# CS 4803 / 7643: Deep Learning

Topics:
- Optimization
- Computing Gradients

Dhruv Batra

Georgia Tech

# Administrativia

- HW0 Reminder
  - Due: 09/05, 11:55pm

- A note on expectations
  - Act like a responsible adult

- Thursday 09/06
  - Guest Lecture by Peter Anderson

- No class next week
  - 09/11, 09/13

- HW1 out next week (09/11)
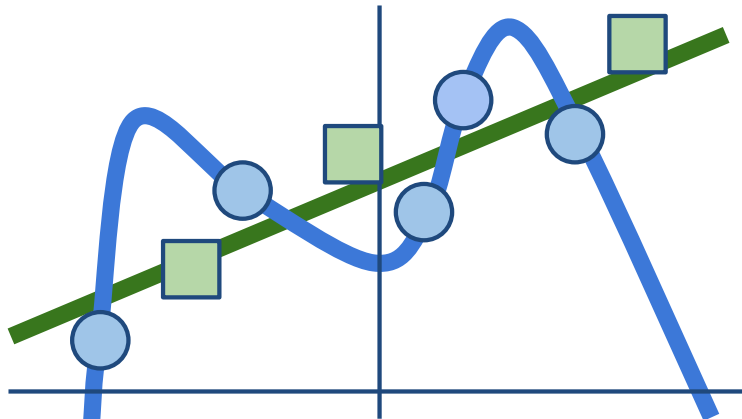
# Recap from last time

# Regularization

$\lambda$ = regularization strength (hyperparameter)

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

**Data loss**: Model predictions should match training data

**Regularization**: Prevent the model from doing *too* well on training data

**Occam's Razor**:
*"Among competing hypotheses, the simplest is the best"*
William of Ockham, 1285 - 1347

# Regularization

$\lambda$ = regularization strength (hyperparameter)

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

**Data loss**: Model predictions should match training data

**Regularization**: Prevent the model from doing *too* well on training data

**Simple examples**

L2 regularization: $R(W) = \sum_k \sum_l W_{k,l}^2$

L1 regularization: $R(W) = \sum_k \sum_l |W_{k,l}|$

Elastic net (L1 + L2): $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

**More complex**:

Dropout

Batch normalization

Stochastic depth, fractional pooling, etc

# Neural networks: without the brain stuff

(**Before**) Linear score function:   $f = Wx$

# Neural networks: without the brain stuff

(**Before**) Linear score function: $\quad f = Wx$

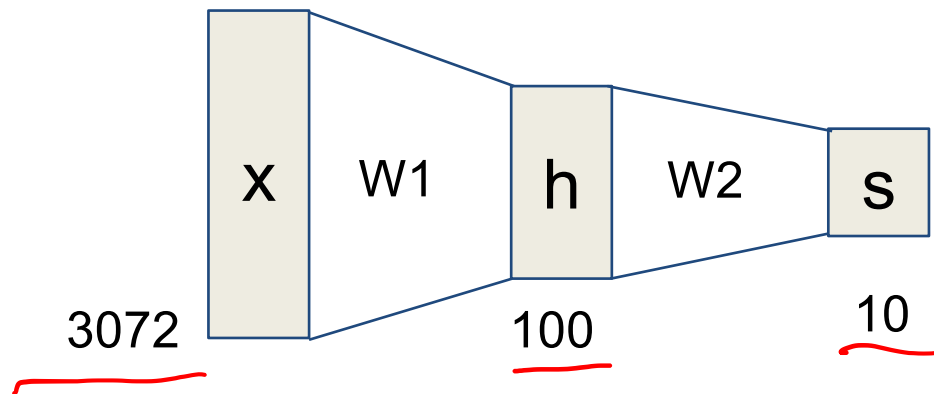(**Now**) 2-layer Neural Network $\quad f = W_2 \max(0, W_1 x)$

# Neural networks: without the brain stuff

(**Before**) Linear score function: $f = Wx$

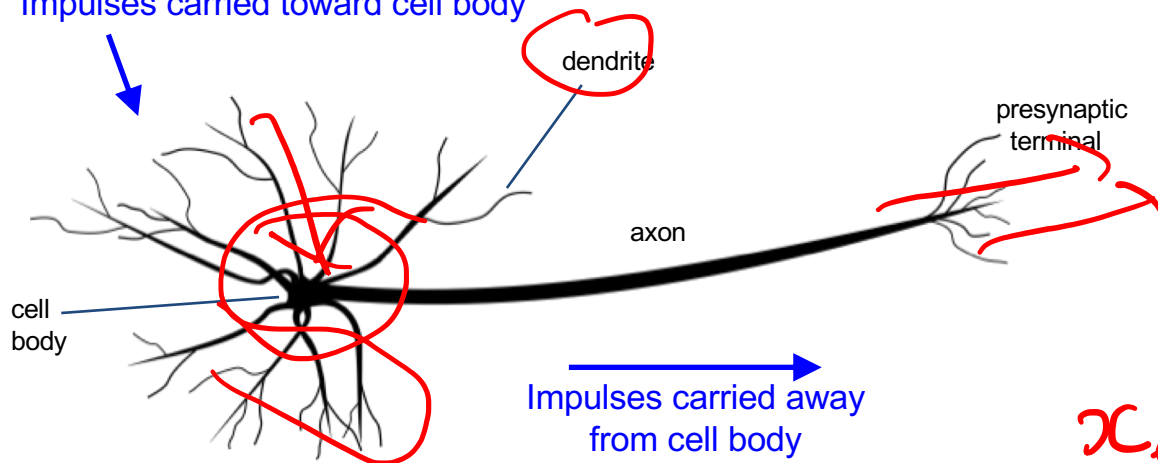(**Now**) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$

# Neural networks: without the brain stuff

(**Before**) Linear score function: $f = Wx$

(**Now**) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$
   or 3-layer Neural Network

$$f = W_3 \max(0, W_2 \max(0, W_1 x))$$

Impulses carried toward cell body

dendrite

presynaptic terminal

axon

cell body

This image by Felipe Perucho is licensed under CC-BY 3.0

Impulses carried away from cell body

$$a = \sum_j w_j x_j$$

$$= \underline{\underline{w^T x}}$$

$$x_0 \quad w_0$$

$$x_1$$

$$x_d \quad w_d \rightarrow$$

$$a$$

$$f$$

$$y = f(\underline{a})$$

Impulses carried toward cell body

dendrite

presynaptic terminal

axon

cell body

Impulses carried away from cell body

$x_0$   $w_0$

synapse

axon from a neuron

$w_0 x_0$

dendrite

cell body

$w_1 x_1$

$\sum_i w_i x_i + b$   $f$

$f\left(\sum_i w_i x_i + b\right)$

output axon

activation function

$w_2 x_2$

sigmoid activation function

$$\frac{1}{1 + e^{-x}}$$

$$f(a) \quad \frac{1}{1 + e^{-a}}$$

Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

# Activation functions

**Sigmoid**

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

**tanh**

$$\tanh(x)$$

**ReLU**

$$\max(0, x)$$

**Leaky ReLU**

$$\max(0.1x, x)$$
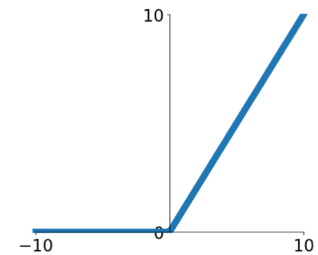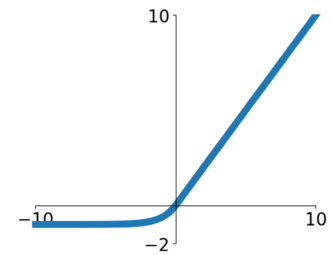
**Maxout**

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

**ELU**

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

# Activation Functions

- sigmoid vs tanh

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

$$\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$$

$-1/2$

$+1$

$-1$

$\sigma(\sigma(\sigma(a)))$

$= 2\sigma(2a) - 1$

# Multilayer Networks

- Cascade Neurons together
- The output from one layer is the input to the next
- Each Layer has its own sets of weights



input layer

hidden layer

output layer

input layer

hidden layer 1　hidden layer 2

output layer

$$h_i^{(1)} = f(w_i^T x)$$

# Plan for Today

- (Finish) Optimization
- Computing Gradients

# Optimization

Strategy: **Follow the slope**

$$\min_{w} \; L(w; D)$$

Strategy: **Follow the slope**

In 1-dimension, the derivative of a function:

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

$$\frac{\partial}{\partial x_i} f(x_1 \cdots x_d) = \frac{f(x_1 \cdots x_i + h, \cdots x_d) f(x)}{h}$$

$$\nabla_n f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix}$$

Strategy: **Follow the slope**

In 1-dimension, the derivative of a function:

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

In multiple dimensions, the **gradient** is the vector of (partial derivatives) along each dimension

The slope in any direction is the **dot product** of the direction with the gradient
The direction of steepest descent is the **negative gradient**

# Gradient Descent

```
# Vanilla Gradient Descent

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

← = backprop

$$W^{(0)} = init$$

$$\text{for } t=1. \cdots \text{ tired}$$

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

$0.00001 w_1^2 + w_2^2$

W_2

W_1

negative gradient direction

original W

$\mathcal{L}(W^{(i)})$

# 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
```

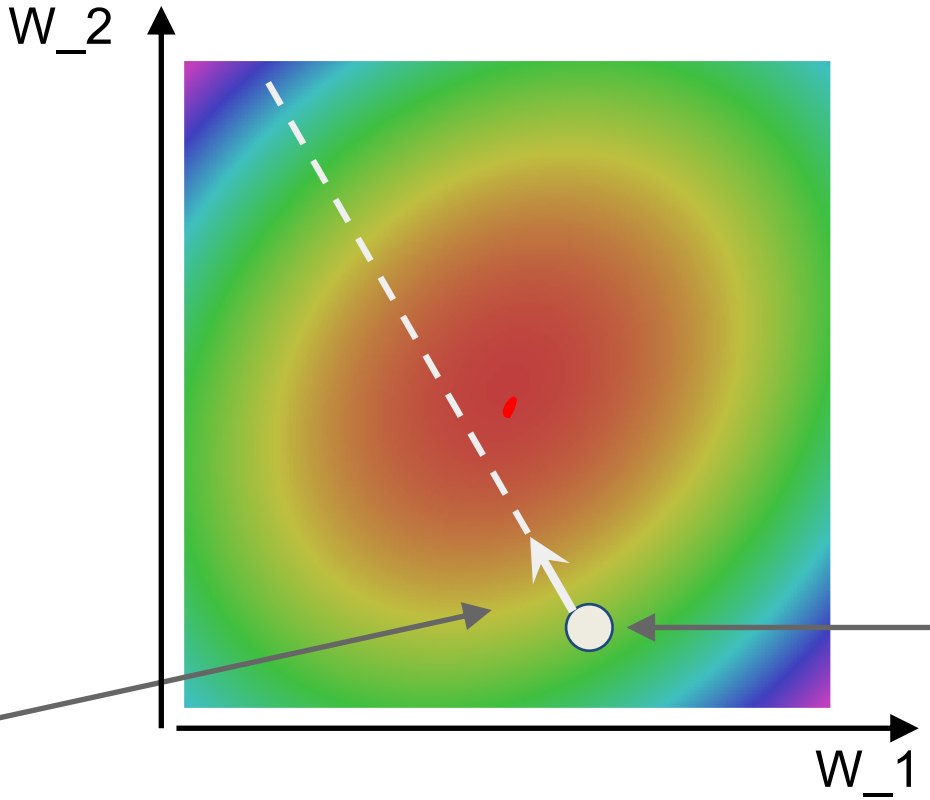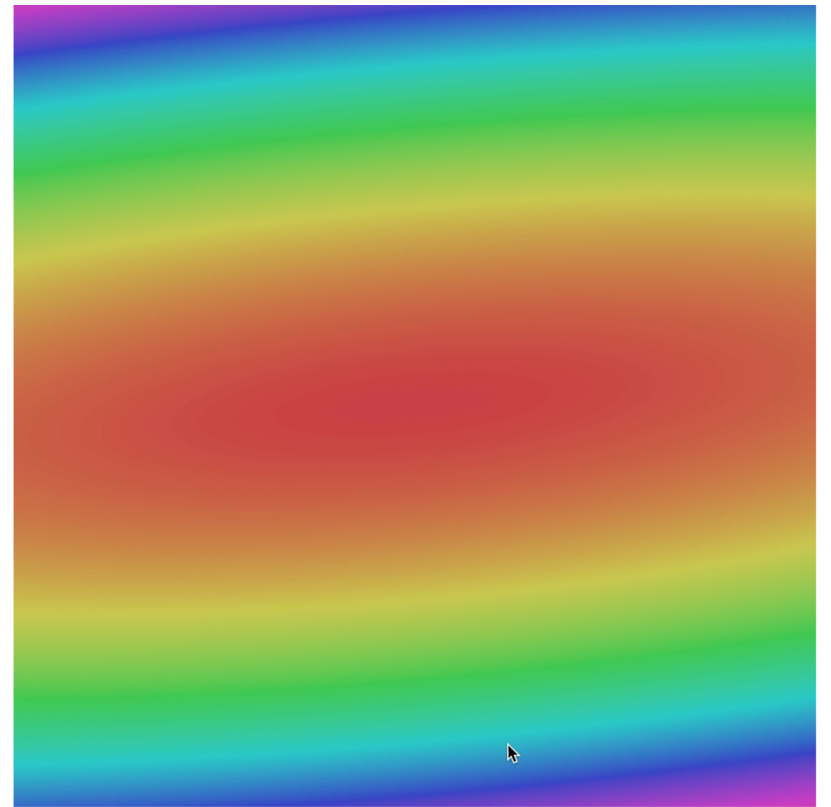Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

