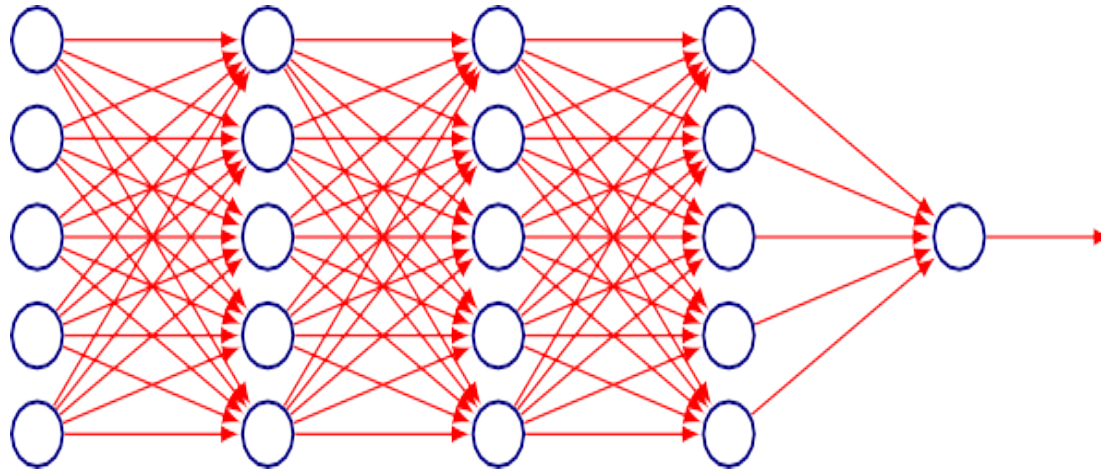# CS60010: Deep Learning

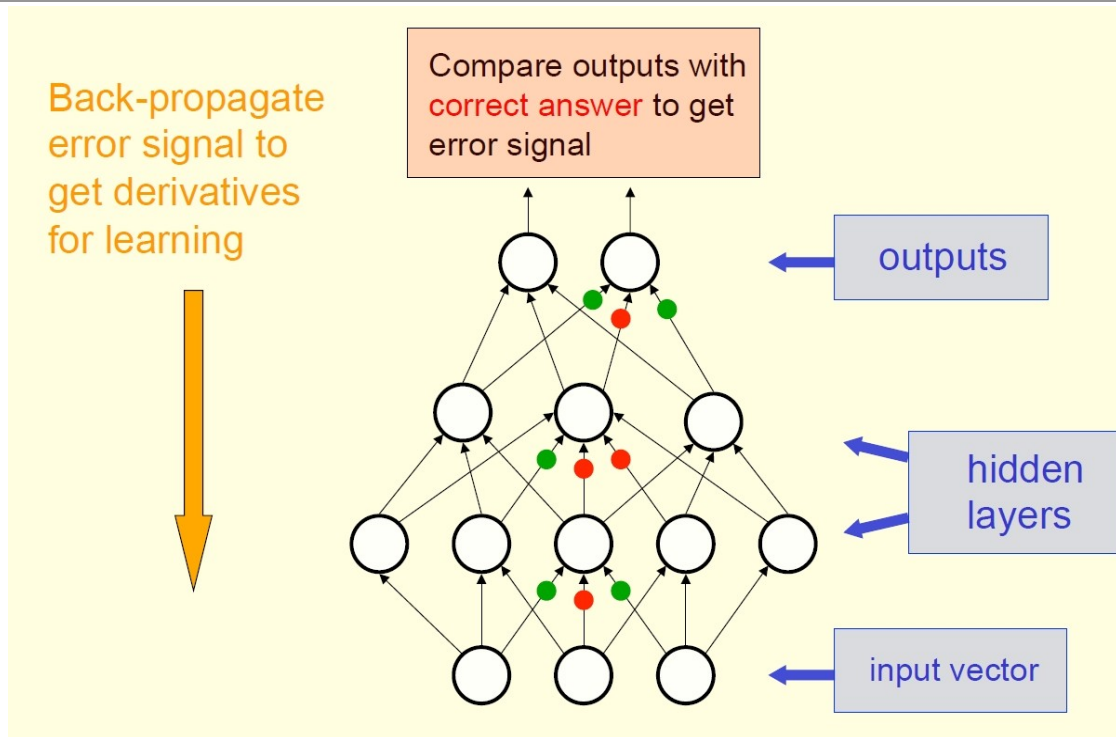**Sudeshna Sarkar**

Spring 2018

16 Jan 2018

# BACKPROPAGATION: INTRODUCTION

# How do we learn weights?



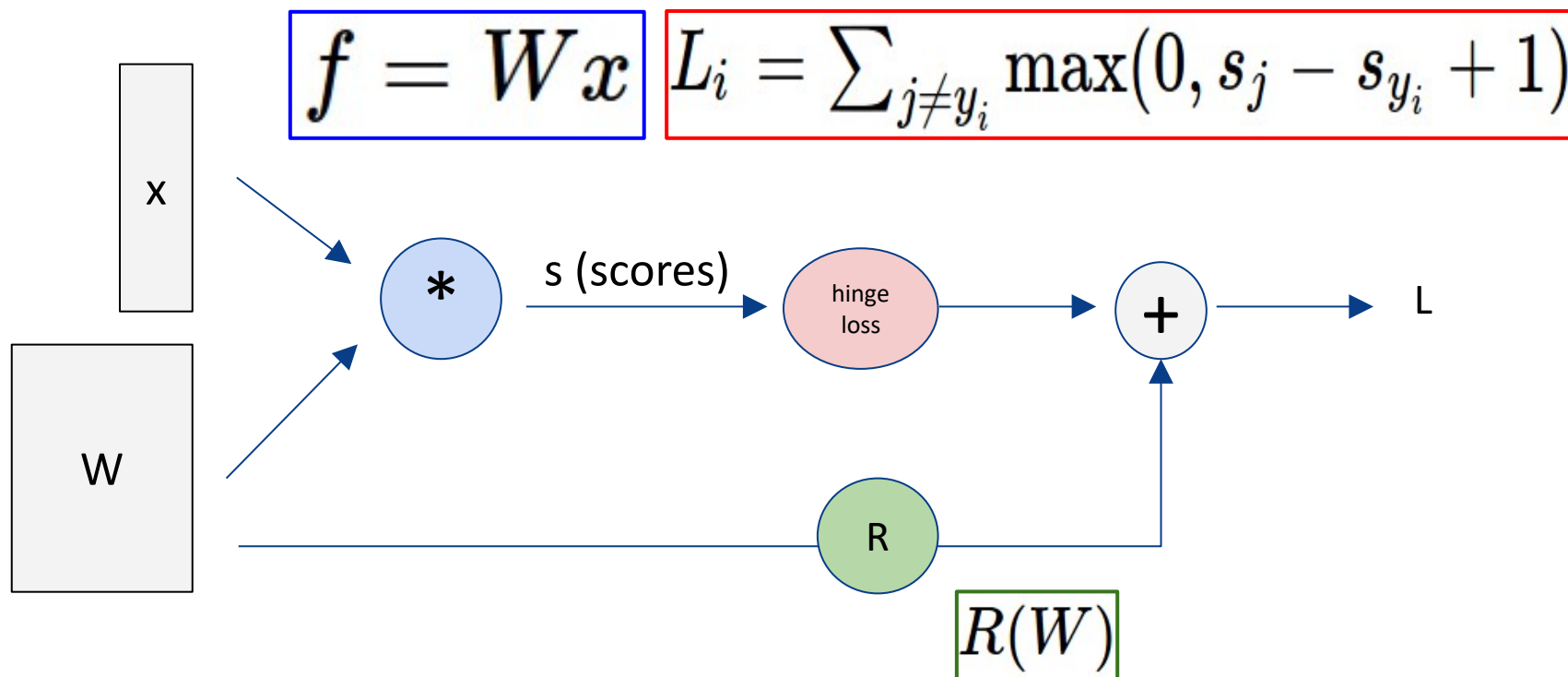- Perturn the weights and check

# Backpropagation



- **Feedforward Propagation**: Accept input $x$, pass through intermediate stages and obtain output $\hat{y}$
- **During Training**: Use $\hat{y}$ to compute a scalar cost $J(\theta)$
- Backpropagation allows information to flow backwards from cost to compute the gradient

# Backpropagation

- From the training data we don't know what the hidden units should do

- But, we can compute how fast the error changes as we change a hidden activity

- Use error derivatives w.r.t hidden activities

- Each hidden unit can affect many output units and have separate effects on error – combine these effects

