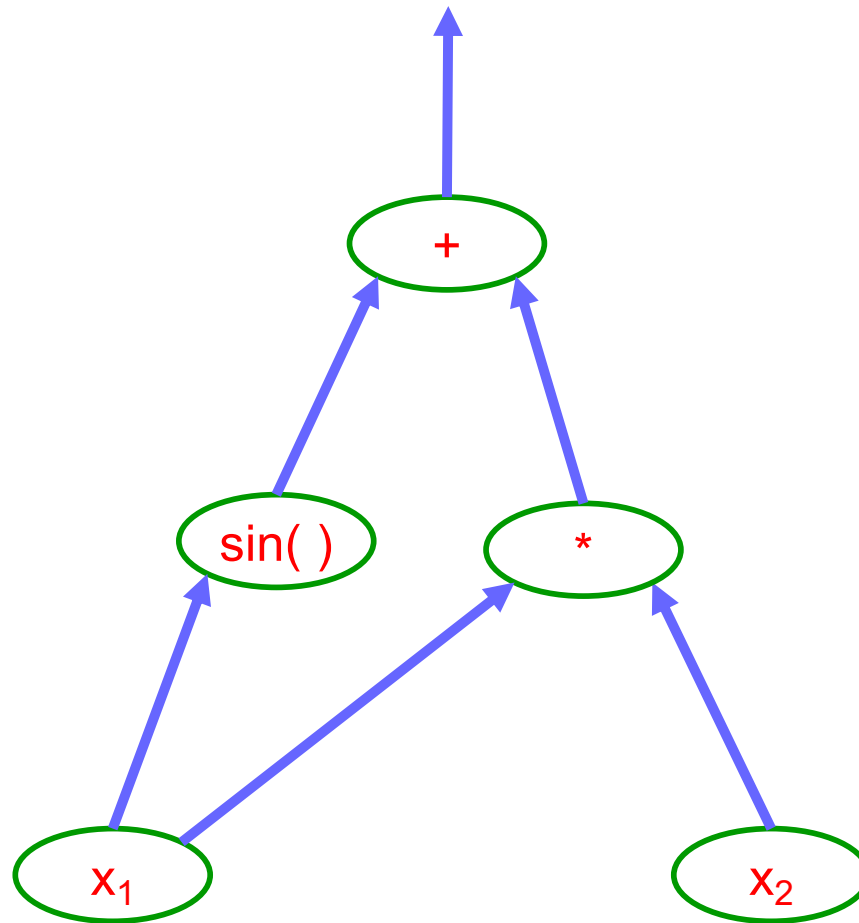
Example: Forward mode AD

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



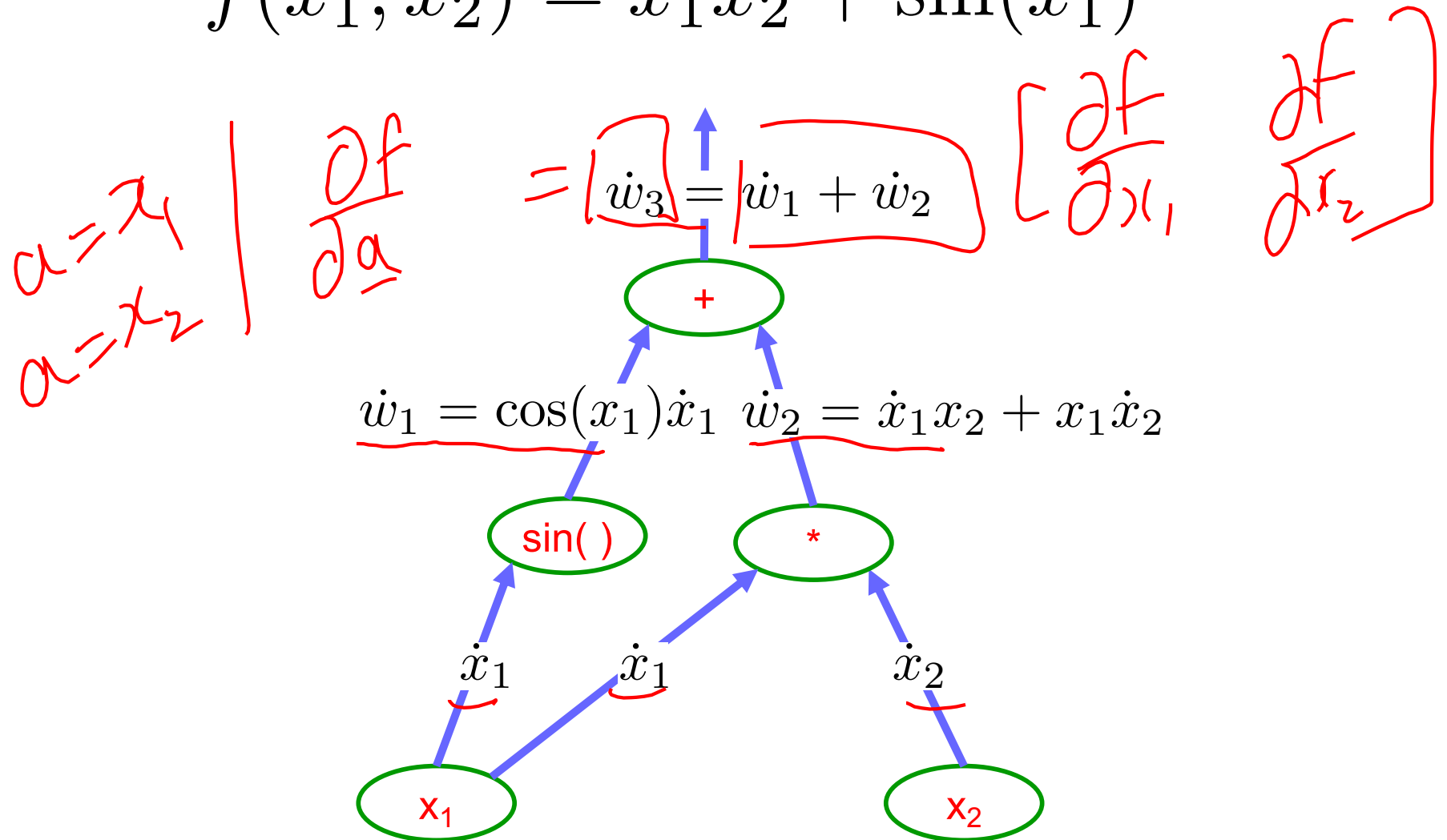
Example: Forward mode AD