# Stochastic Gradient Descent (SGD)

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

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^{N} \nabla_W L_i(x_i, y_i, W) = \sum_i \left(\frac{1}{N}\right) \nabla L_i$$

$$= E_{I \sim U(\cdot)}[\nabla L_i]$$

$$I = i \in \{1, \ldots, N\}$$

$$I \sim U(1, N)$$

# Stochastic Gradient Descent (SGD)

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

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

$$\approx \nabla_W \; E_{x, Y \sim p^*} \left[ L(x, y, W) \right]$$

$$\nabla_W \int \sum_{y} L(x, y, W) \, p^*(x, y) \, dx$$

$$\underset{x}{} \quad \underset{y}{}$$

$$\int \sum_{y} \nabla_W L(x, y, W) \, p^*(x, y) \, dx$$

$$= E\left[ \nabla_W L \right] \approx \frac{1}{N} \sum \nabla_W L$$

# 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
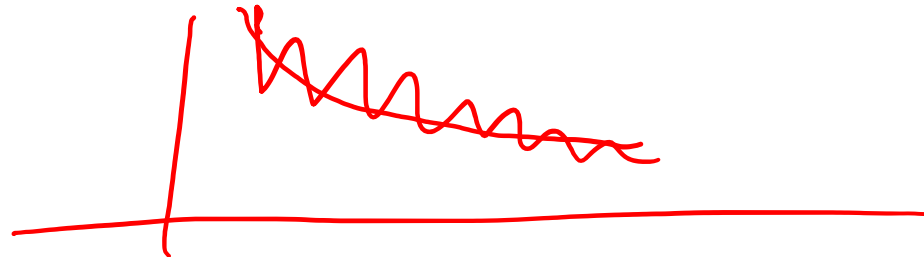32 / 64 / 128 common

```
# 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 + 16 x)(1 - 2
        x)^2 (1 - 8 x + 8 x^2) + 64 (1
        - x)(1 - 2 x)^2 (1 - 8 x + 8
        x^2)^2 - 64x(1 - 2 x)^2 (1 - 8
        x + 8 x^2)^2 - 256x(1 - x)(1 -
        2 x)(1 - 8 x + 8 x^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')
    (v,v')
```
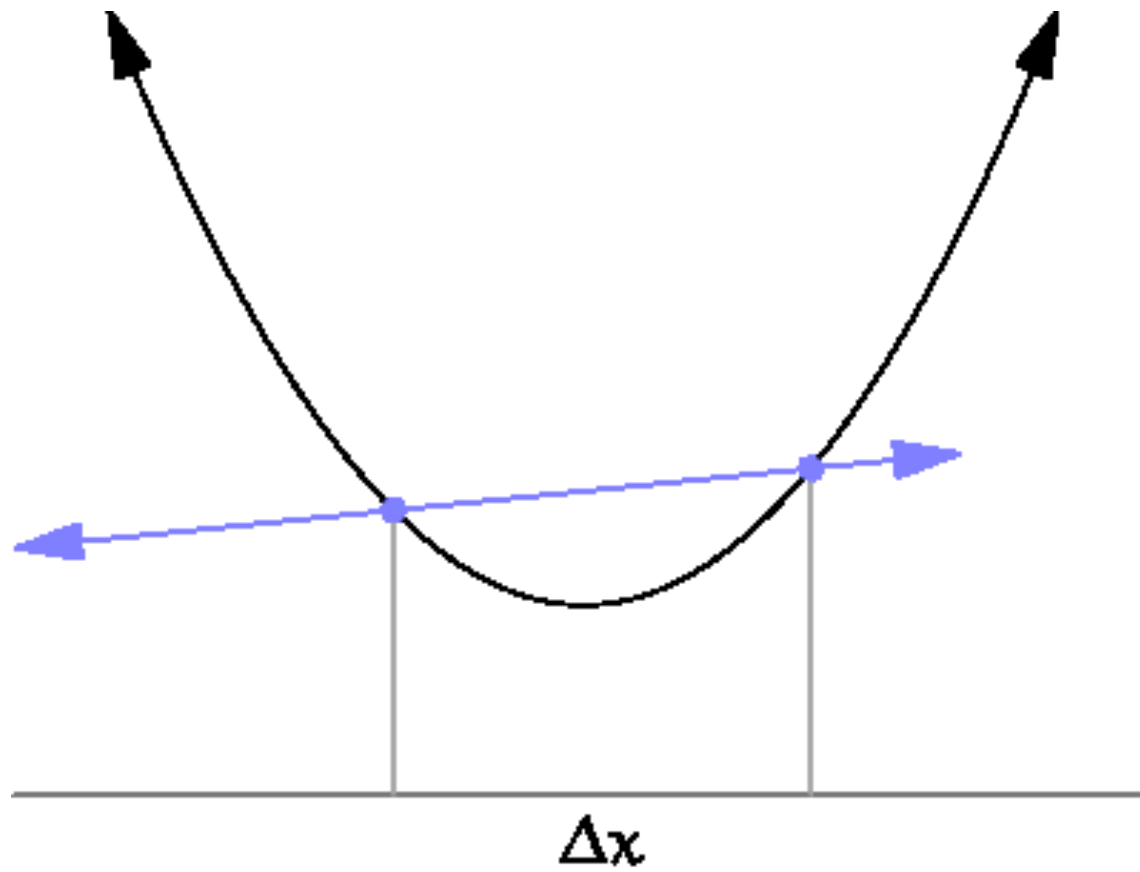
$f'(x_0) = f'(x_0)$
Exact

Numerical Differentiation →

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

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

(C) [

# How do we compute gradients?

- Analytic or "Manual" Differentiation

- Symbolic Differentiation

- Numerical Differentiation

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

$\Delta x$

$L(w)$

**current W:**

$w$

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

**loss 1.25347**

$L(w)$

**gradient dW:**

[?,
?,
?,
?,
?,
?,
?,
?,
?,…]

| current W: | W + h (first dim): | gradient dW: |
| --- | --- | --- |
| [0.34, | [0.34 + **0.0001**, | [?, |
| -1.11, | -1.11, | ?, |
| 0.78, | 0.78, | ?, |
| 0.12, | 0.12, | ?, |
| 0.55, | 0.55, | ?, |
| 2.81, | 2.81, | ?, |
| -3.1, | -3.1, | ?, |
| -1.5, | -1.5, | ?, |
| 0.33,…] | 0.33,…] | ?,…] |
| **loss 1.25347** | **loss 1.25322** | |

**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

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

?,
?,...]

| current W: | W + h (second dim): | gradient dW: |
|---|---|---|
| [0.34, | [0.34, | [-2.5, |
| -1.11, | -1.11 + **0.0001**, | ?, |
| 0.78, | 0.78, | ?, |
| 0.12, | 0.12, | ?, |
| 0.55, | 0.55, | ?, |
| 2.81, | 2.81, | ?, |
| -3.1, | -3.1, | ?, |
| -1.5, | -1.5, | ?, |
| 0.33,…] | 0.33,…] | ?,…] |
| **loss 1.25347** | **loss 1.25353** | |

**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

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

?,…]

| current W: | W + h (third dim): | gradient dW: |
|---|---|---|
| [0.34, | [0.34, | [-2.5, |
| -1.11, | -1.11, | 0.6, |
| 0.78, | 0.78 + **0.0001**, | ?, |
| 0.12, | 0.12, | ?, |
| 0.55, | 0.55, | ?, |
| 2.81, | 2.81, | ?, |
| -3.1, | -3.1, | ?, |
| -1.5, | -1.5, | ?, |
| 0.33,…] | 0.33,…] | ?,…] |
| **loss 1.25347** | **loss 1.25347** | |

**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

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

?,...]

# Numerical vs Analytic Gradients

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

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

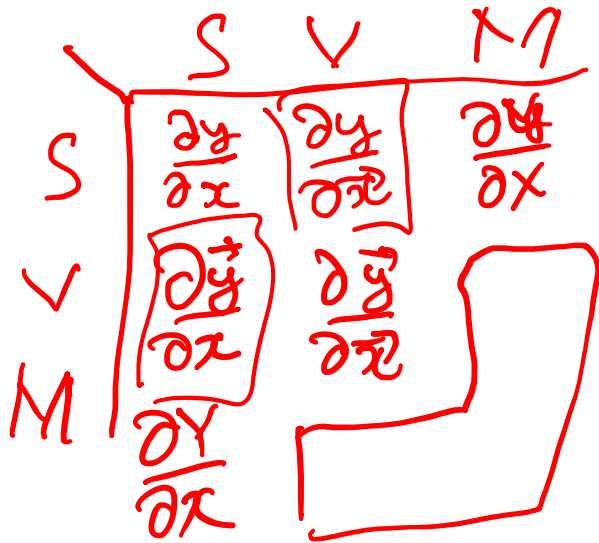In practice: Derive analytic gradient, check your implementation with numerical gradient.
This is called a **gradient check.**

# 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

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

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

$$\frac{\partial \vec{y}}{\partial x} = \begin{bmatrix} \dfrac{\partial y_1}{\partial x} \\ \dfrac{\partial y_2}{\partial x} \\ \cdot \\ \dfrac{\partial y_c}{\partial x} \end{bmatrix} \quad \downarrow num = dim\ 1$$

$$\longrightarrow den = dim\ 2$$

$$\frac{\partial y}{\partial \vec{x}} = \begin{bmatrix} \dfrac{\partial y}{\partial x_1} & \dfrac{\partial y}{\partial x_2} & \cdots & \dfrac{\partial y}{\partial x_d} \end{bmatrix}$$

# Matrix/Vector Derivatives Notation

$$\frac{\partial \vec{y}}{\partial \vec{x}} = \quad i - \left[ \quad \frac{\partial \dot{y}_i}{\partial x_j} \quad \right]_{c \times d}$$

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

$$\sum_i x_i w_i$$

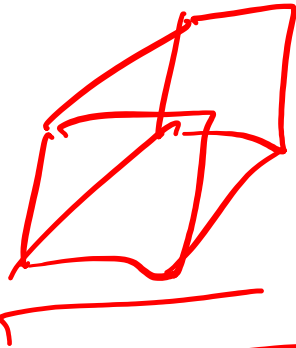$$= x^T \quad \left[ x_1 \quad \cdots \quad x_d \right] = x^T$$

# Vector Derivative Example

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

$$y_i = \sum_k a_{ik} x_k$$

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

$$\frac{\partial \vec{y}}{\partial \vec{x}} = A \quad i \left\{ \left[ \left(\frac{\partial y_i}{\partial x_j}\right) a_{ij} \right] \right.$$

# Extension to Tensors

$$X \in \mathbb{R}^{d_1 \times d_2 \cdots \times d_n}$$

$$Y \in \mathbb{R}^{c_1 \times c_2 \cdots c_m}$$

$$\frac{\partial \, Y[i_1, i_2 \cdots i_m]}{\partial \, X[j \cdots, j_n]}$$

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

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

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

# Chain Rule: Composite Functions

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

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

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

# Chain Rule: Scalar Case

$$x \rightarrow z \rightarrow y = f(g(x))$$

$$z = g(x)$$
$$y = f(z)$$

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial z} \frac{\partial z}{\partial x}$$
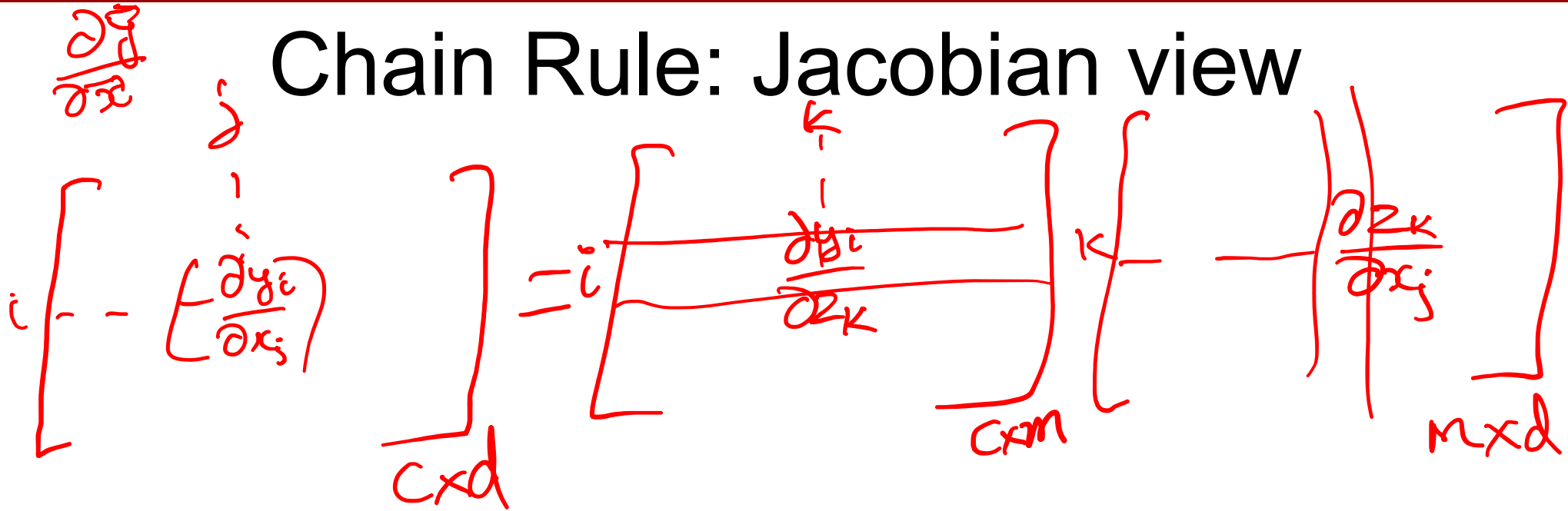
# Chain Rule: Vector Case

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

$$\vec{z} = g(\vec{x})$$
$$g: \mathbb{R}^d \to \mathbb{R}^m$$

$$\vec{y} = f(\vec{z})$$
$$f: \mathbb{R}^m \to \mathbb{R}^c$$

$$\underbrace{\frac{\partial \vec{y}}{\partial \vec{x}}}_{J_{f \circ g}} = \underbrace{\frac{\partial \vec{y}}{\partial \vec{z}}}_{J_f} \cdot \underbrace{\frac{\partial \vec{z}}{\partial \vec{x}}}_{J_g}$$
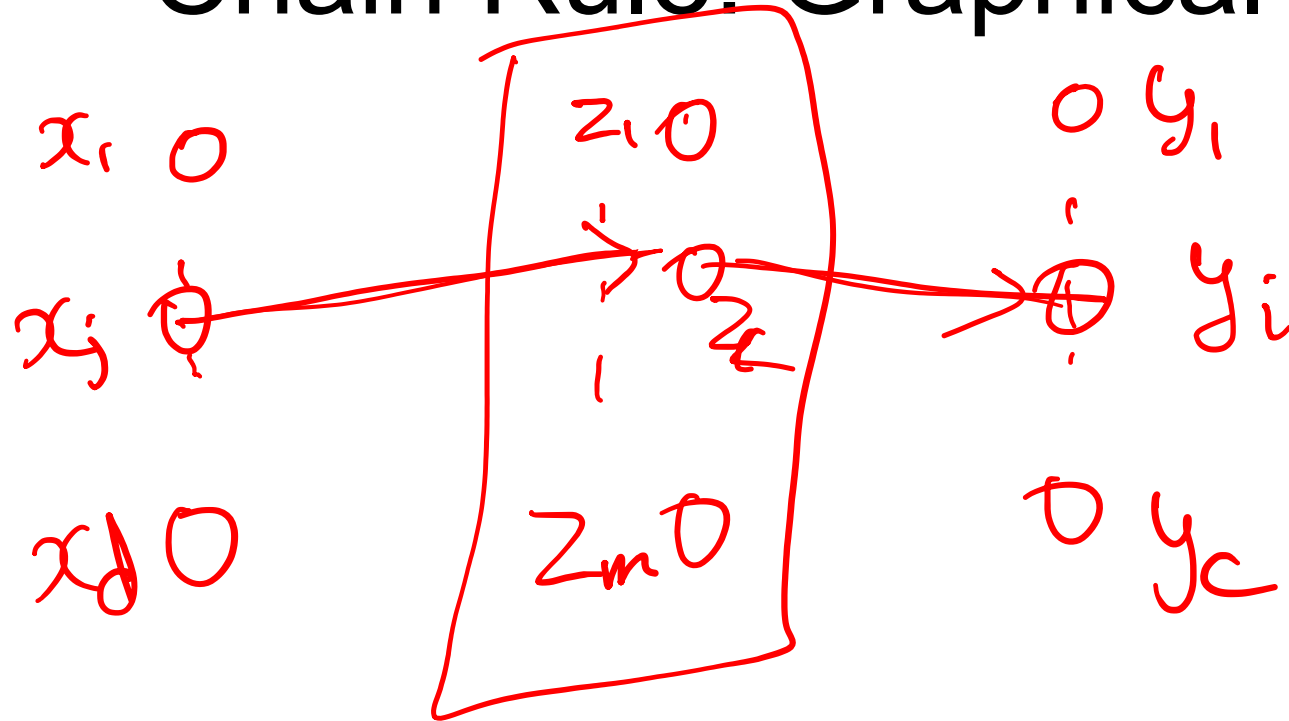
# Chain Rule: Jacobian view