- Can compute error derivatives for hidden units efficiently (and once we have error derivatives for hidden activities, easy to get error derivatives for weights going in)
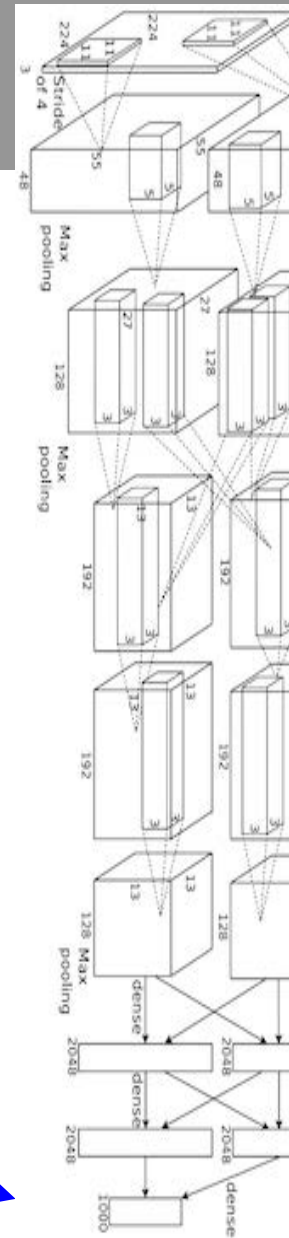
# Computational Graph

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

x

W

\*

s (scores)

hinge loss

+

L

R

$$R(W)$$
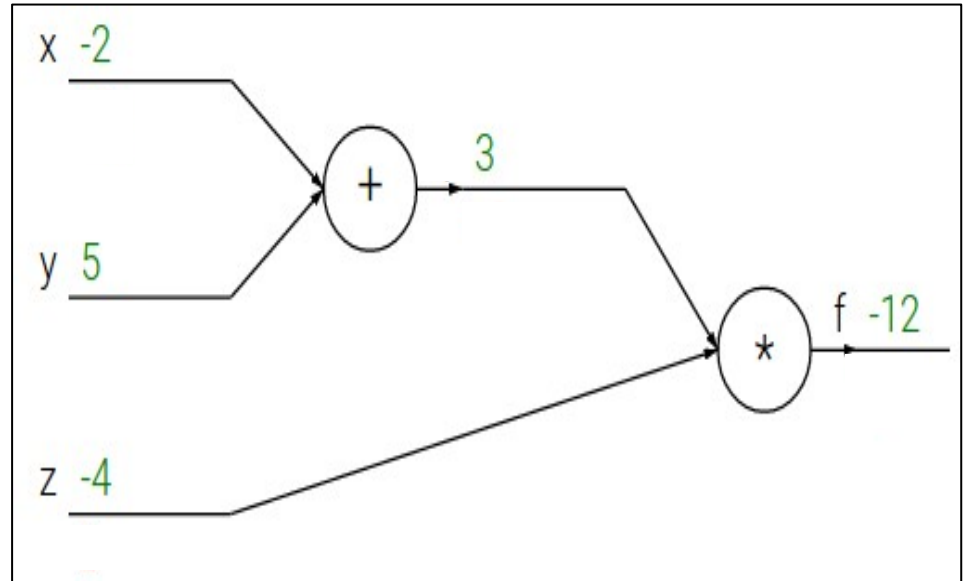
# Convolutional Network (AlexNet)



input image

weights

loss

# Differentiating a Computation Graph

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

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

We can write

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

Where $h(x, y) = x + y$, and $g(a, b) = a * b$

By the chain rule, $\dfrac{df}{dx} = \dfrac{dg}{dh}\dfrac{dh}{dx}$ and $\dfrac{df}{dy} = \dfrac{dg}{dh}\dfrac{dh}{dy}$
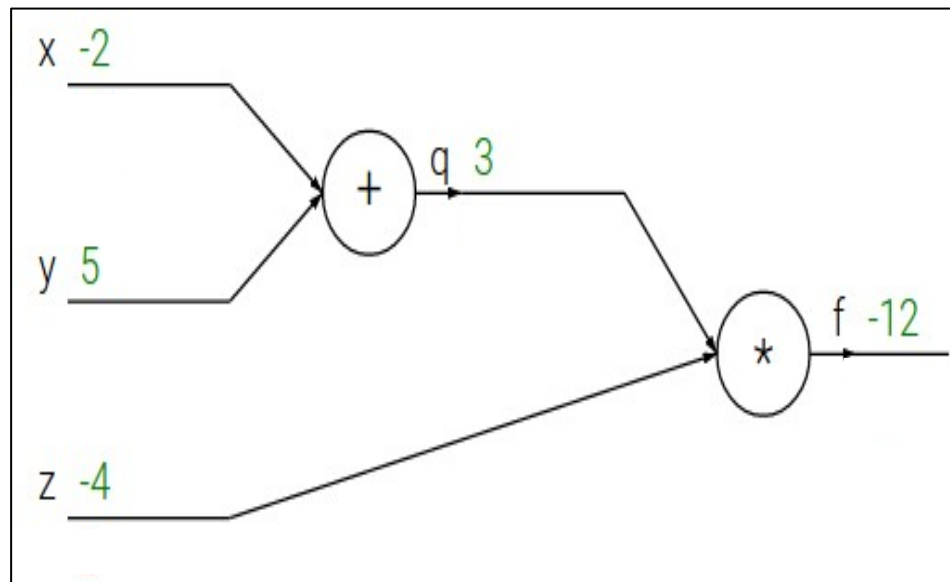
# Differentiating a Computation Graph

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

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

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

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



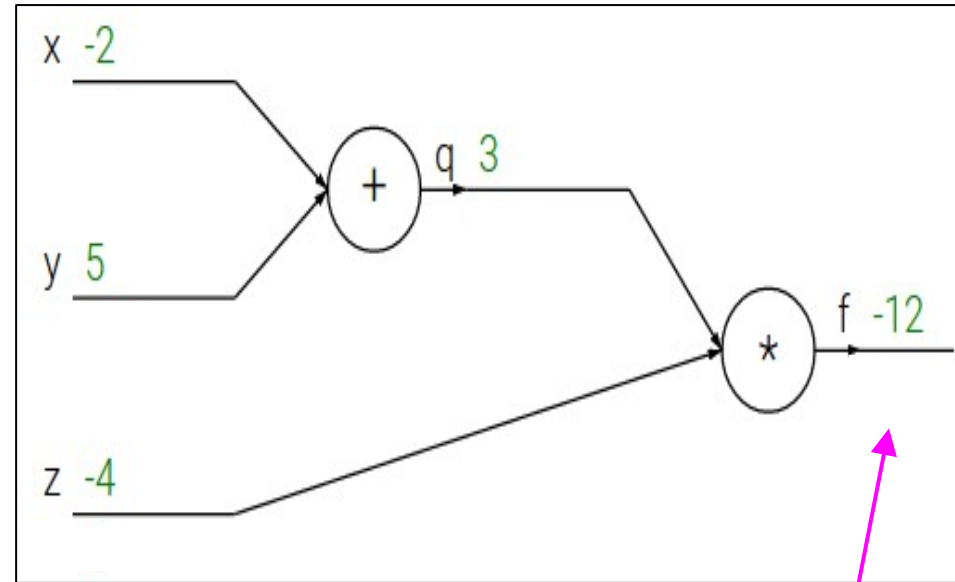Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$

# Differentiating a Computation Graph

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

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

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

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

x -2

y 5

+   q 3

z -4

*    f -12

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

Want:   $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
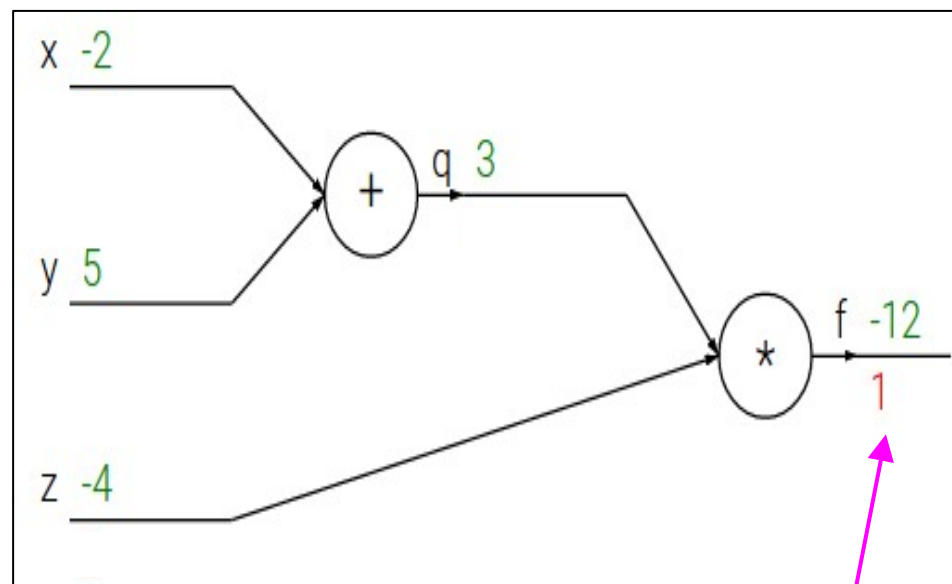
# Differentiating a Computation Graph