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



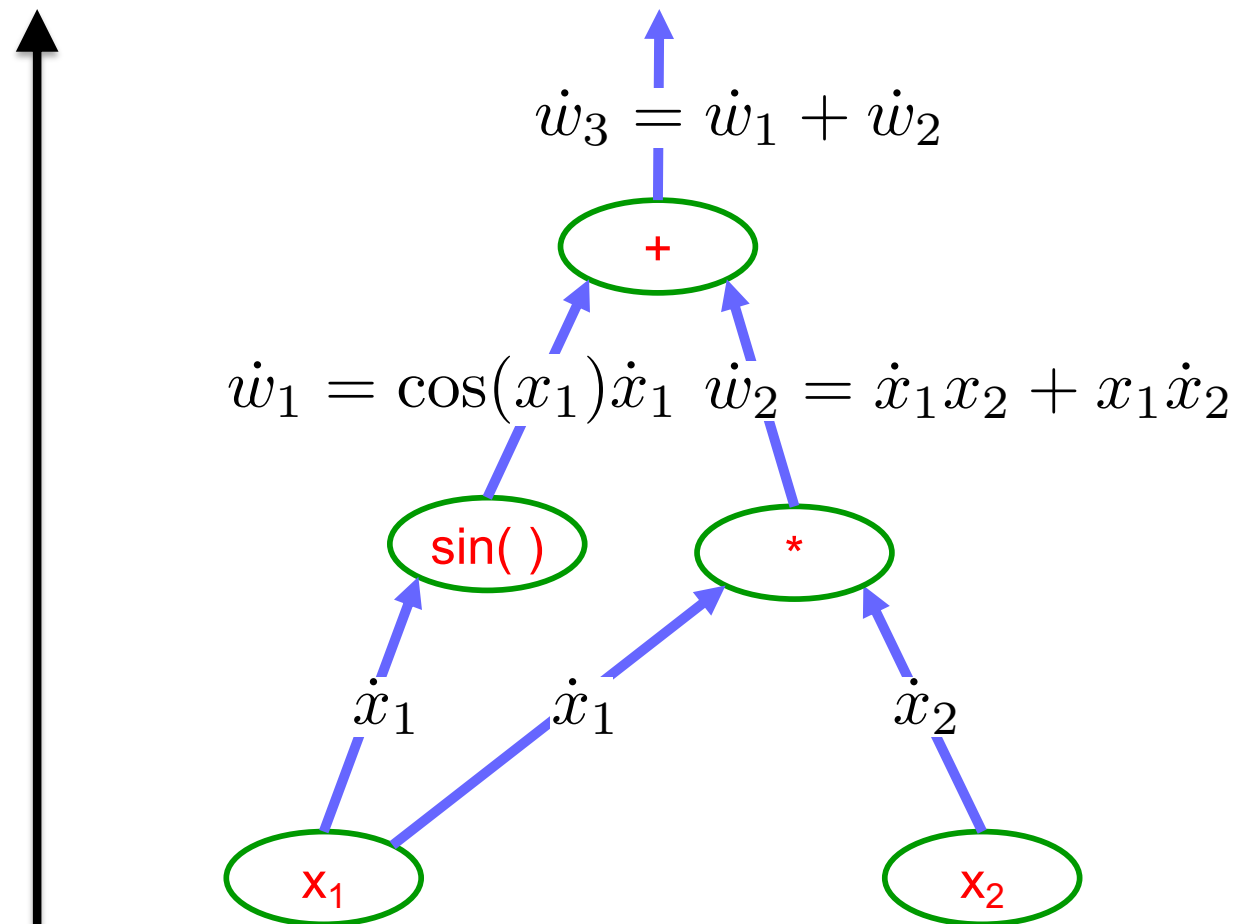
Example: Forward mode AD

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