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

$$i\left[-- \left[\frac{\partial y_i}{\partial x_j}\right] \right]_{c\times d} = i\left[\frac{\partial y_i}{\partial z_k}\right]_{c\times m} \left[\frac{\partial z_k}{\partial x_j}\right]_{m\times d}$$

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

# Chain Rule: Graphical view



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

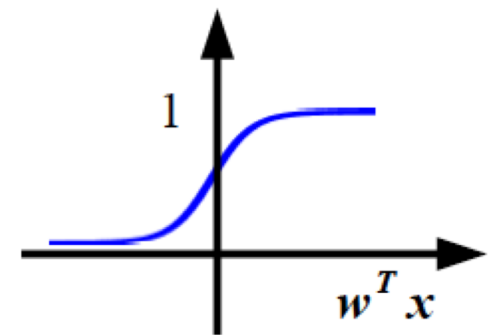# Linear Classifier: Logistic Regression

Input: $x \in R^D$

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

Parameters: $w \in R^D$

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

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

Log Loss

Ranzato

# Logistic Regression Derivatives

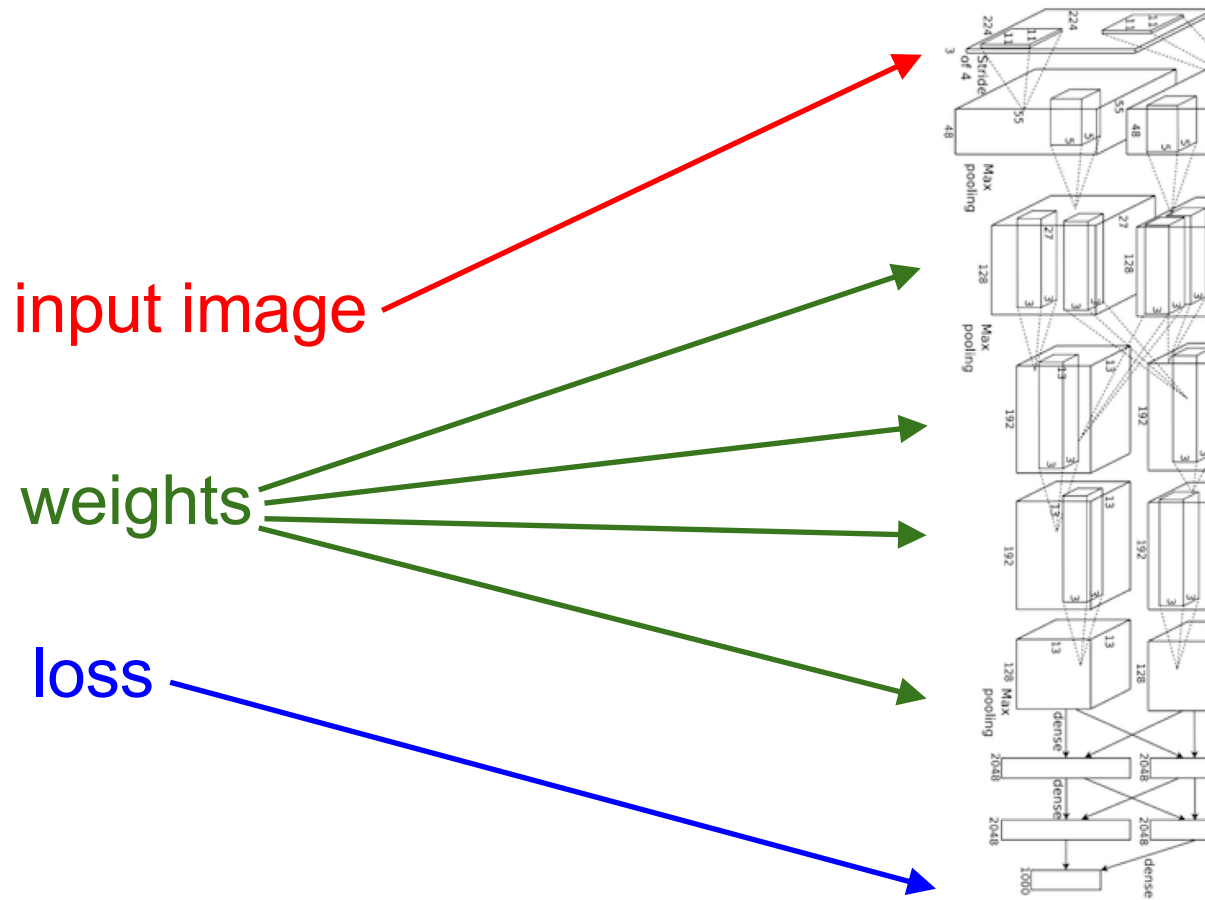# Logistic Regression Derivatives

# Convolutional network (AlexNet)



input image

weights

loss

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.
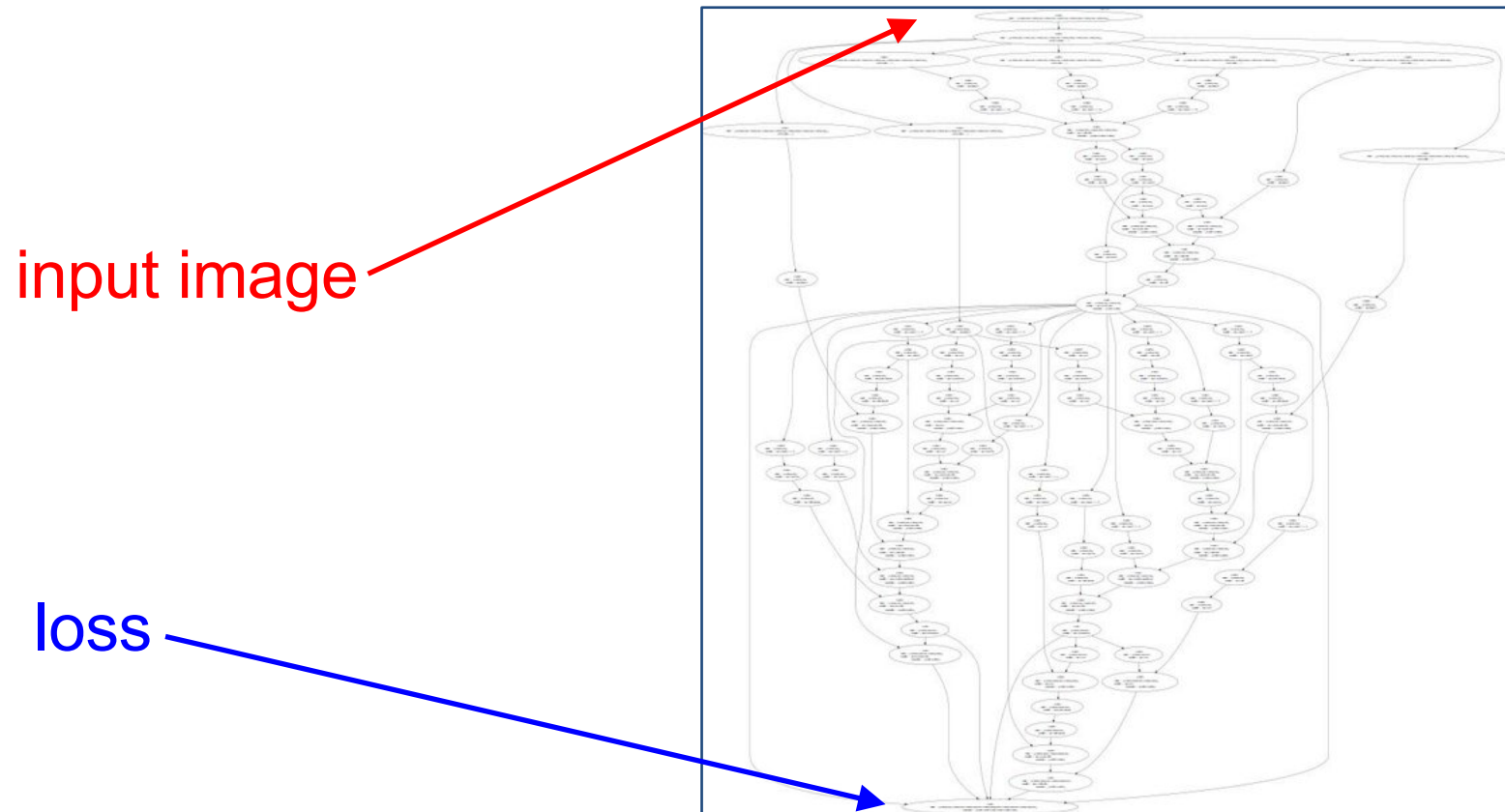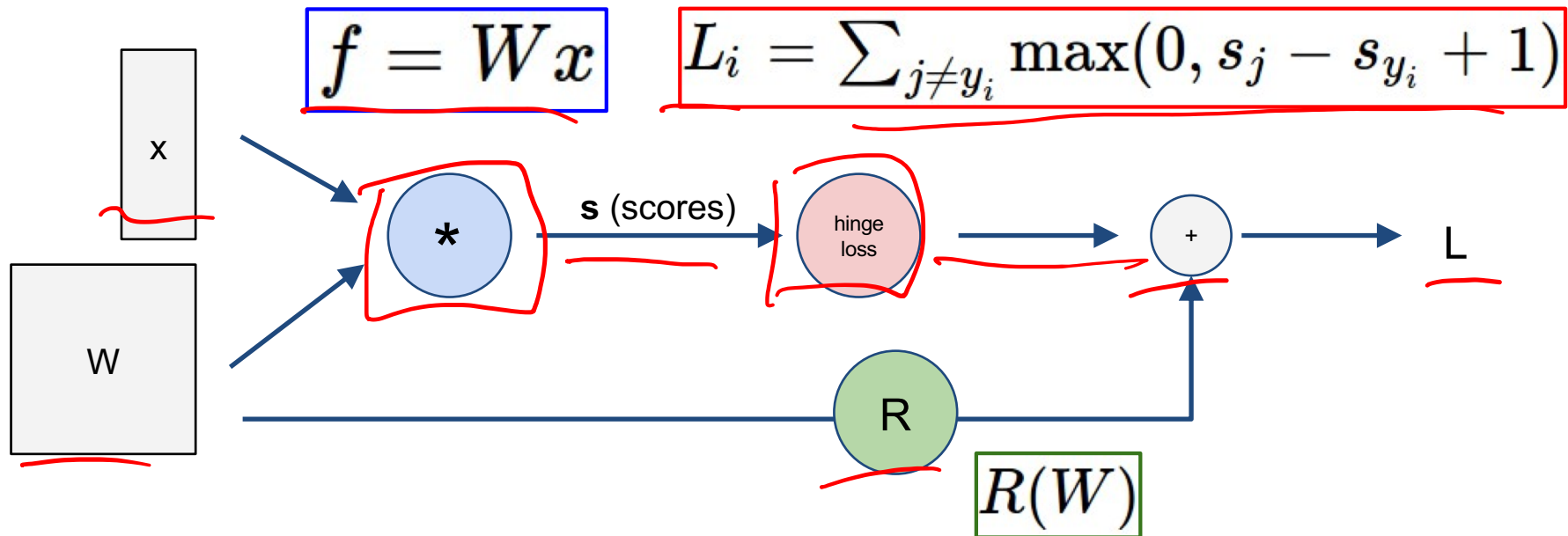
# Neural Turing Machine



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
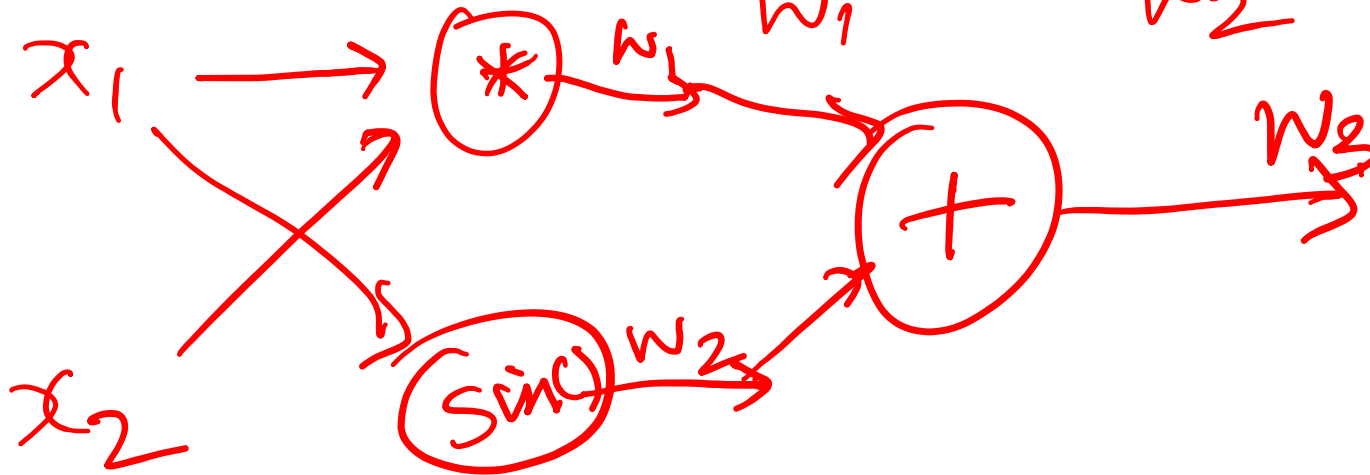    - aka "backprop"