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

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

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

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



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

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
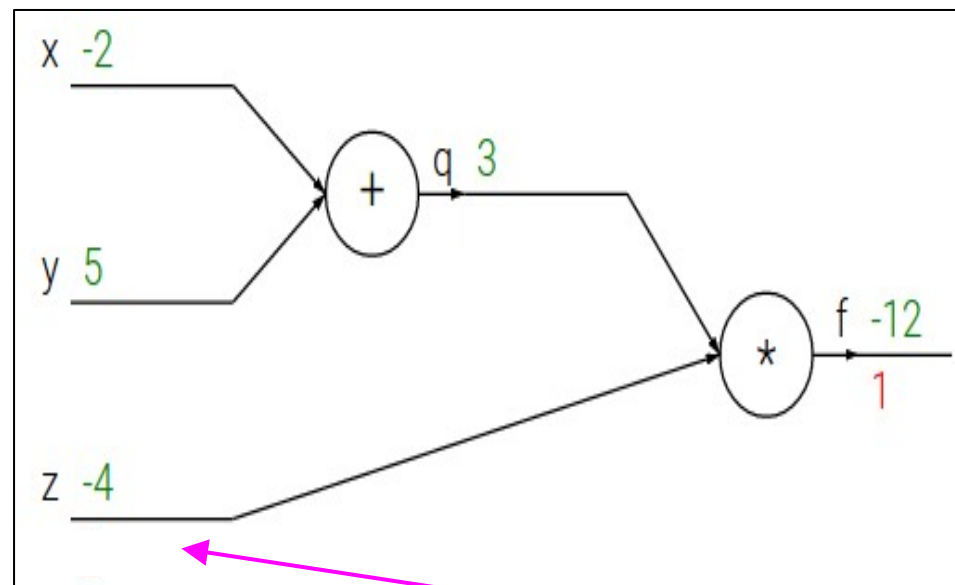
# Differentiating a Computation Graph

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

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

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

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



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

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
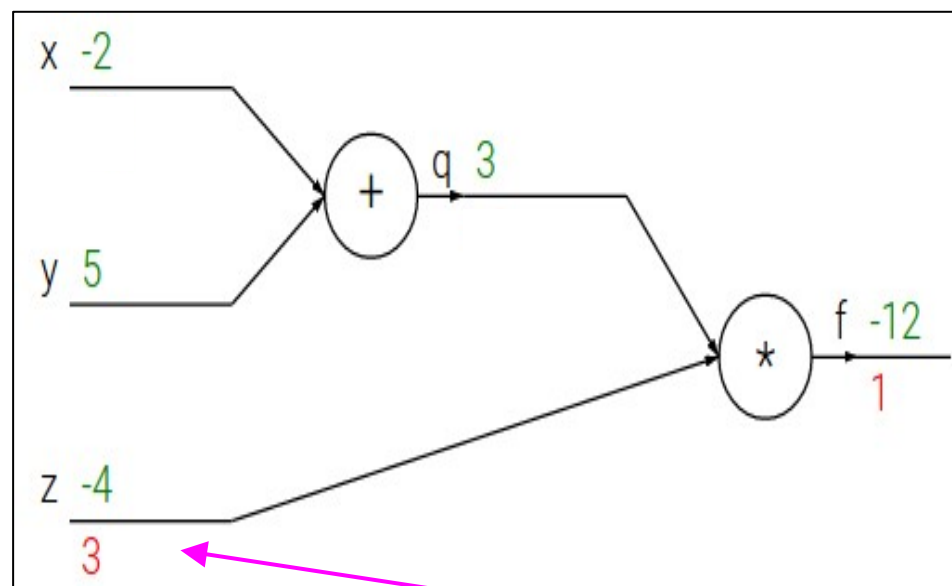
# Differentiating a Computation Graph

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

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

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

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



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

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

# Differentiating a Computation Graph

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

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

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

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



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

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
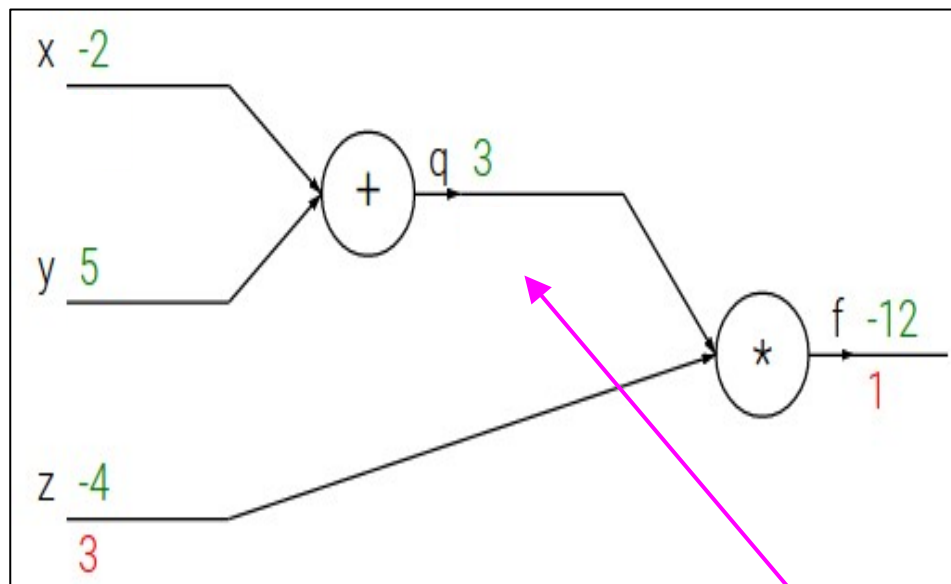
# Differentiating a Computation Graph

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

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

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

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



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

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
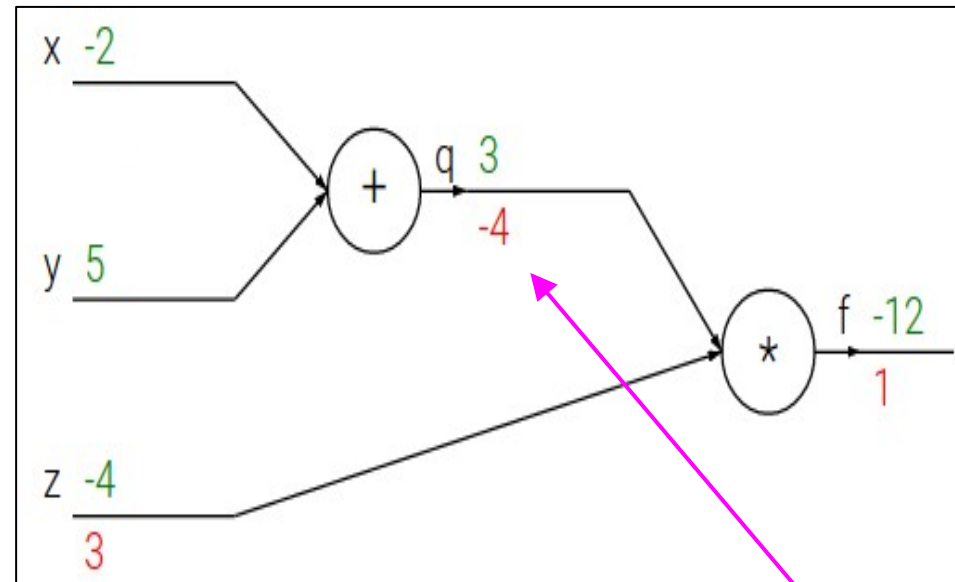
# Differentiating a Computation Graph