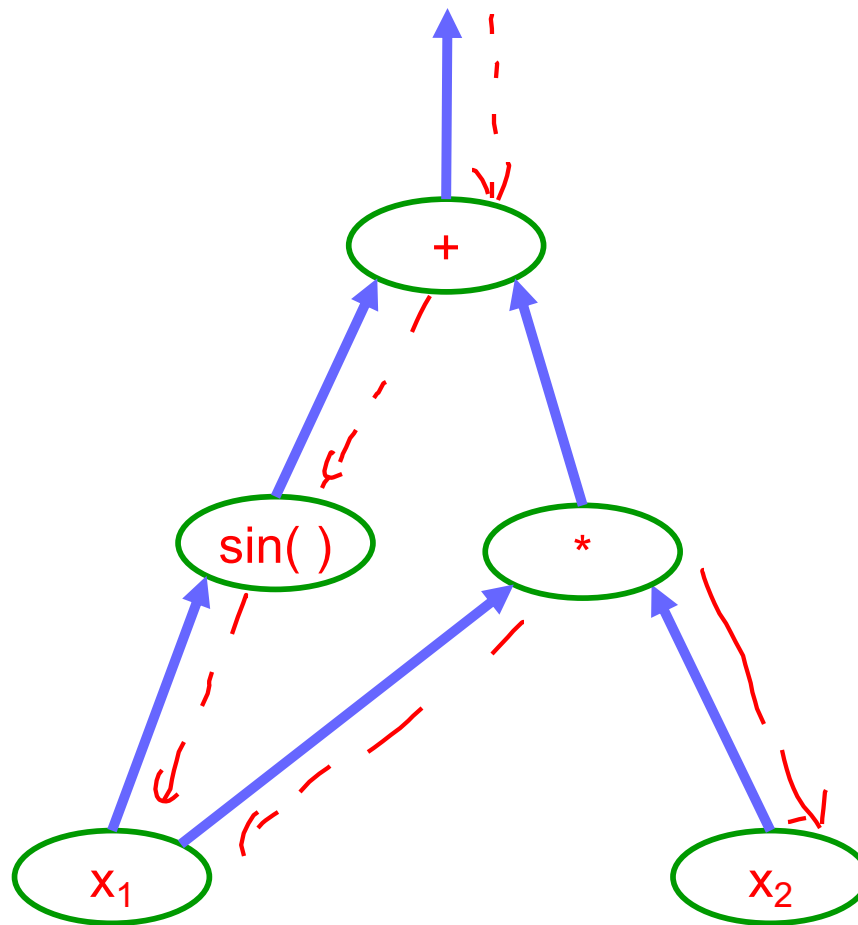
Example: Forward mode AD

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



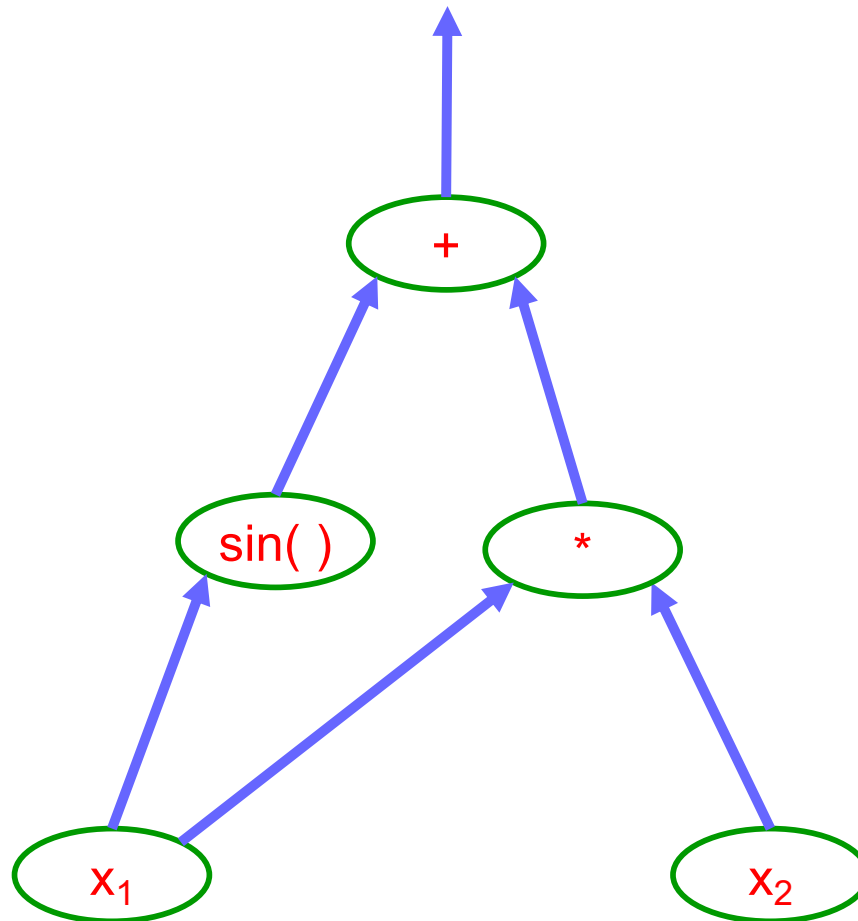
Example: Reverse mode AD