# Computational Graph

$$f = Wx$$

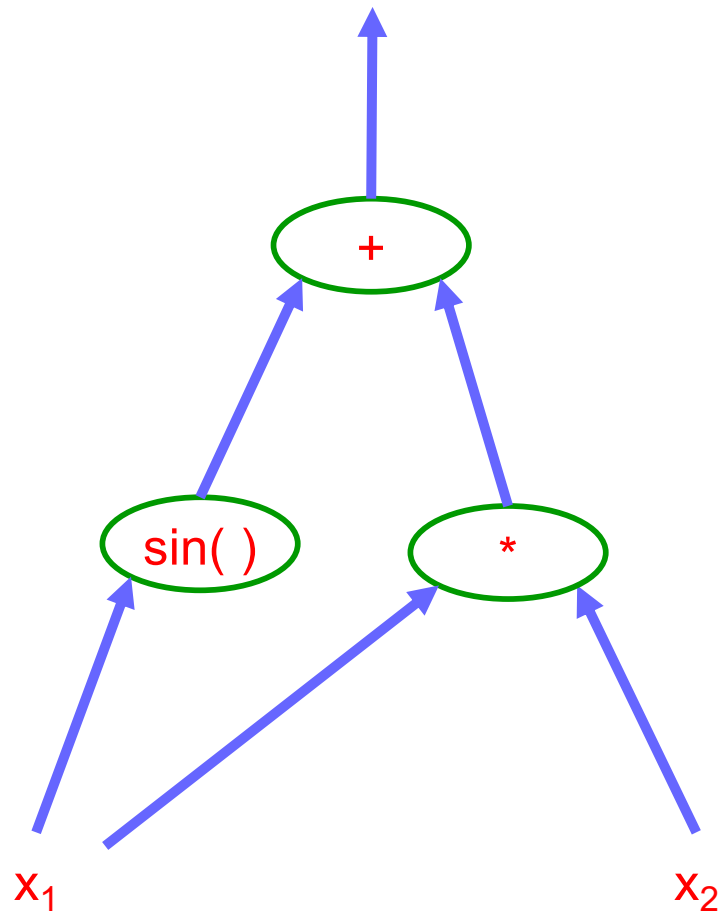$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$



**s** (scores)

*

hinge loss

+

L

x

W

R

$$R(W)$$

# Computational Graphs

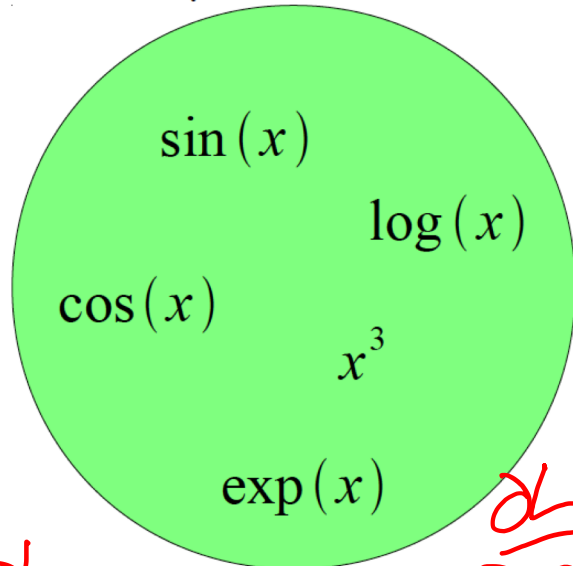- Notation

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

# Example

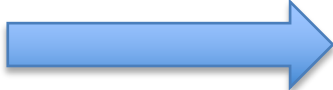$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

# Logistic Regression as a Cascade

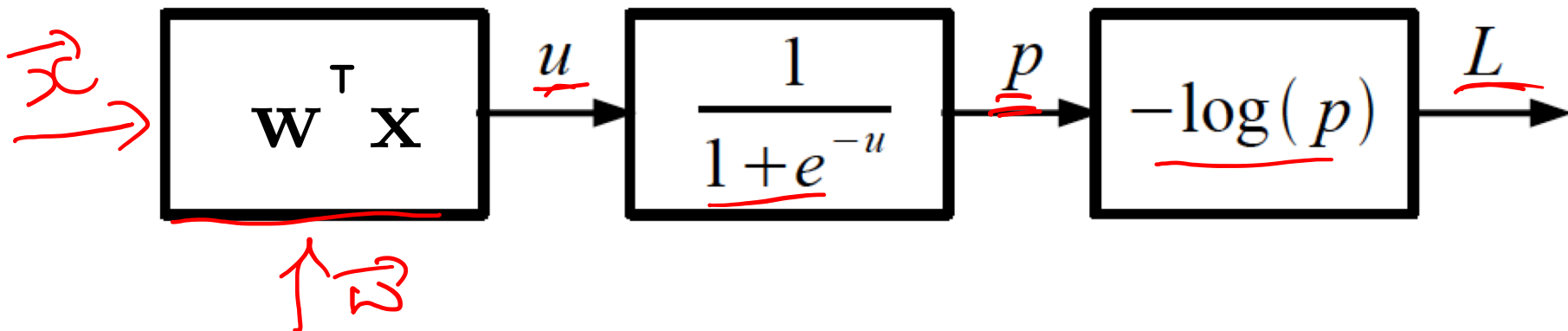Given a library of simple functions

$$P\left(Y=1 \mid \vec{x}_i, \vec{w}\right)$$

Compose into a complicate function
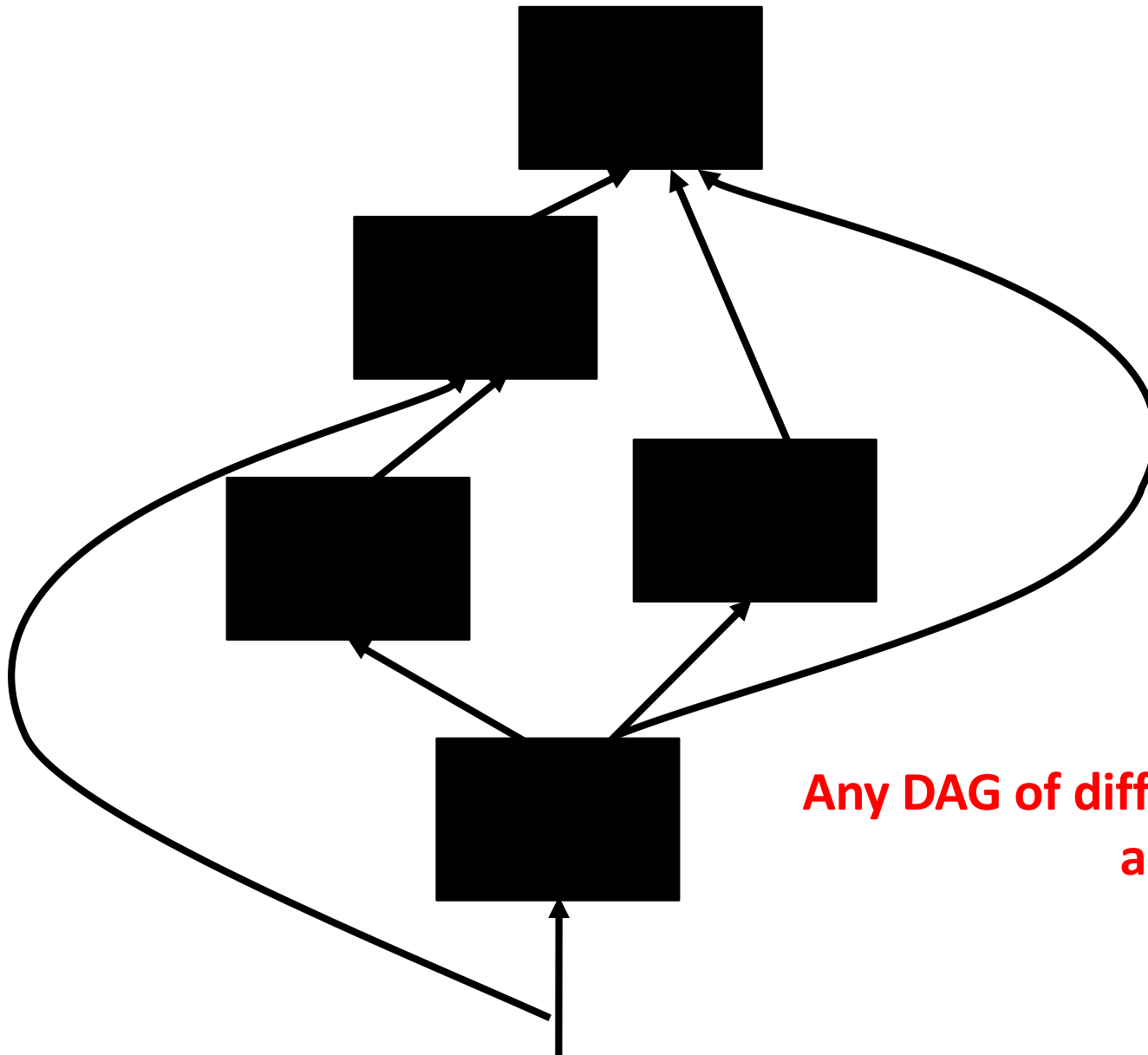
$$-\log\left(\frac{1}{1+e^{-\mathbf{w}^{\mathsf{T}}\mathbf{x}}}\right)$$

$$\sin(x)$$
$$\log(x)$$
$$\cos(x)$$
$$x^3$$
$$\exp(x)$$

$$\frac{d}{dw} = \frac{d}{dp}\frac{\partial p}{\partial v}$$

$$\frac{\partial L}{\partial p} = \frac{-1}{p} \quad \frac{\partial L}{\partial L} = 1$$

$$\frac{d}{dw}$$

$$\vec{x} \rightarrow \boxed{\mathbf{w}^{\mathsf{T}}\mathbf{x}} \xrightarrow{u} \boxed{\frac{1}{1+e^{-u}}} \xrightarrow{p} \boxed{-\log(p)} \xrightarrow{L}$$

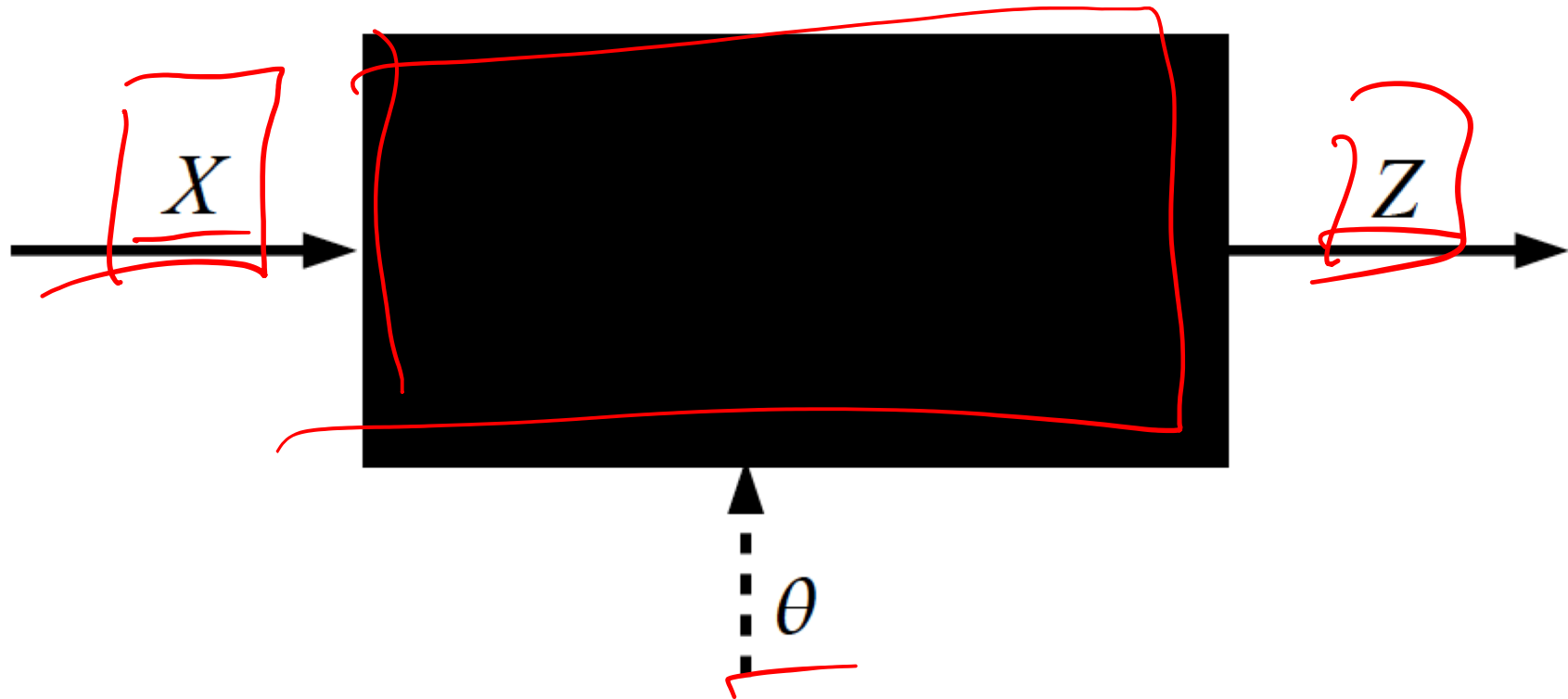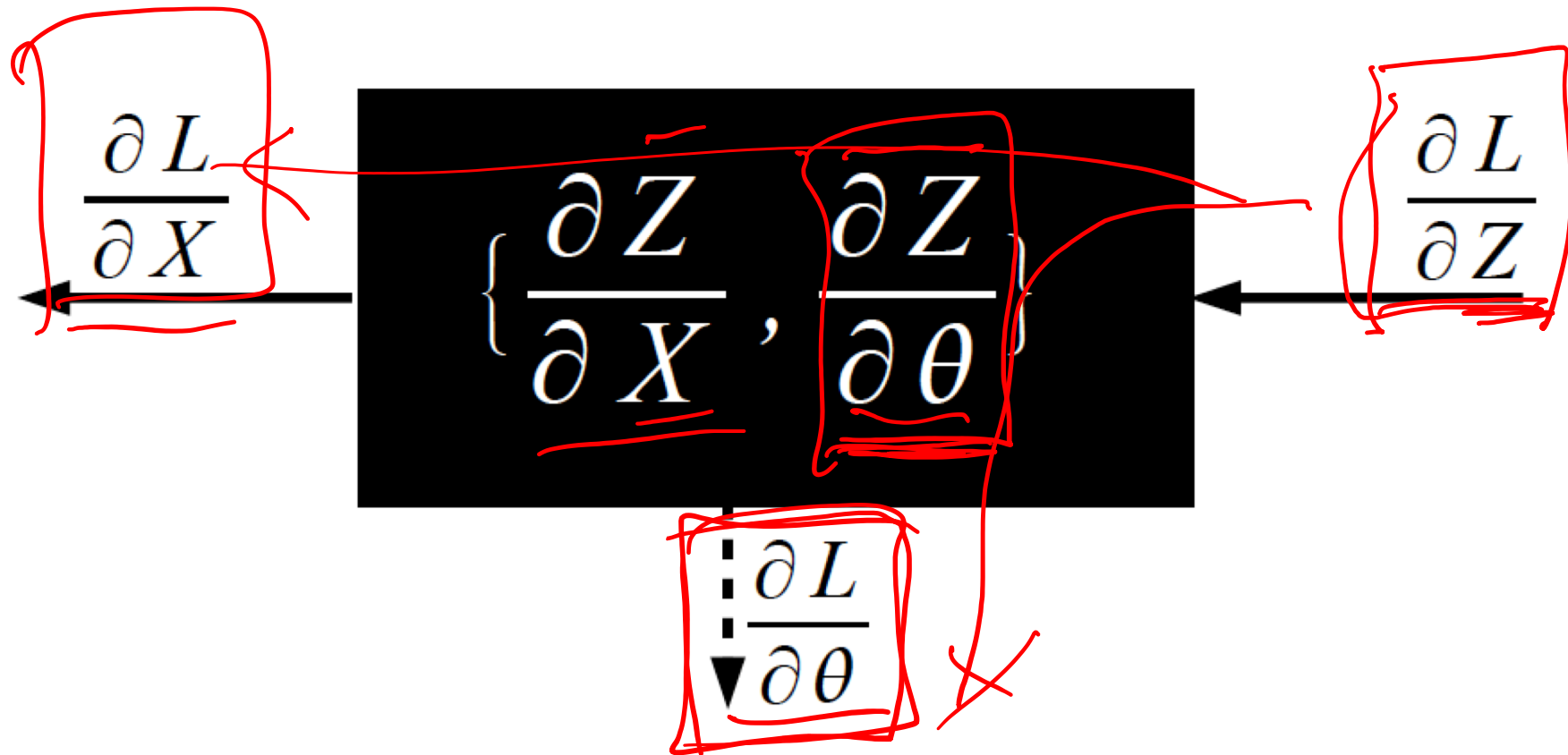$$\uparrow \vec{w}$$