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

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

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

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



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

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

Based on cs231n by Fei-Fei Li & Andrej Karpathy & Justin Johnson
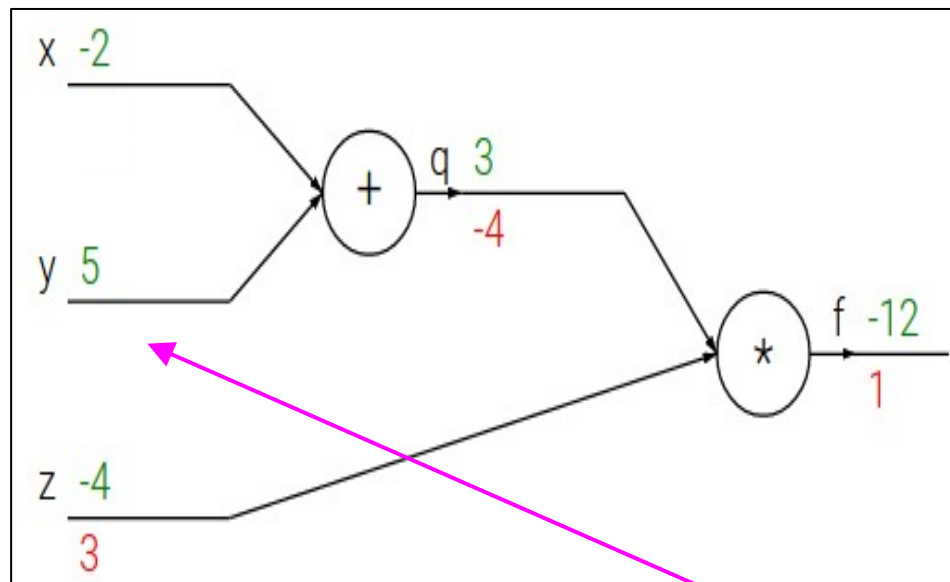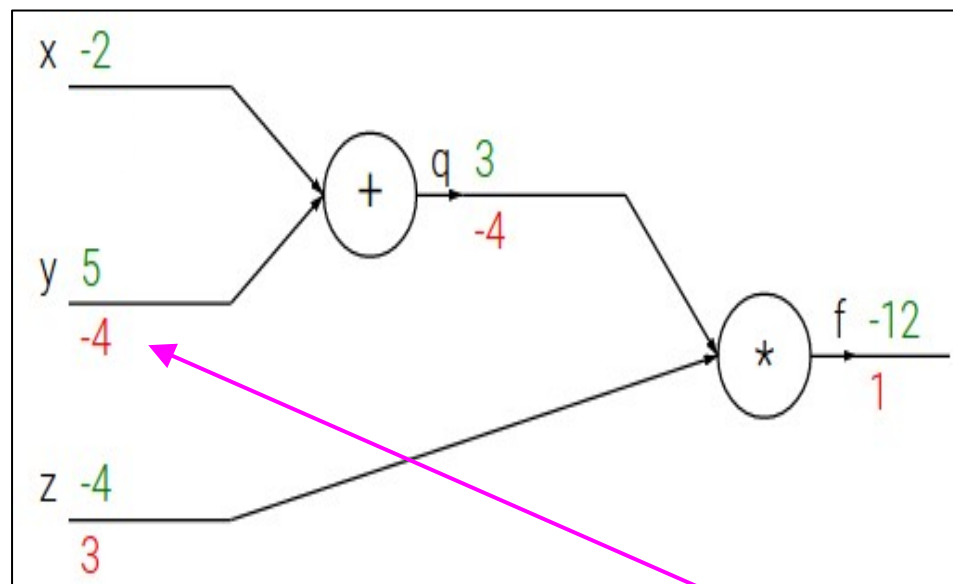
# Differentiating a Computation Graph

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

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

Chain rule:

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

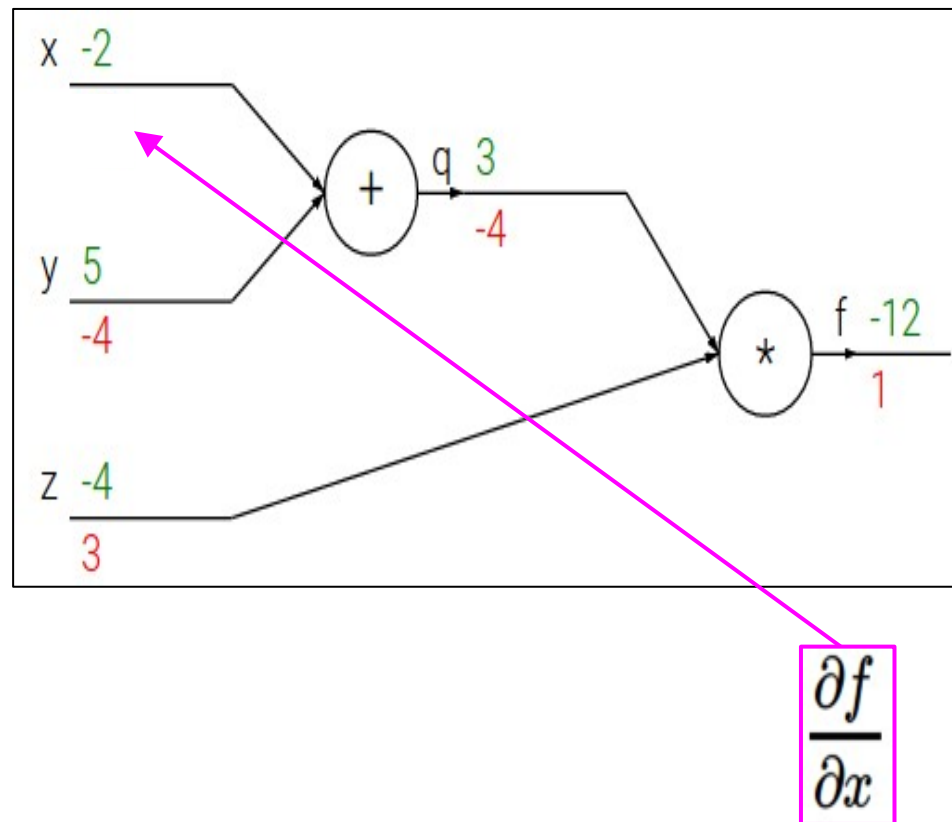$$\frac{\partial f}{\partial y}$$

# Differentiating a Computation Graph

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

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

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

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



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

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$
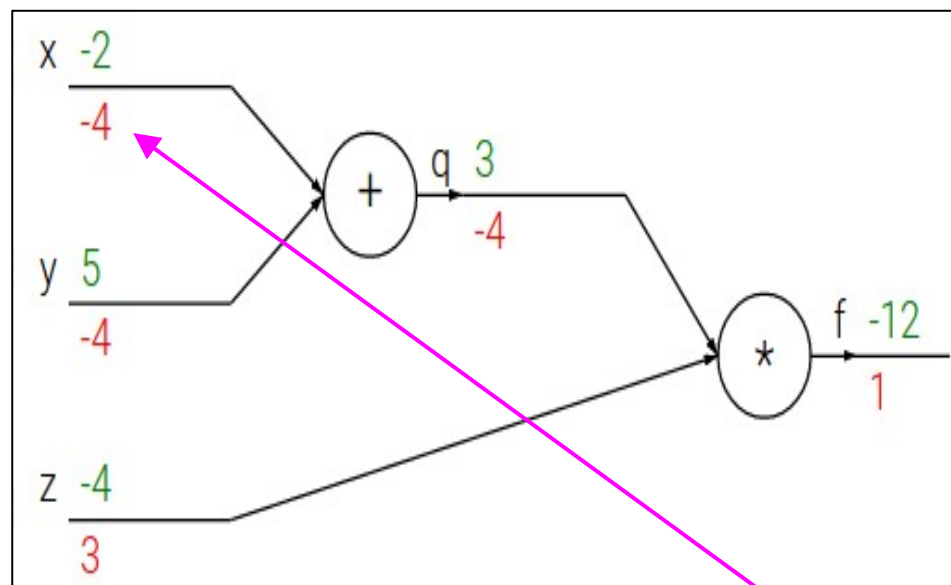
# Differentiating a Computation Graph