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



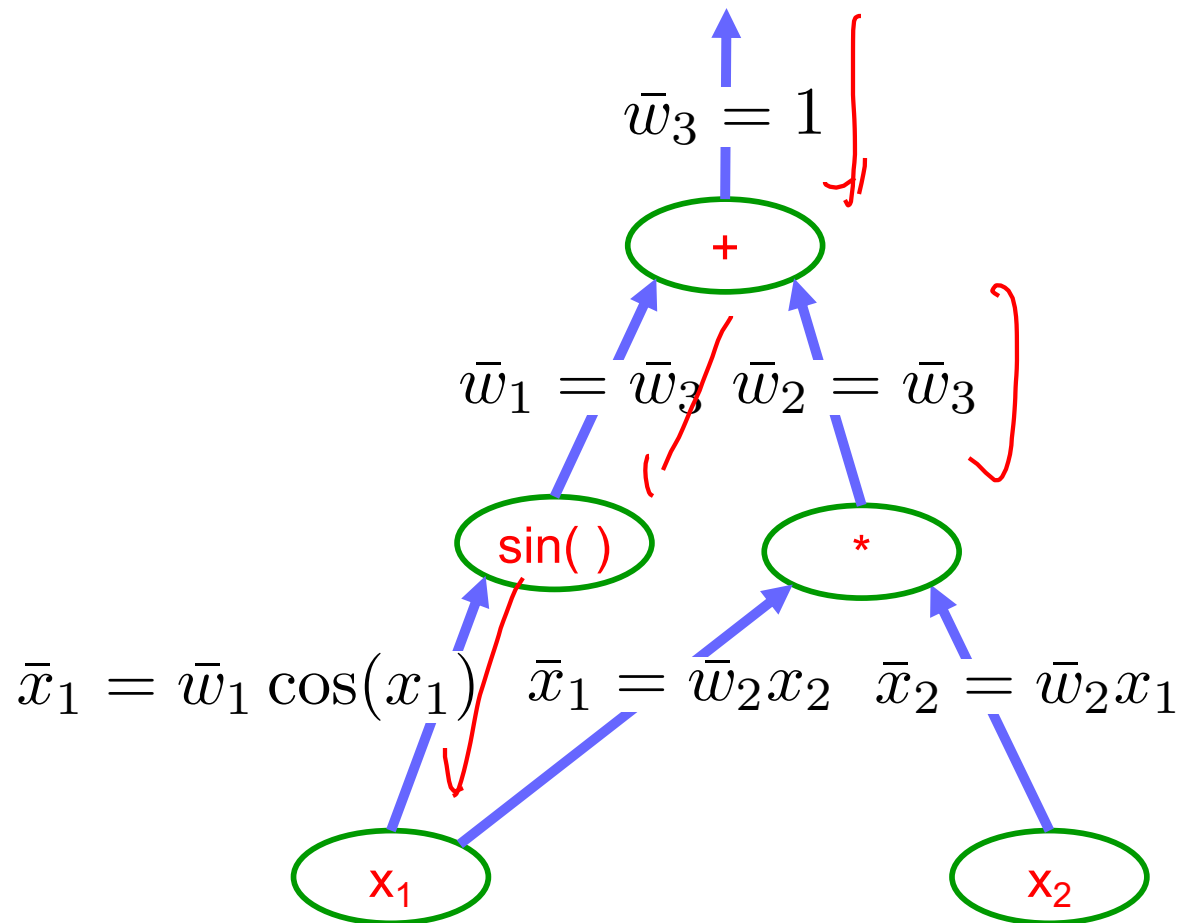
Example: Reverse mode AD

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



Example: Reverse mode AD