Slide Credit: Marc'Aurelio Ranzato, Yann LeCun

# Computational Graph



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

# Key Computation: Forward-Prop

# Key Computation: Back-Prop



$$\frac{\partial L}{\partial X} \qquad \left\{ \frac{\partial Z}{\partial X}, \frac{\partial Z}{\partial \theta} \right\} \qquad \frac{\partial L}{\partial Z}$$

$$\frac{\partial L}{\partial \theta}$$
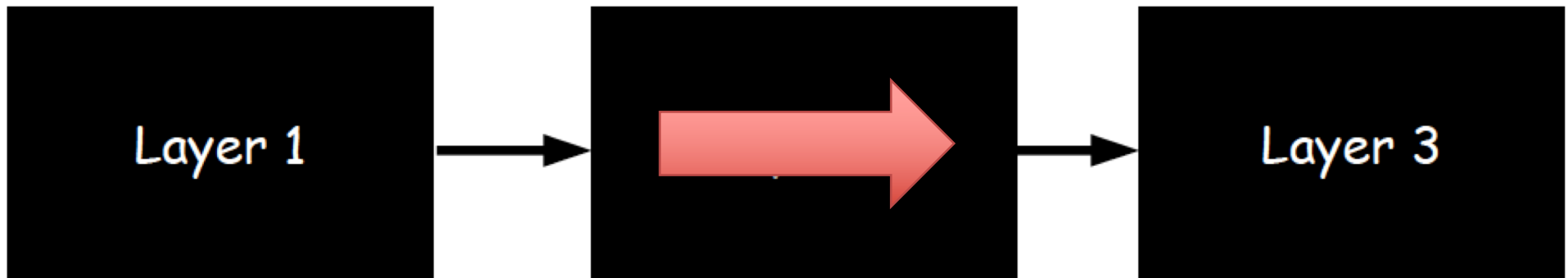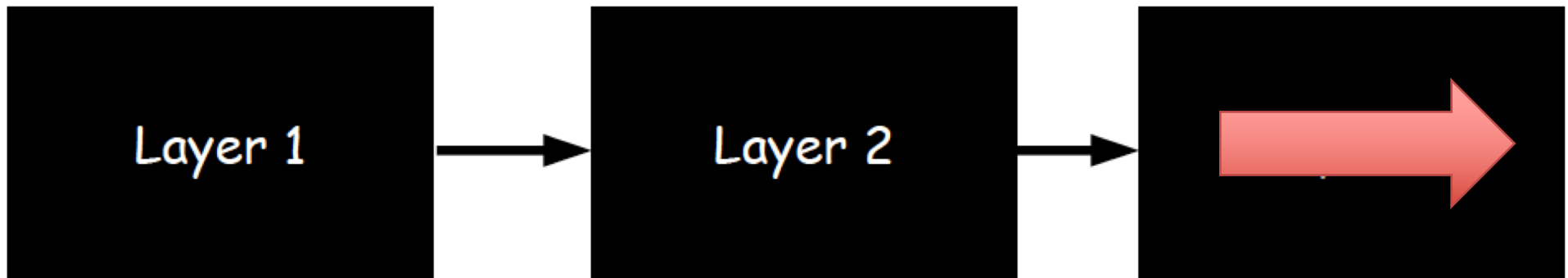
# Neural Network Training

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

# Neural Network Training

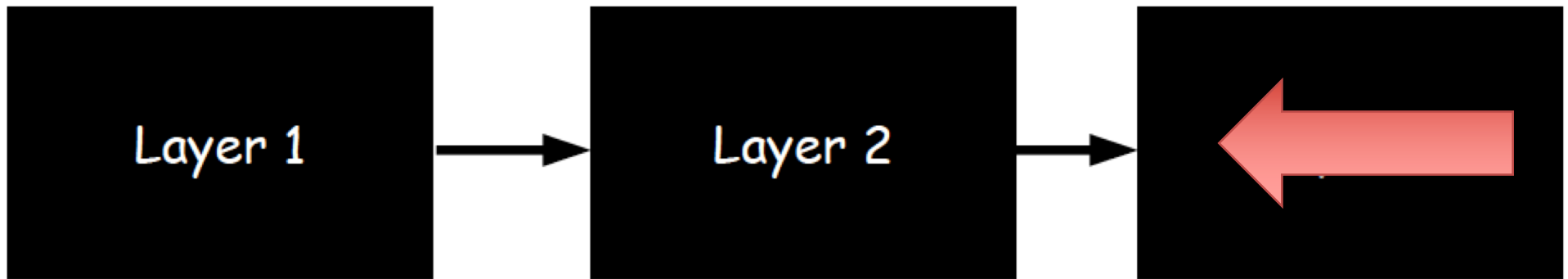- Step 1: Compute Loss on mini-batch        [F-Pass]

# Neural Network Training

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

# Neural Network Training

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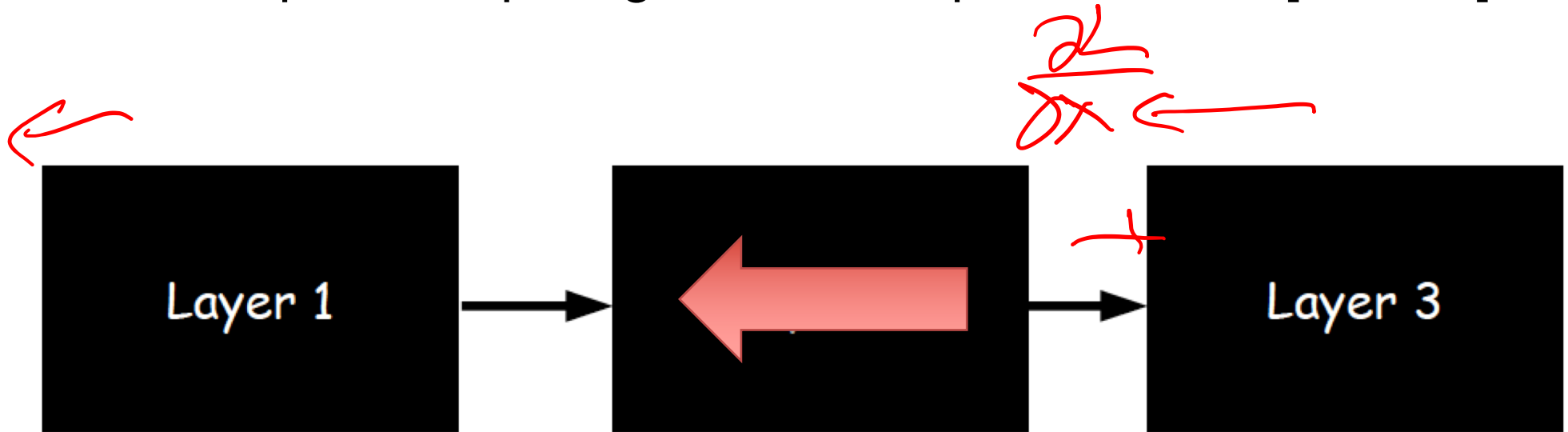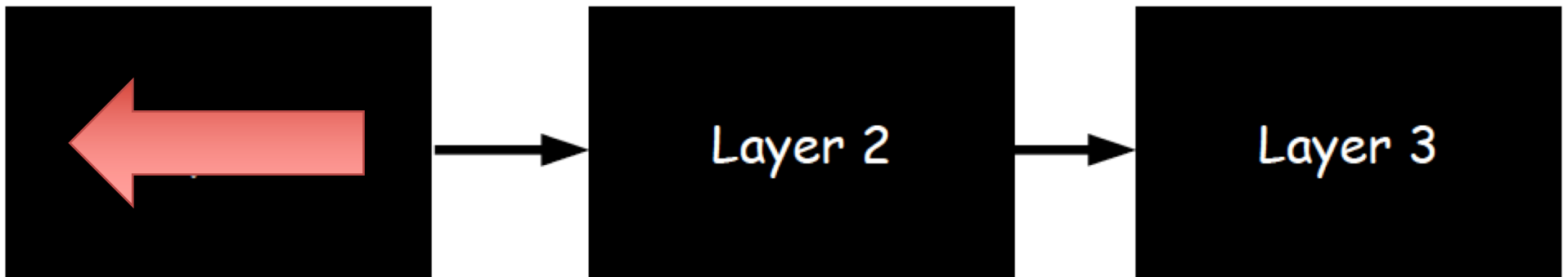- Step 2: Compute gradients wrt parameters    [B-Pass]

# 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]

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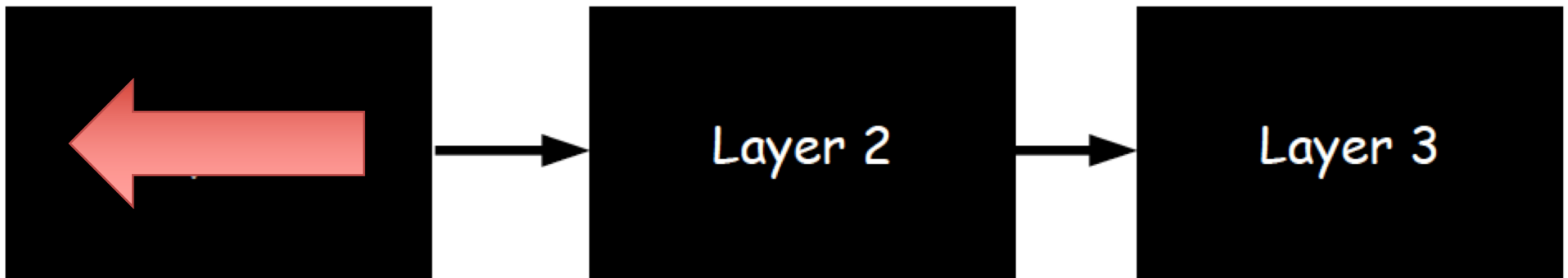
- Step 1: Compute Loss on mini-batch        [F-Pass]
- Step 2: Compute gradients wrt parameters    [B-Pass]
- Step 3: Use gradient to update parameters



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

Slide Credit: Marc'Aurelio Ranzato, Yann LeCun

# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$
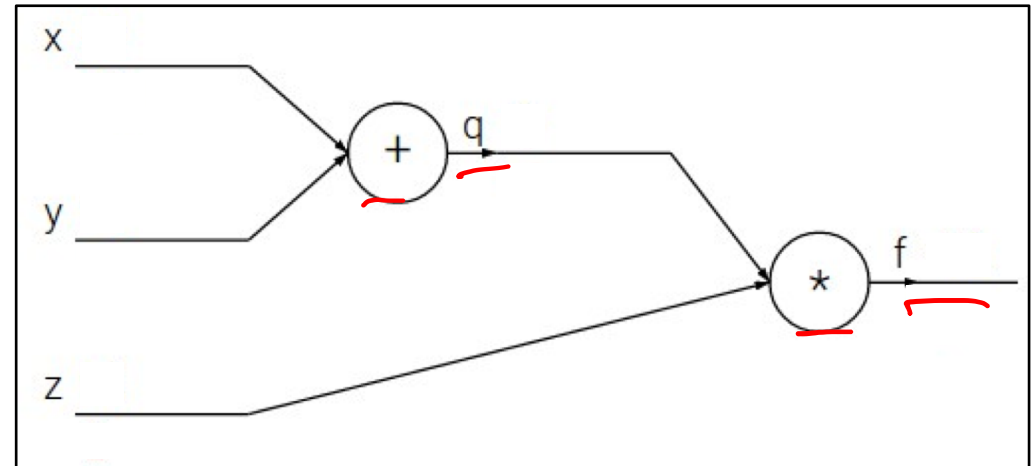
# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4



Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$

# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$

x  -2
y  5
q  3
z  -4
f  -12

# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
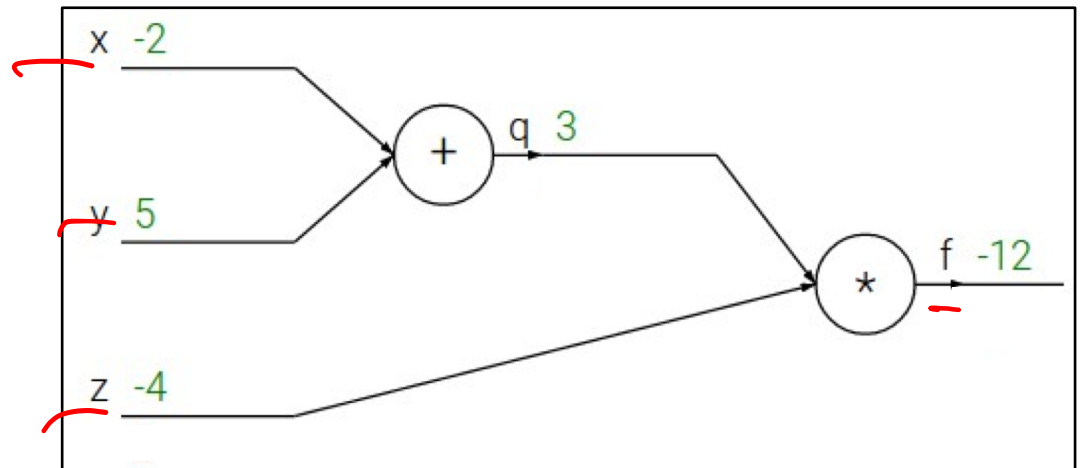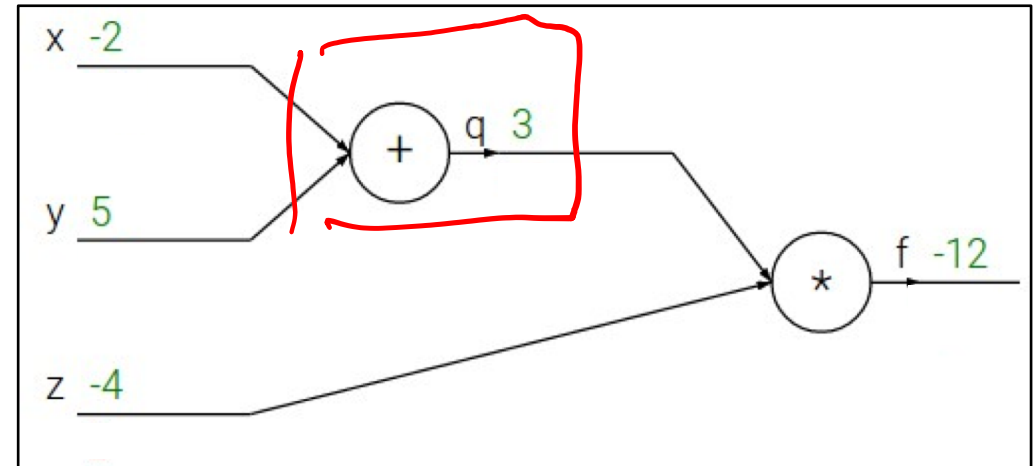
# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$

x  -2

y  5

q  3

z  -4

f  -12

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

# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$

x  -2

q  3

y  5

f  -12
1

z  -4

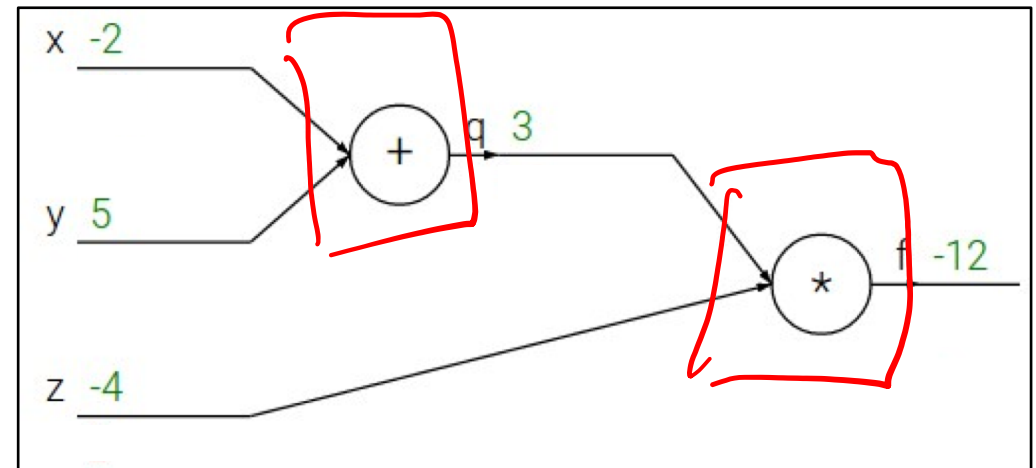$\dfrac{\partial f}{\partial f}$

# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



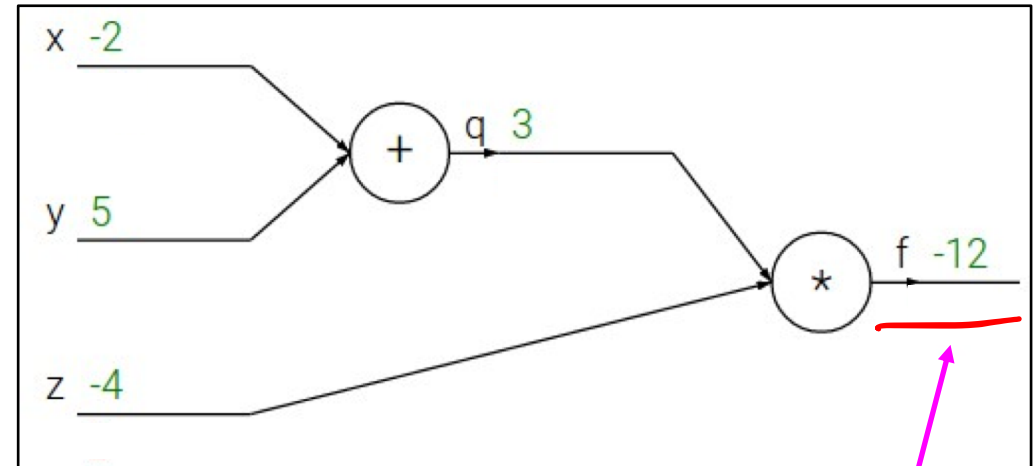$$\frac{\partial f}{\partial z}$$
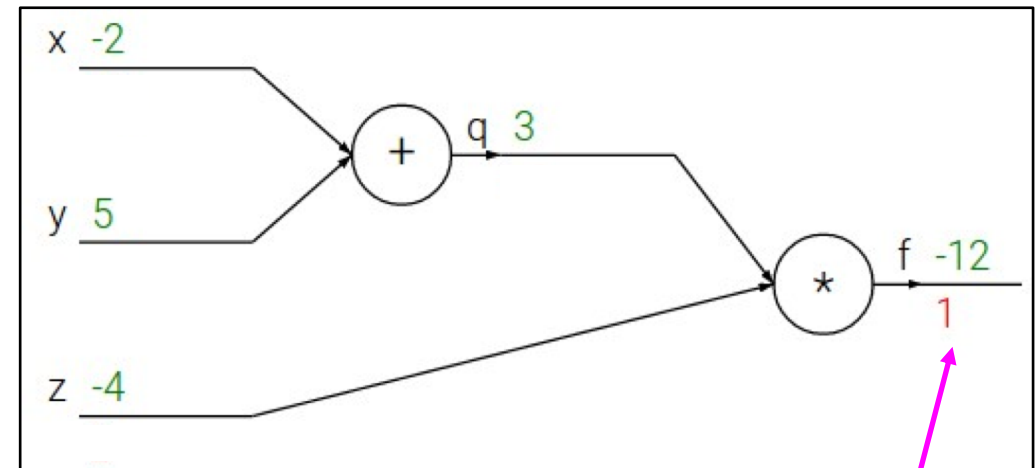
# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$

x  -2
y  5
q  3
z  -4
3
f  -12
1
$$\frac{\partial f}{\partial z}$$

# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



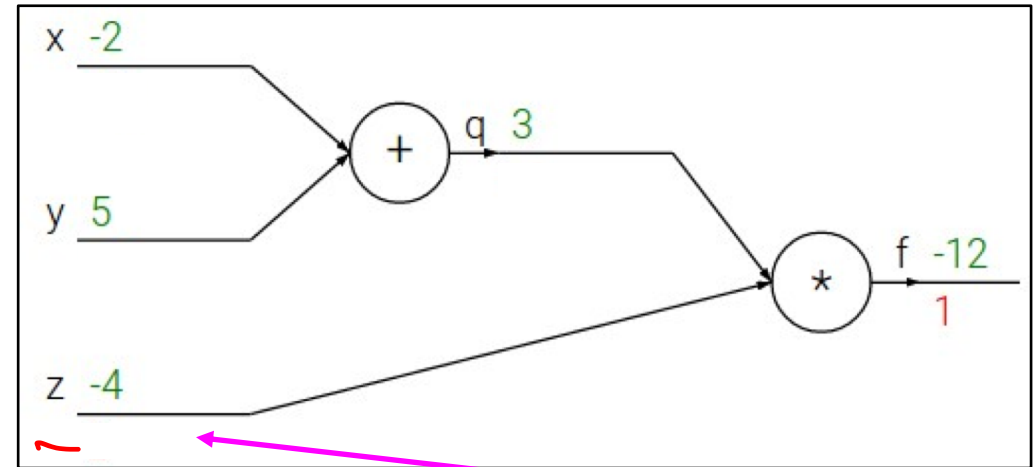$$\frac{\partial f}{\partial q}$$

# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

x  -2

y  5

q  3
-4

z  -4
3

f  -12
1

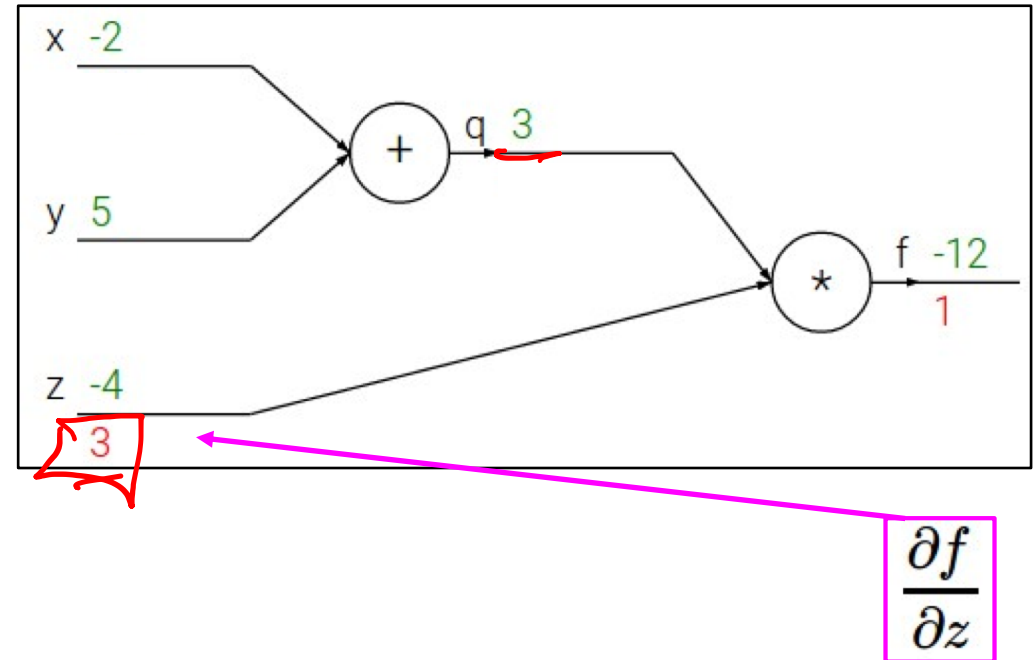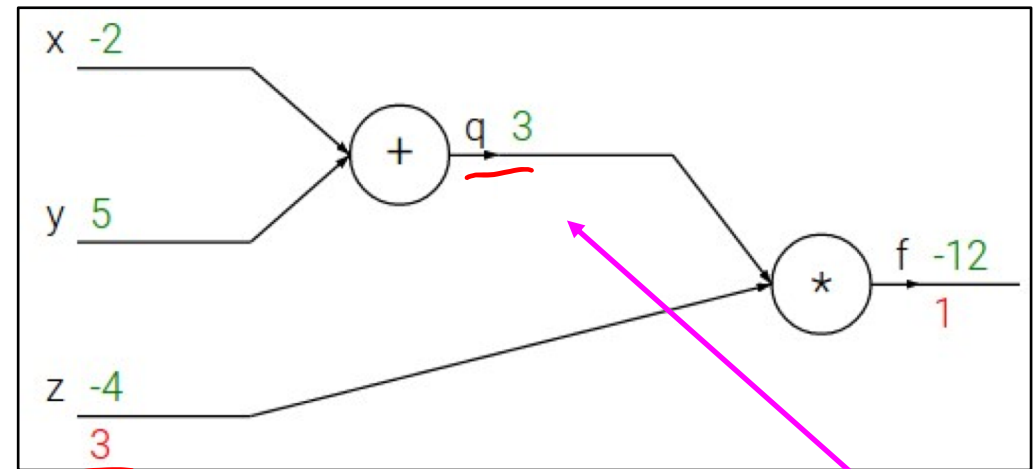$\frac{\partial f}{\partial q}$

# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial y}$$

Chain rule:

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$

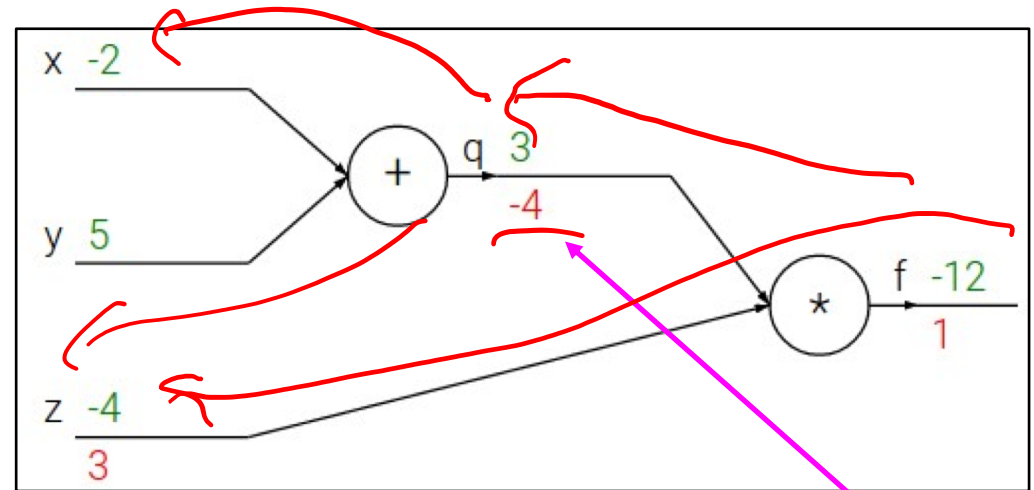Upstream gradient    Local gradient

# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



x  -2

q  3
-4

y  5
-4

z  -4
3

f  -12
1

$\dfrac{\partial f}{\partial y}$

Chain rule:

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$

Upstream gradient    Local gradient

# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

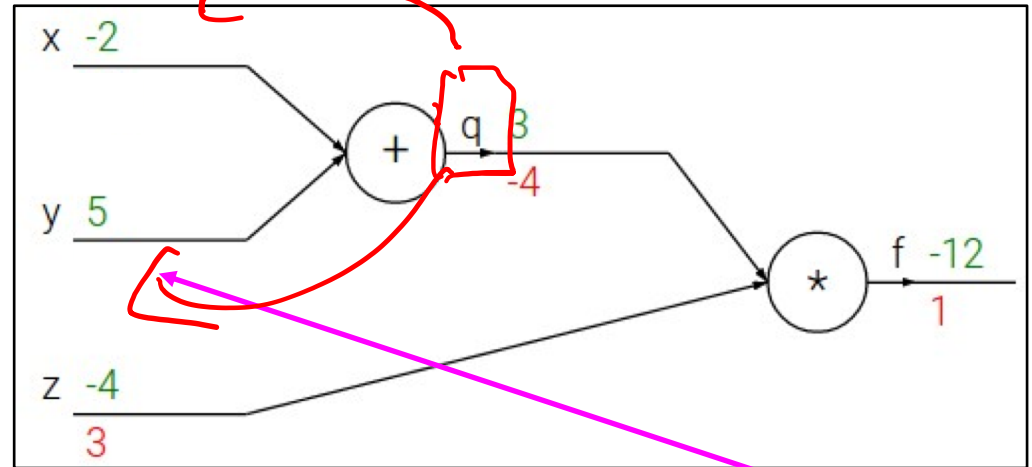Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



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

Chain rule:

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

Upstream gradient    Local gradient

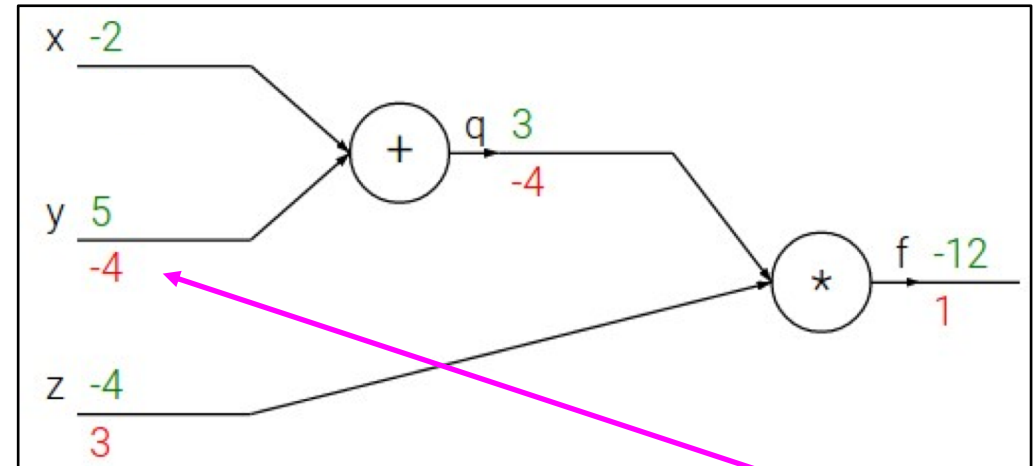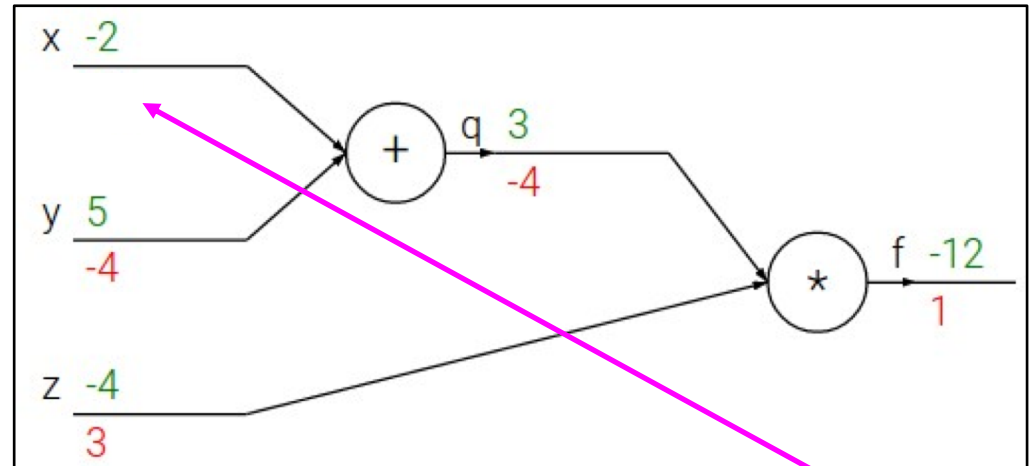# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$q = x + y \qquad \dfrac{\partial q}{\partial x} = 1, \dfrac{\partial q}{\partial y} = 1$

$f = qz \qquad \dfrac{\partial f}{\partial q} = z, \dfrac{\partial f}{\partial z} = q$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



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

Chain rule:

$$\dfrac{\partial f}{\partial x} = \dfrac{\partial f}{\partial q} \dfrac{\partial q}{\partial x}$$

Upstream gradient    Local gradient

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

w0 2.00

x0 -1.00

-2.00

w1 -3.00

x1 -2.00

6.00

4.00

1.00

-1.00

0.37

1.37

0.73

*-1

exp

+1

1/x

w2 -3.00

# Another example:

$$f(w,x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \qquad \Bigg| \qquad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \Bigg| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \quad \Bigg| \quad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \quad \Bigg| \quad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



Upstream gradient    Local gradient

$$(1.00)(\frac{-1}{1.37^2}) = -0.53$$

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



w0 2.00

x0 -1.00

w1 -3.00

x1 -2.00

w2 -3.00

-2.00

4.00

6.00

1.00

-1.00

0.37 / -0.53

1.37 / -0.53

0.73 / 1.00

Upstream gradient    Local gradient

$$(-0.53)(1) = -0.53$$

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$
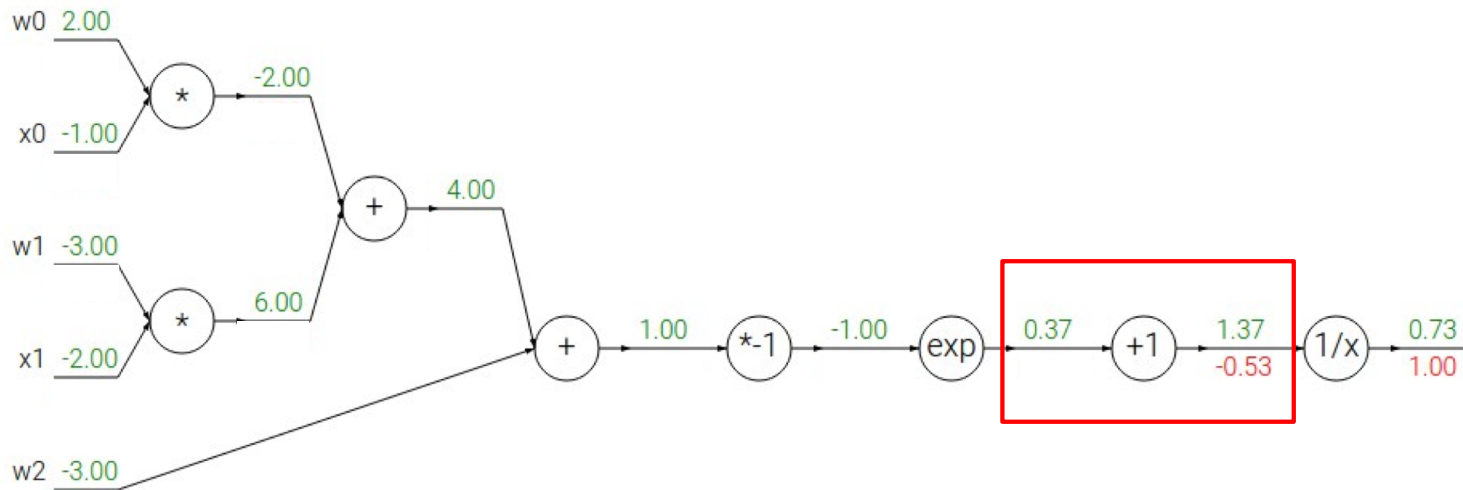
# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

w0 2.00

x0 -1.00

-2.00

w1 -3.00

x1 -2.00

6.00

w2 -3.00

4.00

1.00

Upstream gradient

Local gradient

$$(-0.53)(e^{-1}) = -0.20$$

-1.00    0.37
-0.20    -0.53

1.37
-0.53

0.73
1.00

$f(x) = e^x$ $\rightarrow$ $\dfrac{df}{dx} = e^x$

$f_a(x) = ax$ $\rightarrow$ $\dfrac{df}{dx} = a$
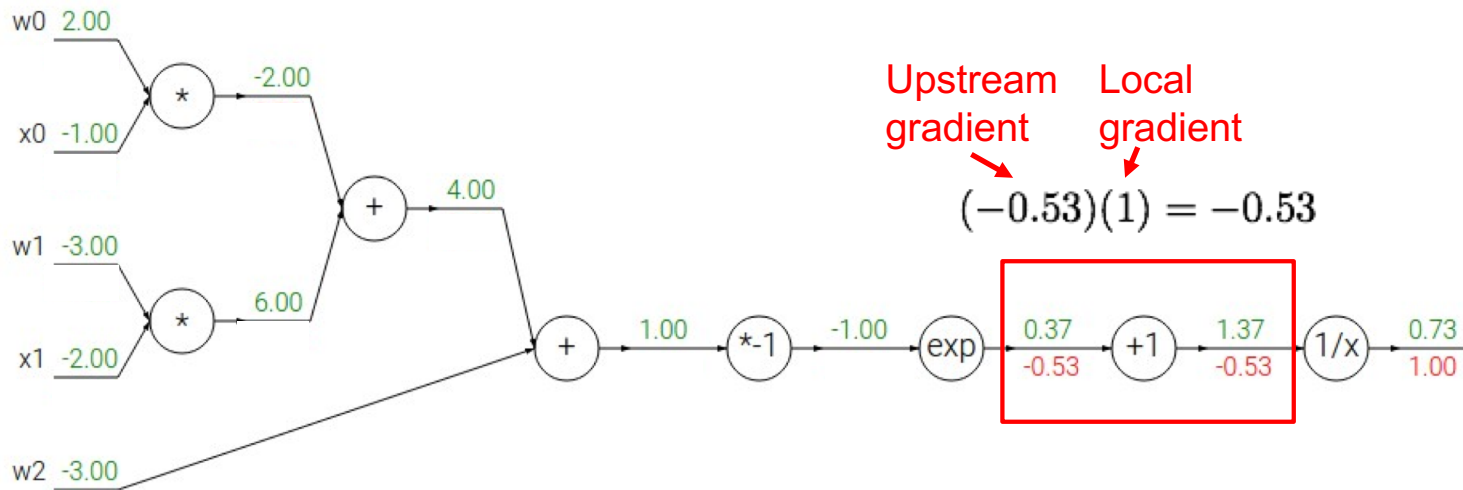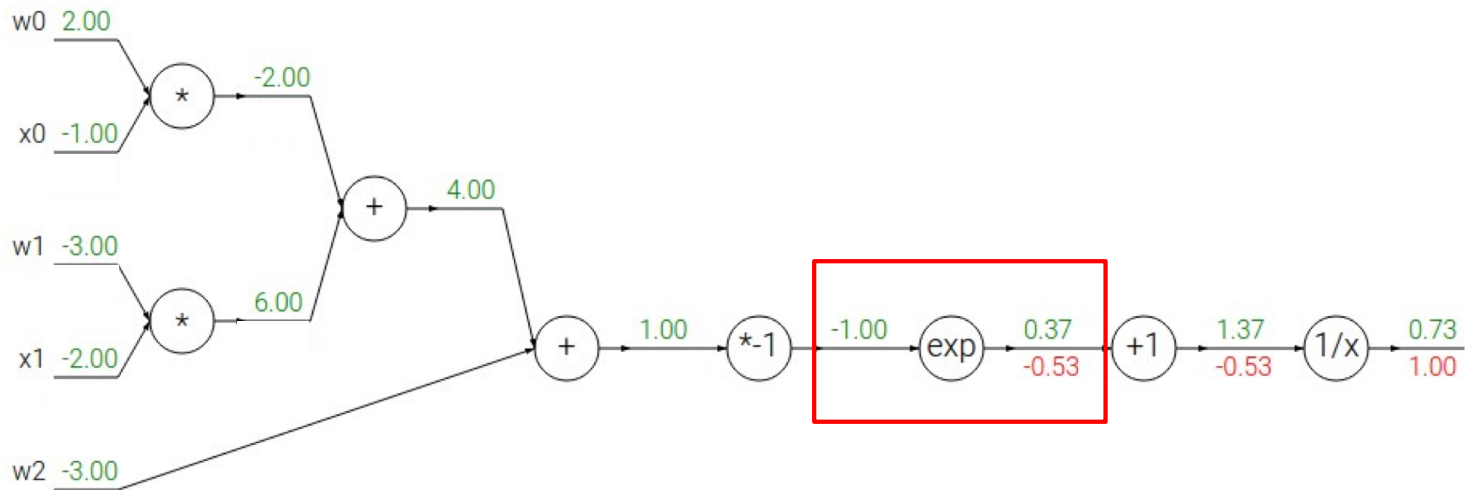
$f(x) = \dfrac{1}{x}$ $\rightarrow$ $\dfrac{df}{dx} = -1/x^2$