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

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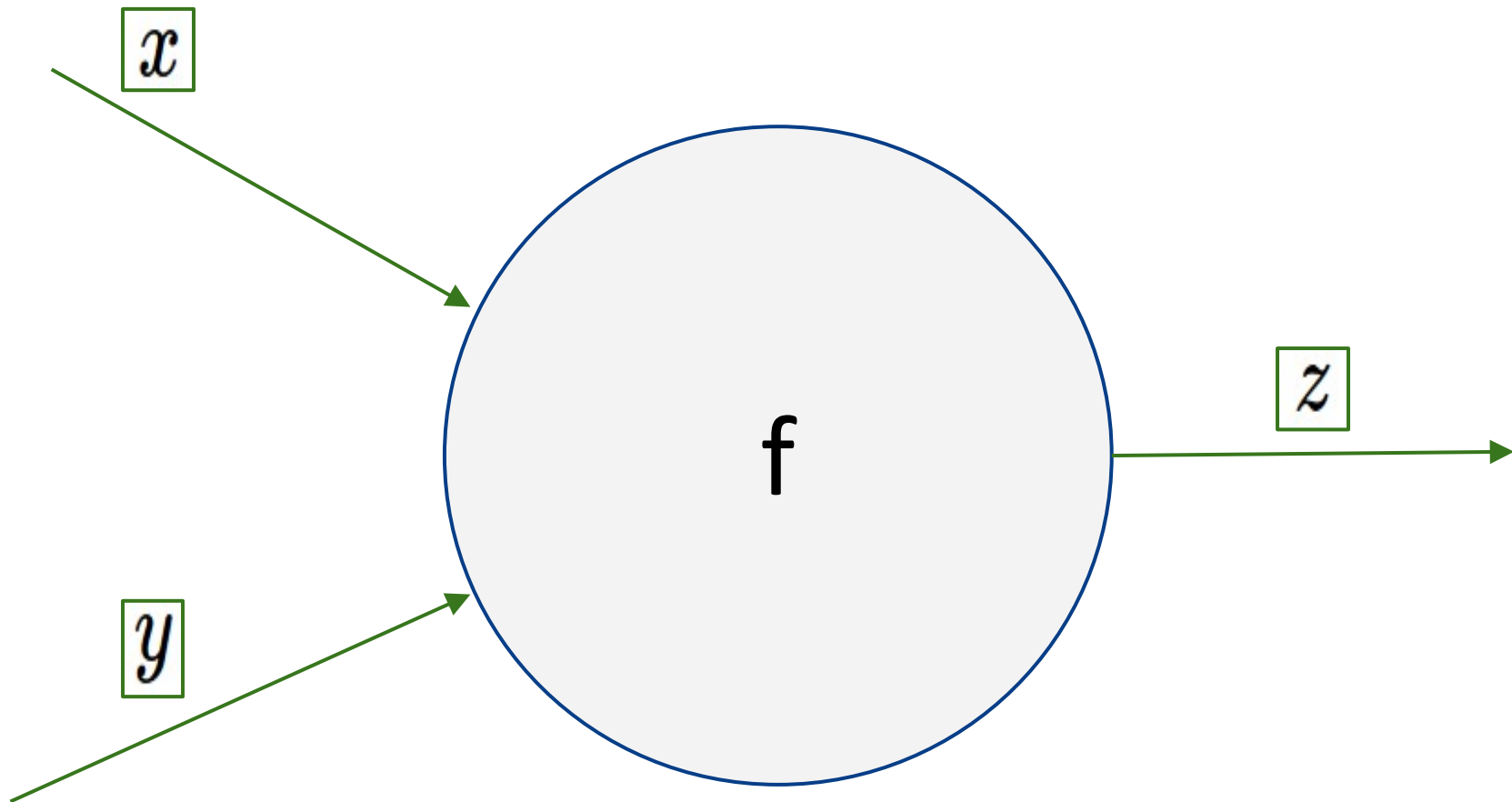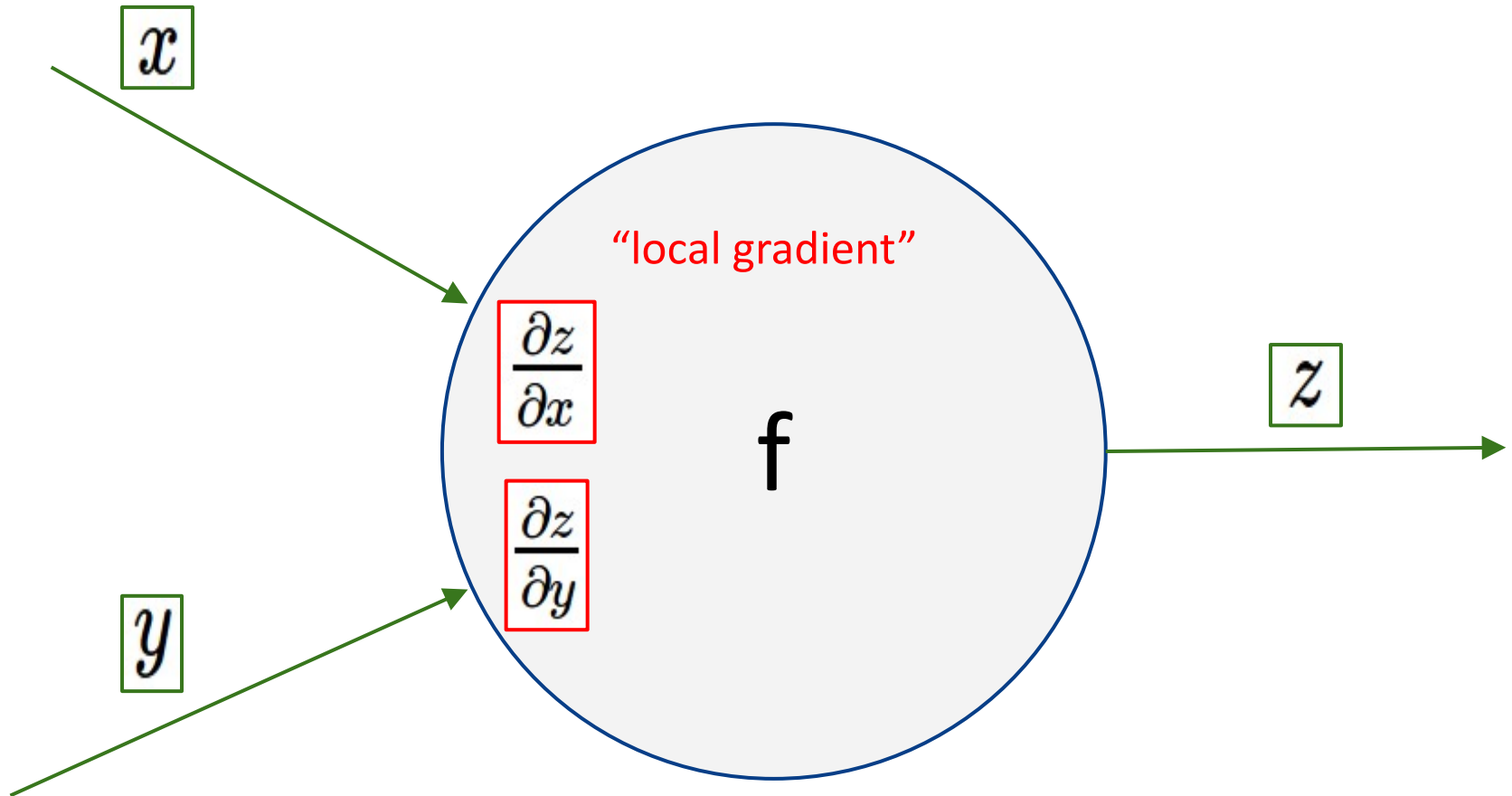


Chain rule:

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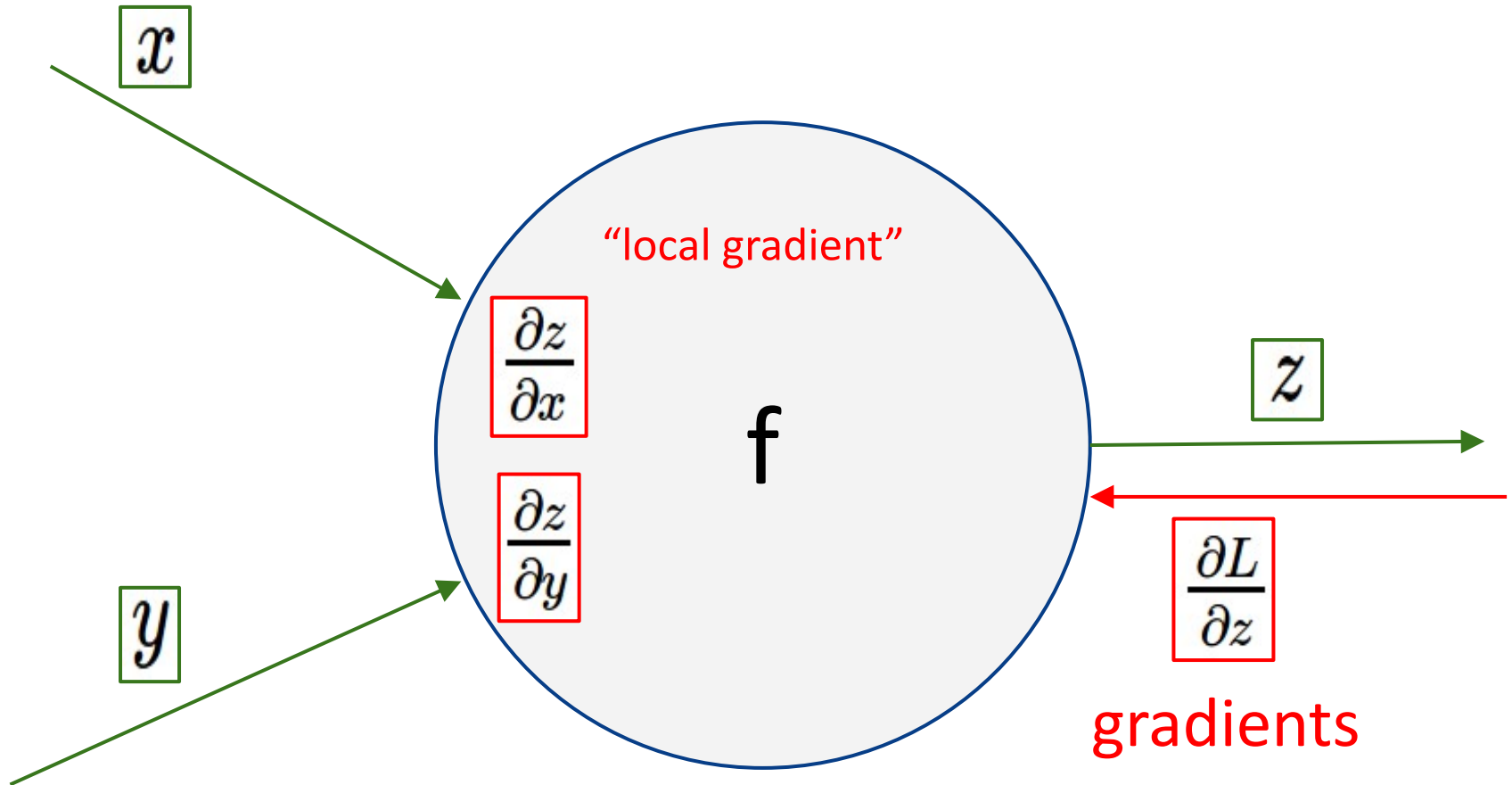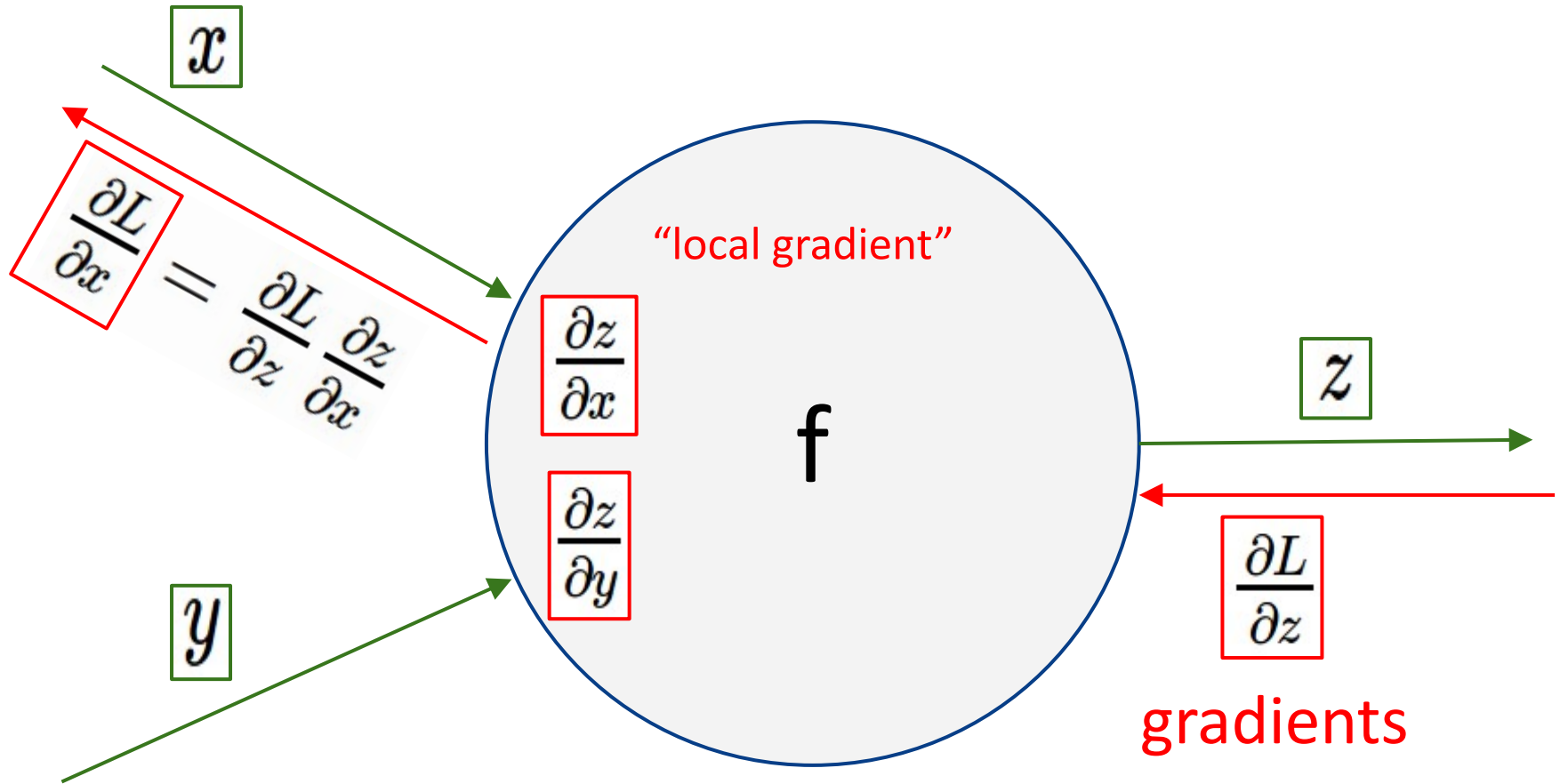
$$\frac{\partial f}{\partial x}$$

# activations

# activations

"local gradient"

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

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

f

$x$

$y$

$z$

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

gradients

$$x$$

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

"local gradient"

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

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

$$f$$

$$z$$

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

gradients

$$y$$

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

# Another backprop example:



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

# Another example:



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

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

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

# Another example:



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

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

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

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

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

# Another example:



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

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

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

Based on cs231n by Fei-Fei Li & Andrej Karpathy & Justin Johnson

# 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:



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

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

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

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

# Another example:



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

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

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

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

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

# Another example:



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

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

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

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

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

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

# 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:



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

(-1) * (-0.20) = 0.20

| | | |
|---|---|---|
| $f(x) = e^x$ | $\rightarrow$ | $\dfrac{df}{dx} = e^x$ |
| $f_a(x) = ax$ | $\rightarrow$ | $\dfrac{df}{dx} = a$ |

| | | |
|---|---|---|
| $f(x) = \dfrac{1}{x}$ | $\rightarrow$ | $\dfrac{df}{dx} = -1/x^2$ |
| $f_c(x) = c + x$ | $\rightarrow$ | $\dfrac{df}{dx} = 1$ |