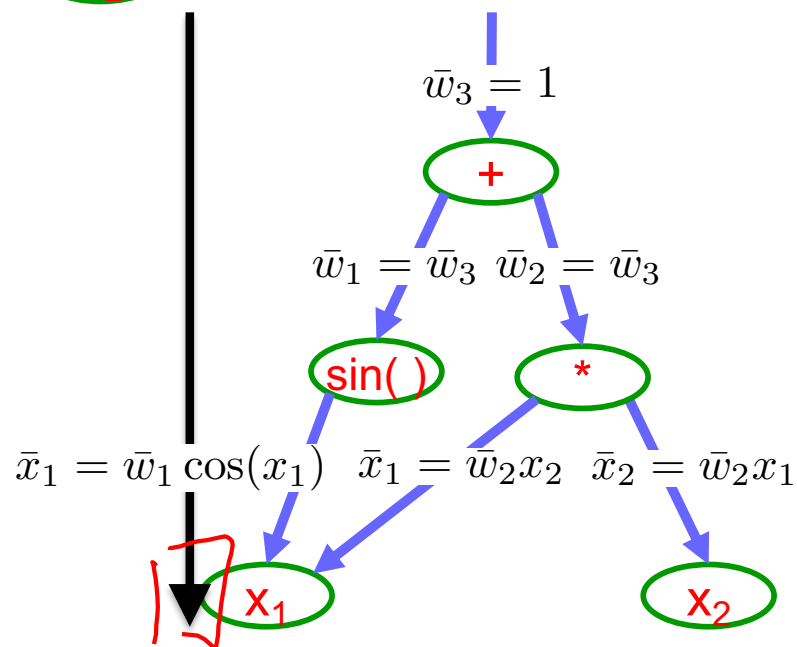
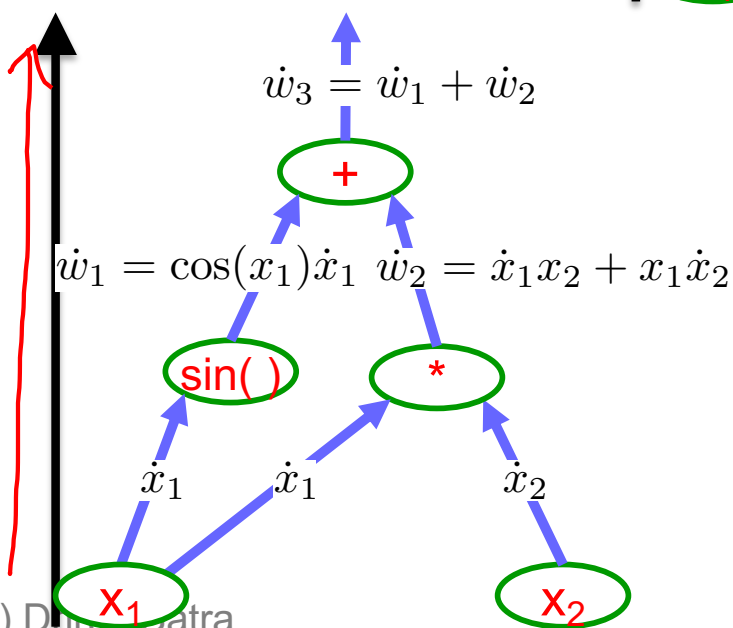
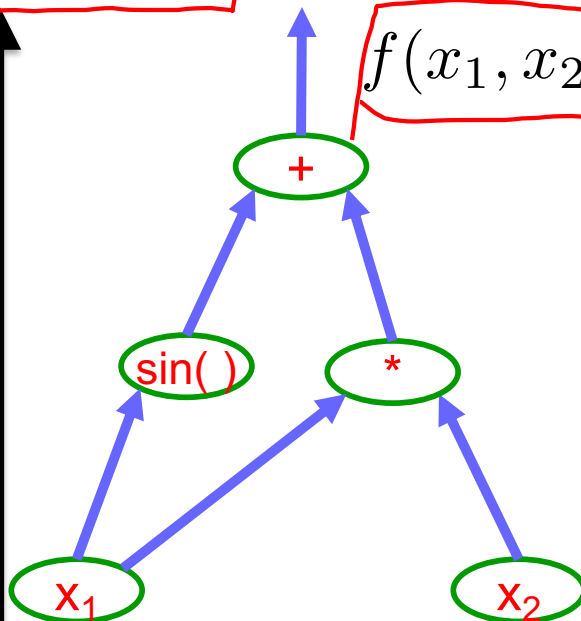
$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$