$f_c(x) = c + x$ $\rightarrow$ $\dfrac{df}{dx} = 1$

# Another example:

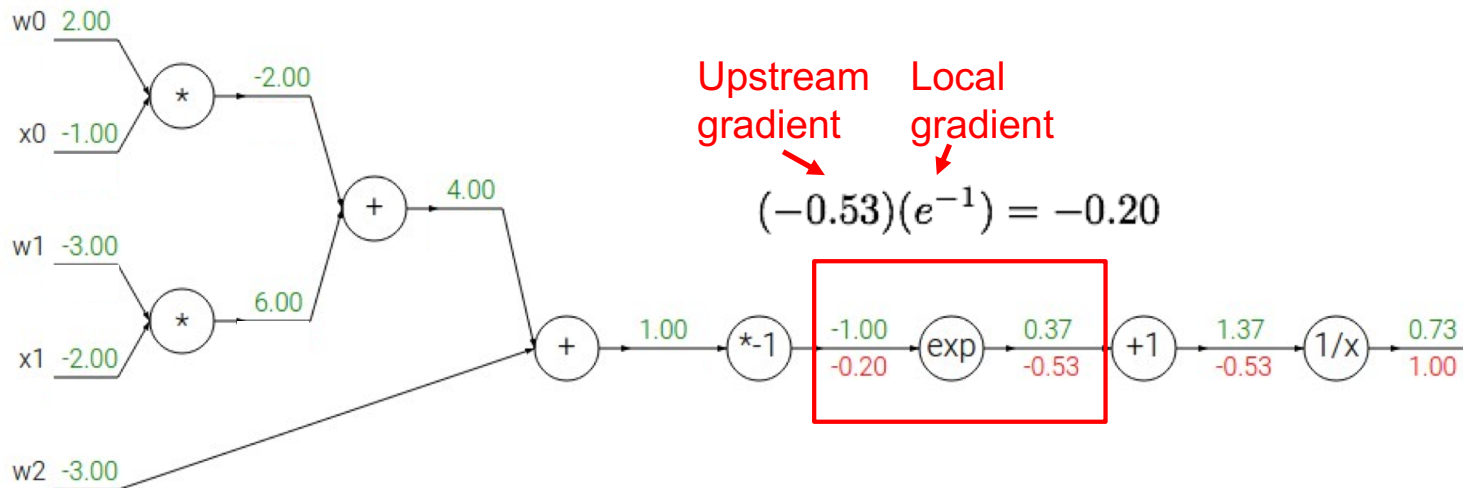$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Another example:

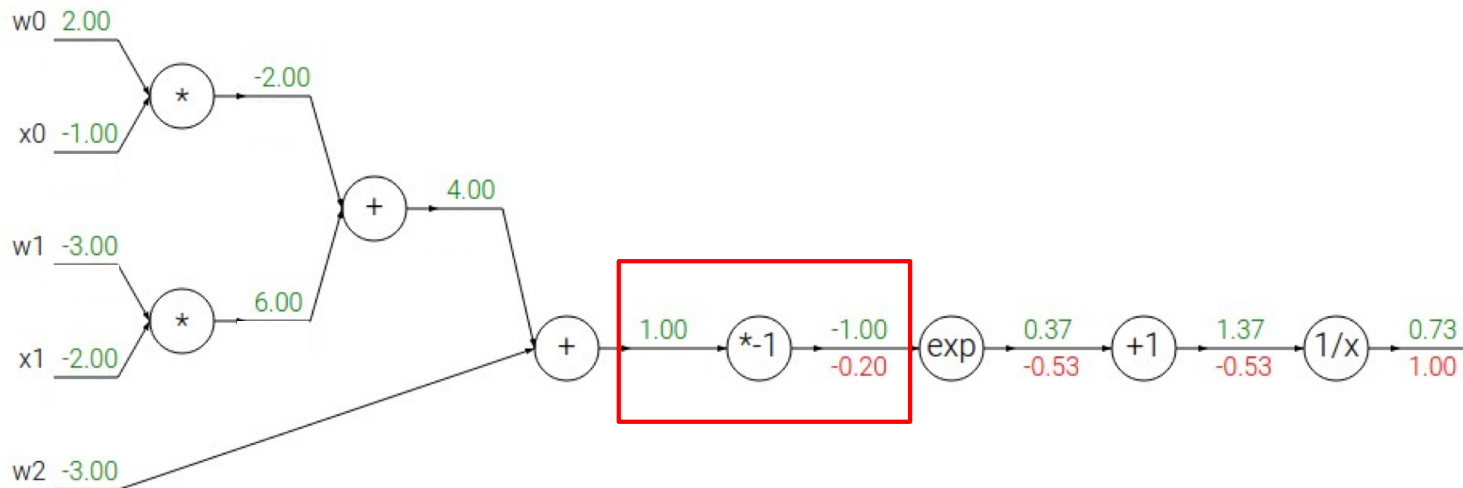$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



Upstream gradient    Local gradient

$$(-0.20)(-1) = 0.20$$

| $f(x) = e^x$ | $\rightarrow$ | $\frac{df}{dx} = e^x$ | $f(x) = \frac{1}{x}$ | $\rightarrow$ | $\frac{df}{dx} = -1/x^2$ |
| $f_a(x) = ax$ | $\rightarrow$ | $\frac{df}{dx} = a$ | $f_c(x) = c + x$ | $\rightarrow$ | $\frac{df}{dx} = 1$ |

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$


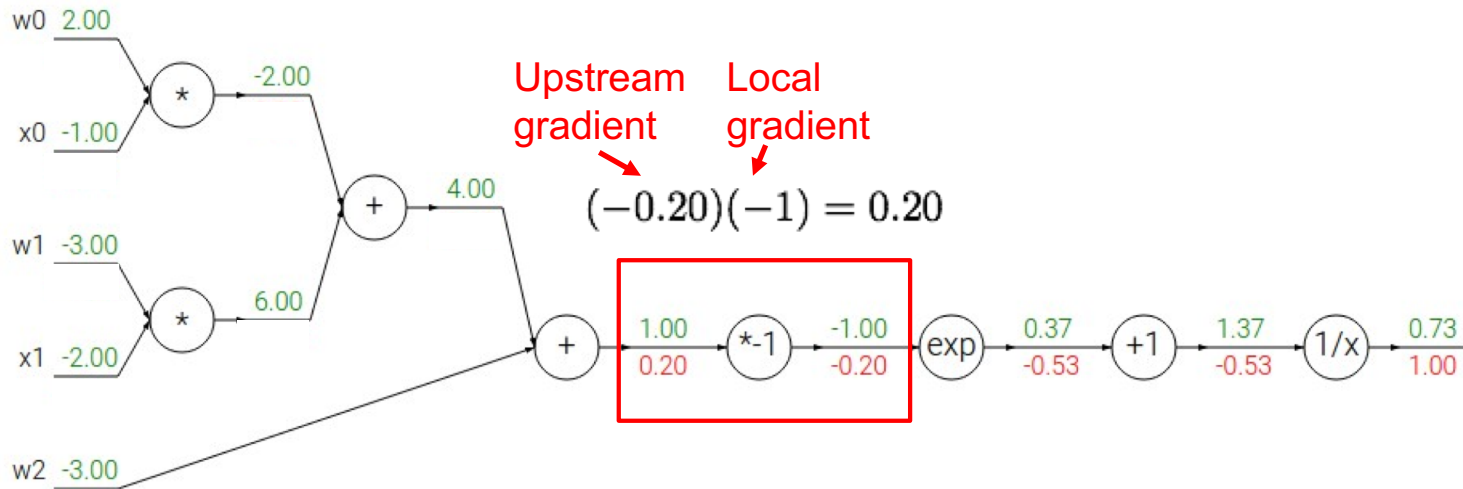
$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \qquad \bigg| \qquad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

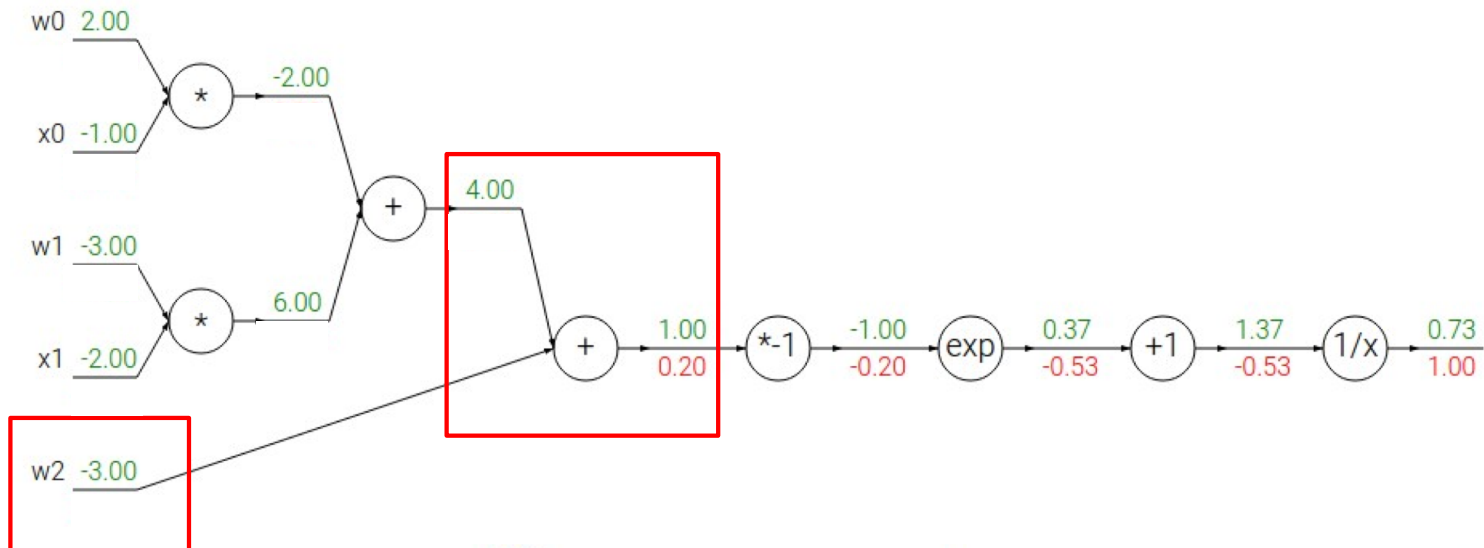$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \bigg| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



[upstream gradient] x [local gradient]
[0.2] x [1] = 0.2
[0.2] x [1] = 0.2  (both inputs!)

$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \qquad \Bigg| \qquad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$

$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \Bigg| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \qquad \Bigg| \qquad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \Bigg| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



[upstream gradient] x [local gradient]
x0: [0.2] x [2] = 0.4
w0: [0.2] x [-1] = -0.2

$$f(x) = e^x \quad \rightarrow \quad \frac{df}{dx} = e^x \quad \Big| \quad f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \quad \rightarrow \quad \frac{df}{dx} = a \quad \Big| \quad f_c(x) = c + x \quad \rightarrow \quad \frac{df}{dx} = 1$$
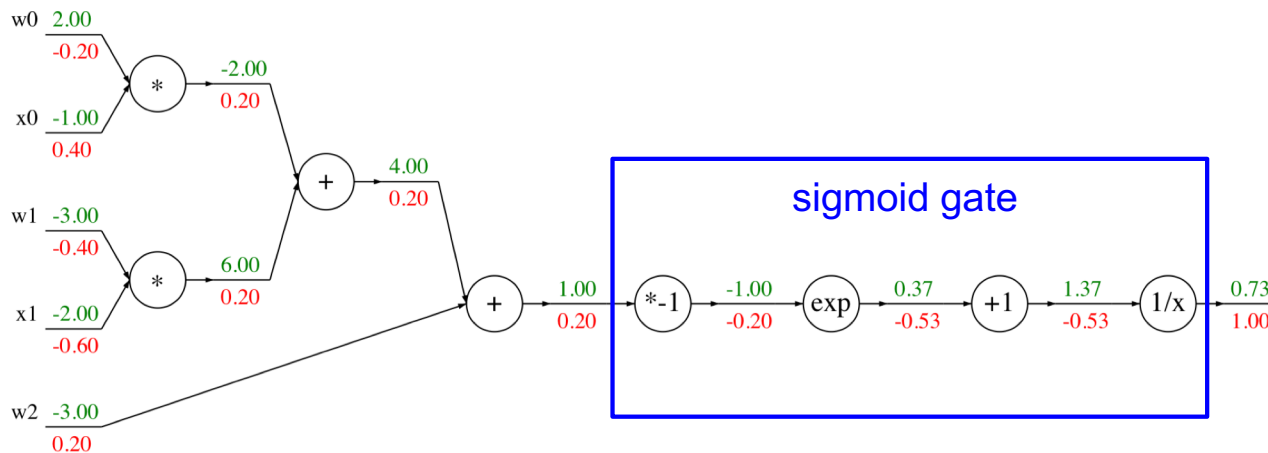
$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

Computational graph representation may not be unique. Choose one where local gradients at each node can be easily expressed!

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

sigmoid function

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}}\right)\left(\frac{1}{1 + e^{-x}}\right) = (1 - \sigma(x))\sigma(x)$$

sigmoid gate

w0  2.00
    -0.20

x0  -1.00
    0.40

w1  -3.00
    -0.40

x1  -2.00
    -0.60

w2  -3.00
    0.20

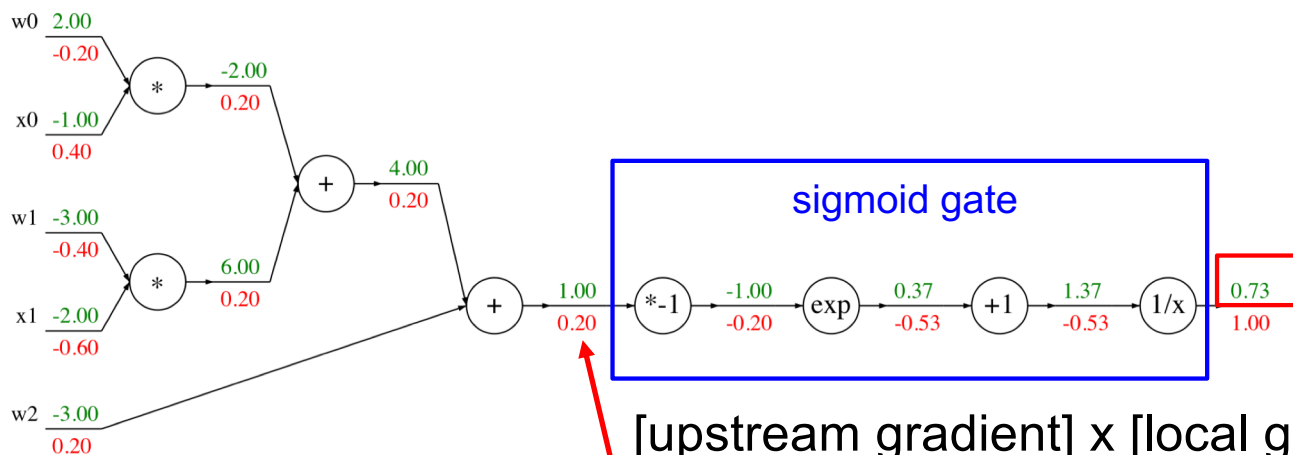Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

$$f(w,x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Computational graph representation may not be unique. Choose one where local gradients at each node can be easily expressed!

$$\boxed{\sigma(x) = \frac{1}{1 + e^{-x}}} \qquad \frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}}\right)\left(\frac{1}{1 + e^{-x}}\right) = (1 - \sigma(x))\,\sigma(x)$$

sigmoid function

sigmoid gate

[upstream gradient] x [local gradient]
[1.00] x [(1 - 0.73) (0.73)]= 0.2