# 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:

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

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

w0 2.00
x0 -1.00
w1 -3.00
x1 -2.00
w2 -3.00
  0.20

-2.00
6.00
4.00
0.20
1.00
0.20
-1.00
-0.20
0.37
-0.53
1.37
-0.53
0.73
1.00

\*-1   exp   +1   1/x

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

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

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

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

# 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:



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

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

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

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

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

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

sigmoid function

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



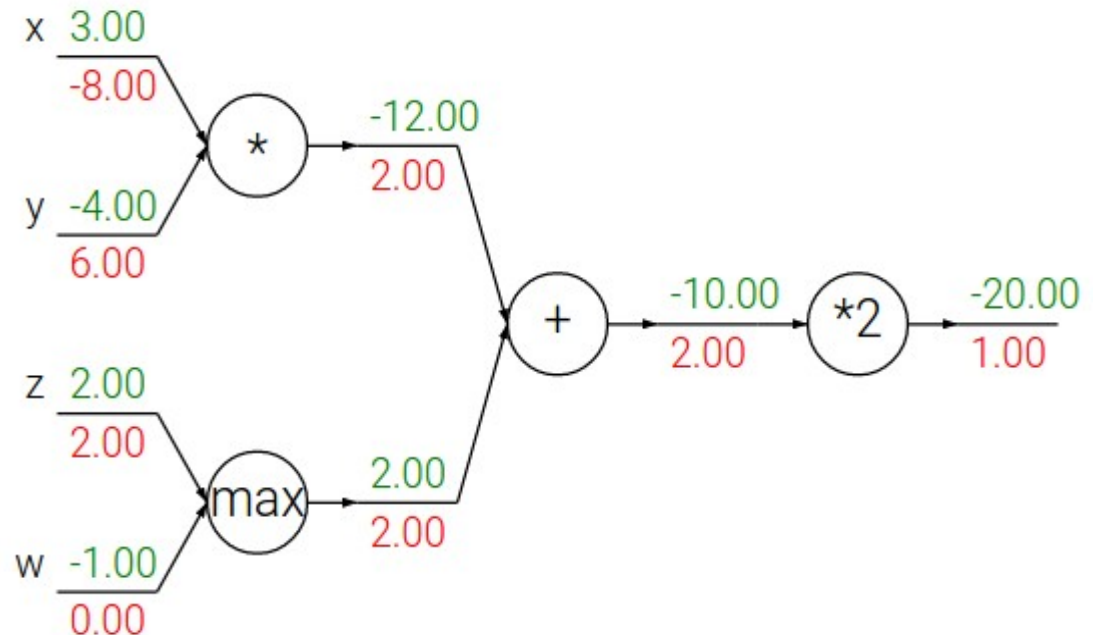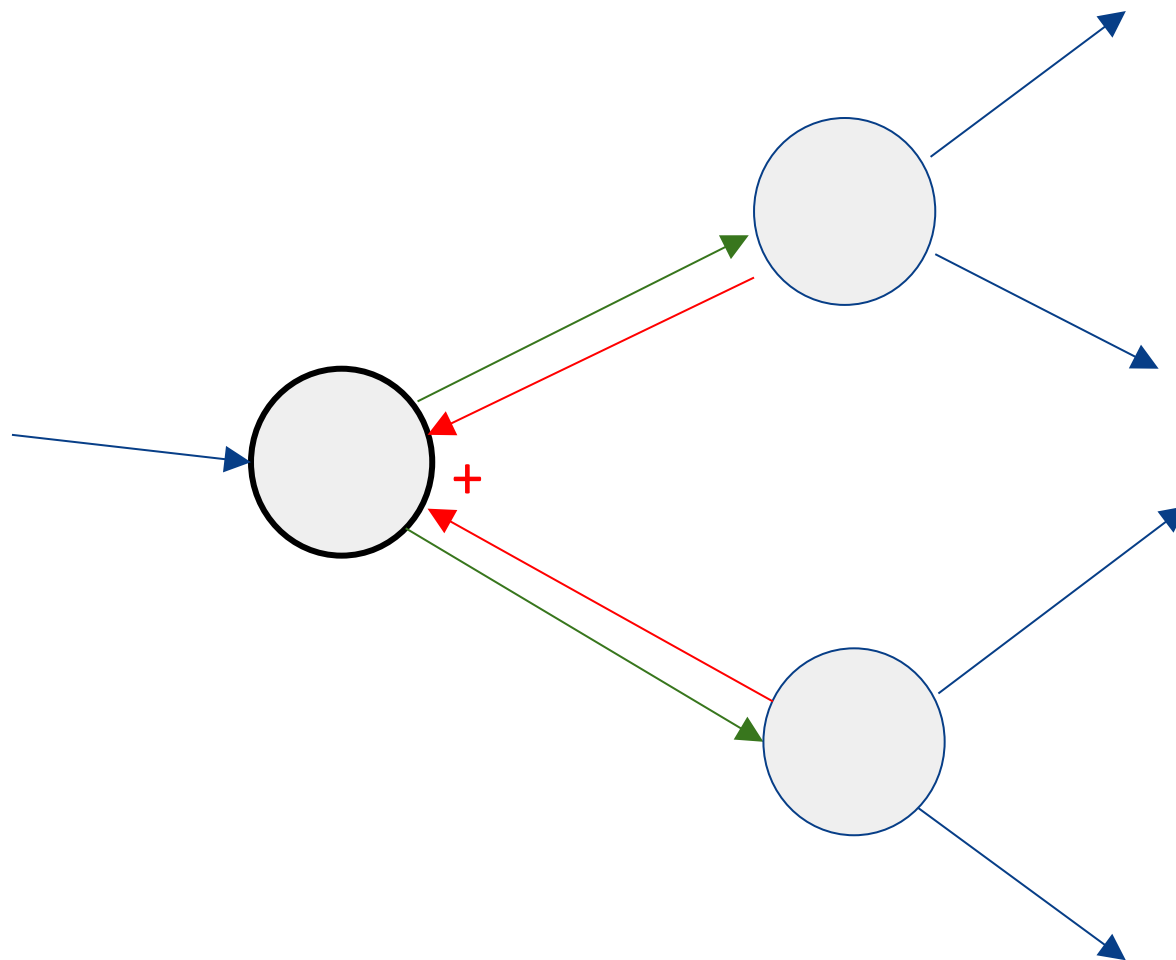sigmoid gate

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

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

(0.73) * (1 - 0.73) = 0.2

# Patterns in backward flow

**add** gate: gradient distributor
**max** gate: gradient router
**mul** gate: gradient... "switcher"?

# Gradients add at branches

Graph (or Net) object. *(Rough psuedo code)*



```python
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
```