Forward Pass vs

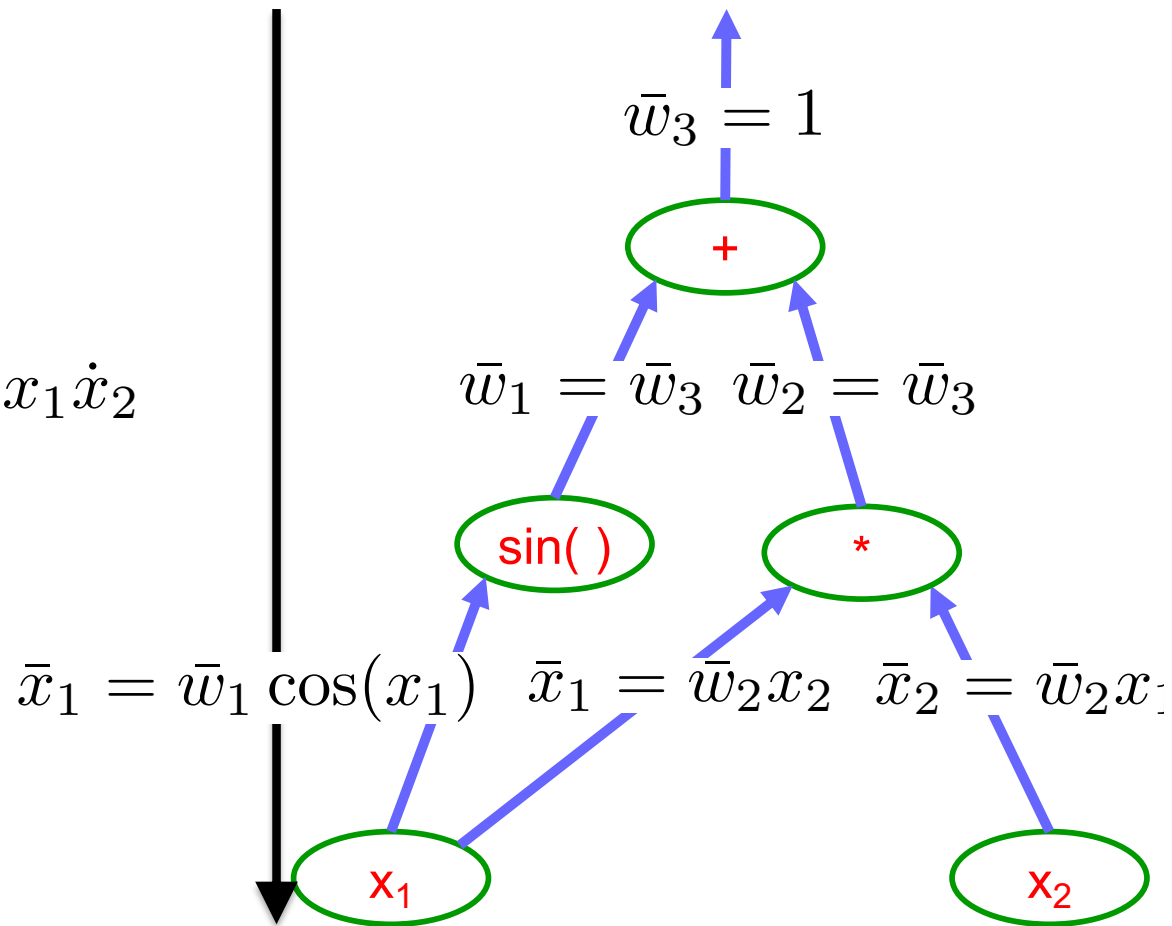
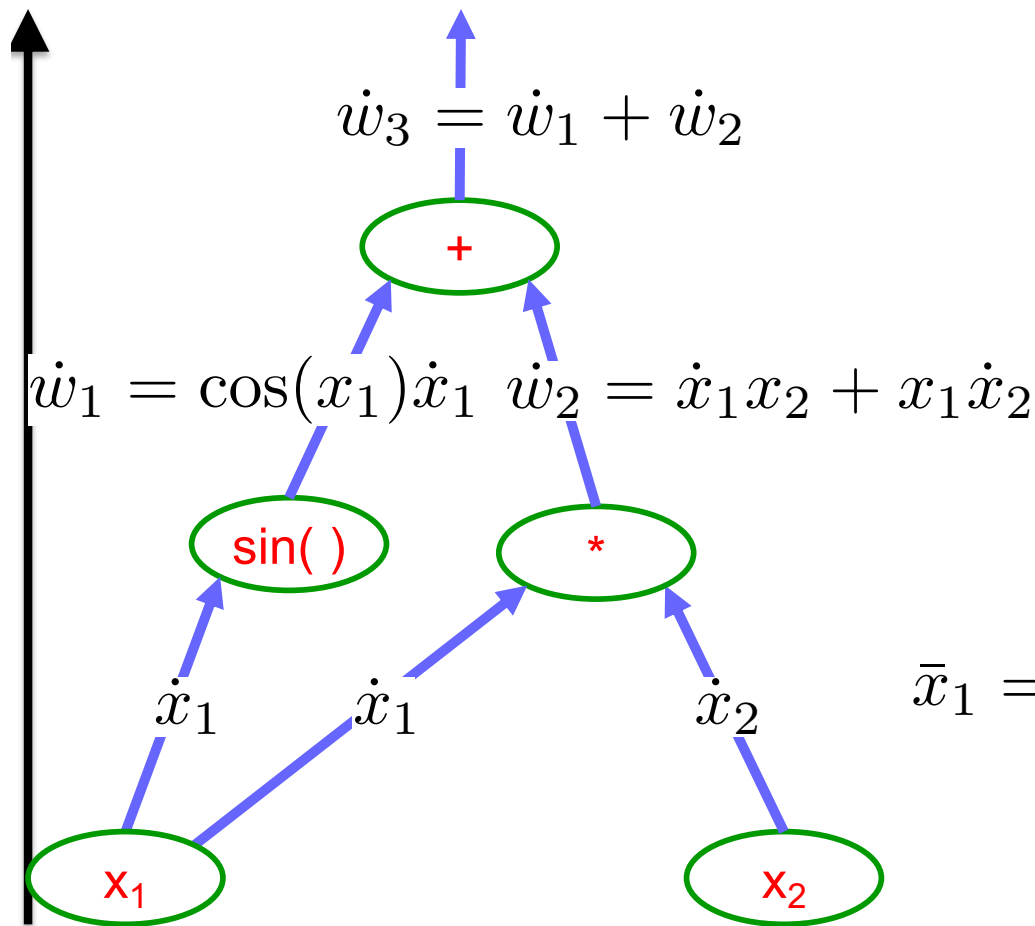
Forward mode AD vs Reverse Mode AD

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




Forward mode vs Reverse Mode

- What are the differences?



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?