# Implementation: forward/backward API



x

z

*

y

(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]
```
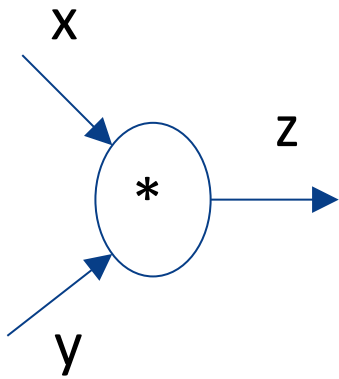
$$\frac{\partial L}{\partial z}$$

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

# Implementation: forward/backward API

x

z

*
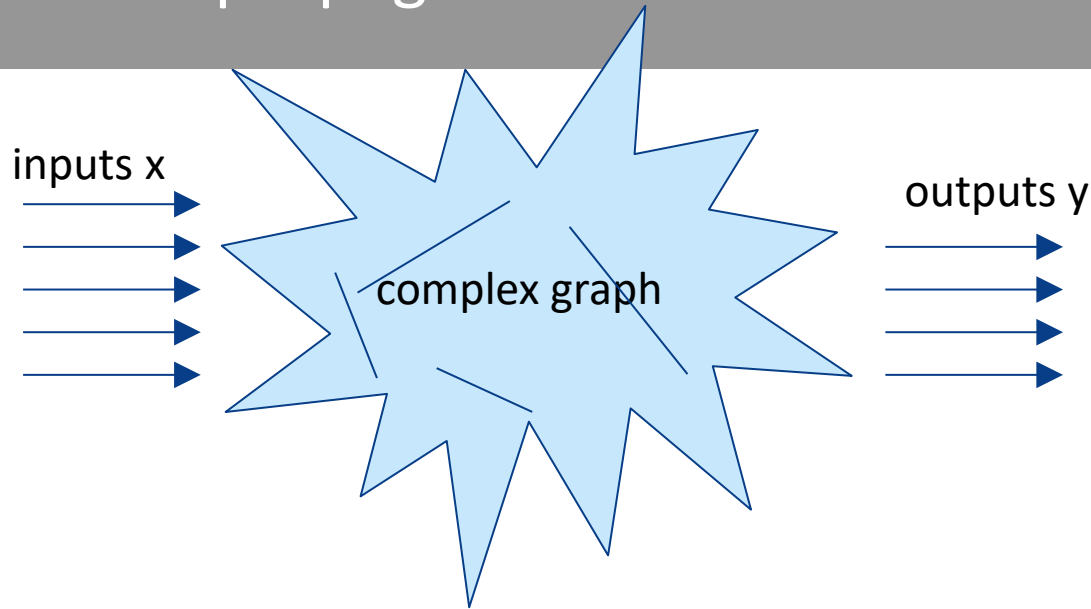
y

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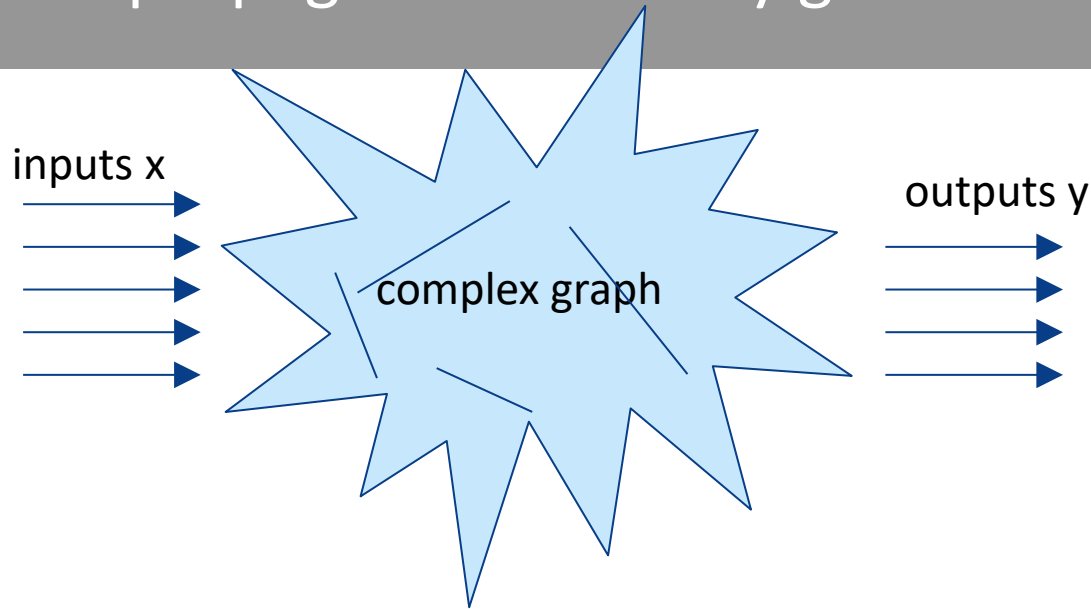
```python
class MultiplyGate(object):
    def forward(x,y):
        z = x*y
        self.x = x # must keep these around!
        self.y = y
        return z
    def backward(dz):
        dx = self.y * dz # [dz/dx * dL/dz]
        dy = self.x * dz # [dz/dy * dL/dz]
        return [dx, dy]
```

inputs x

complex graph

outputs y

inputs x

complex graph

outputs y
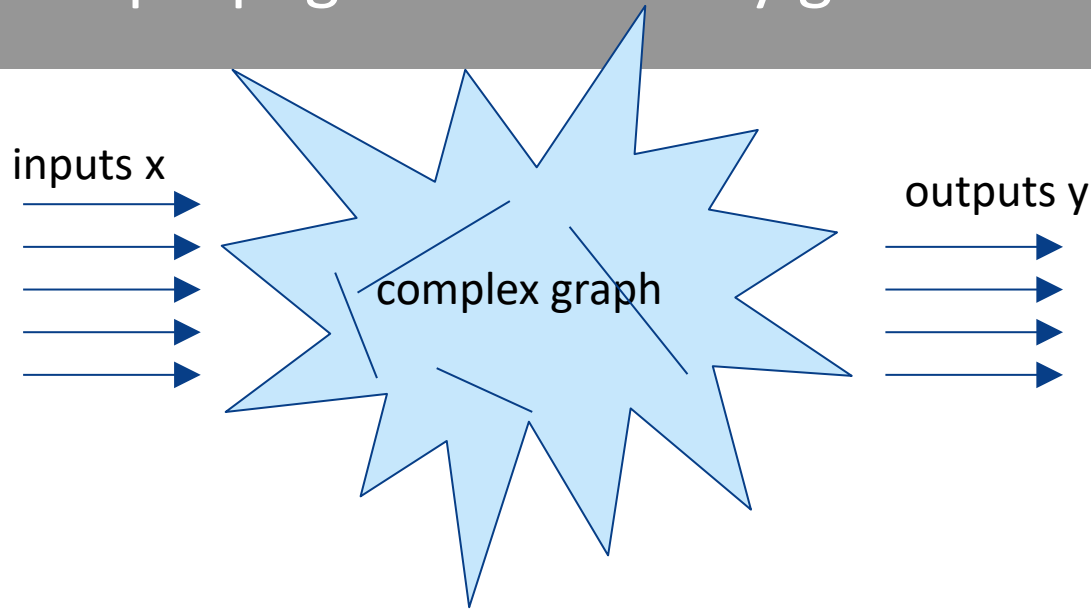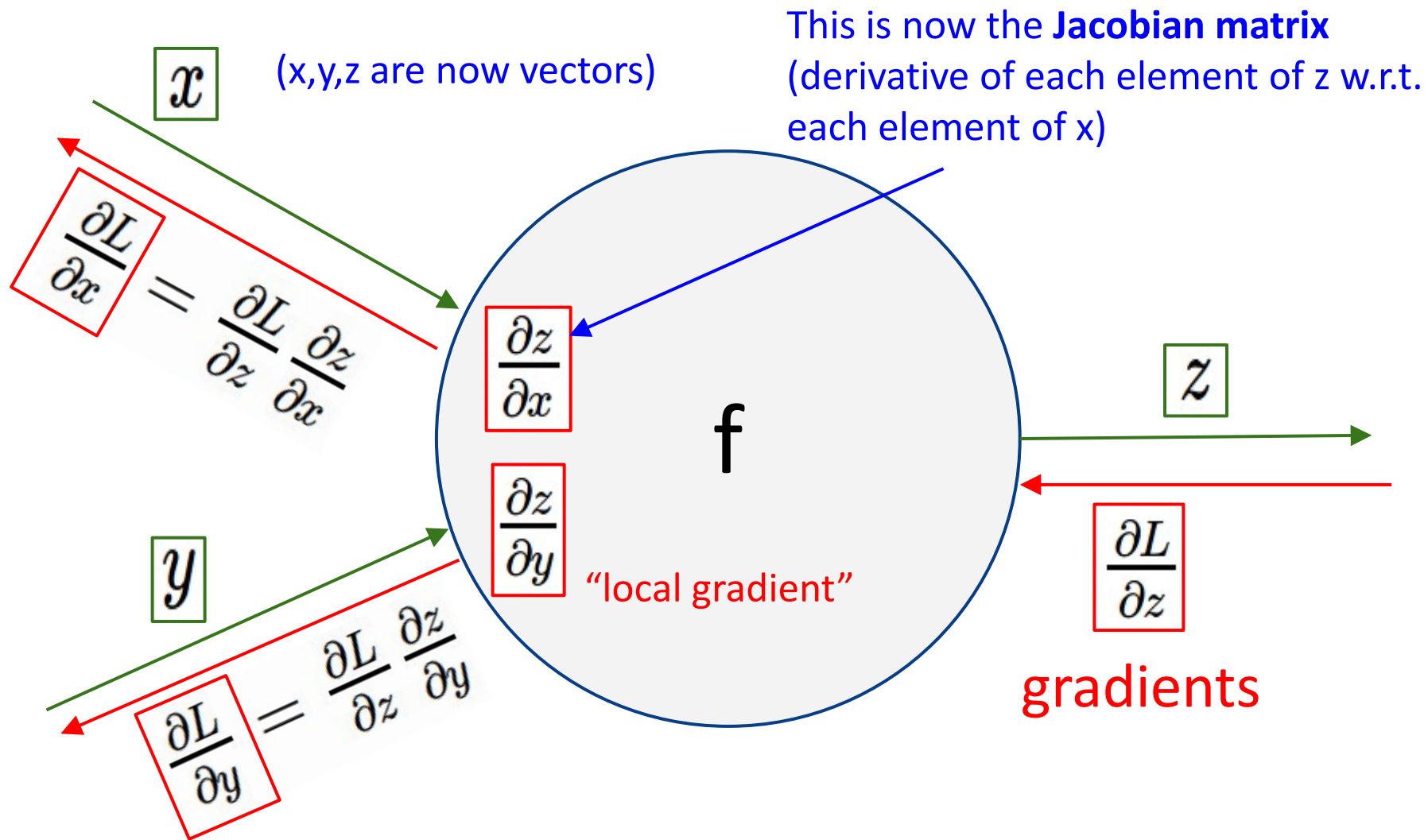
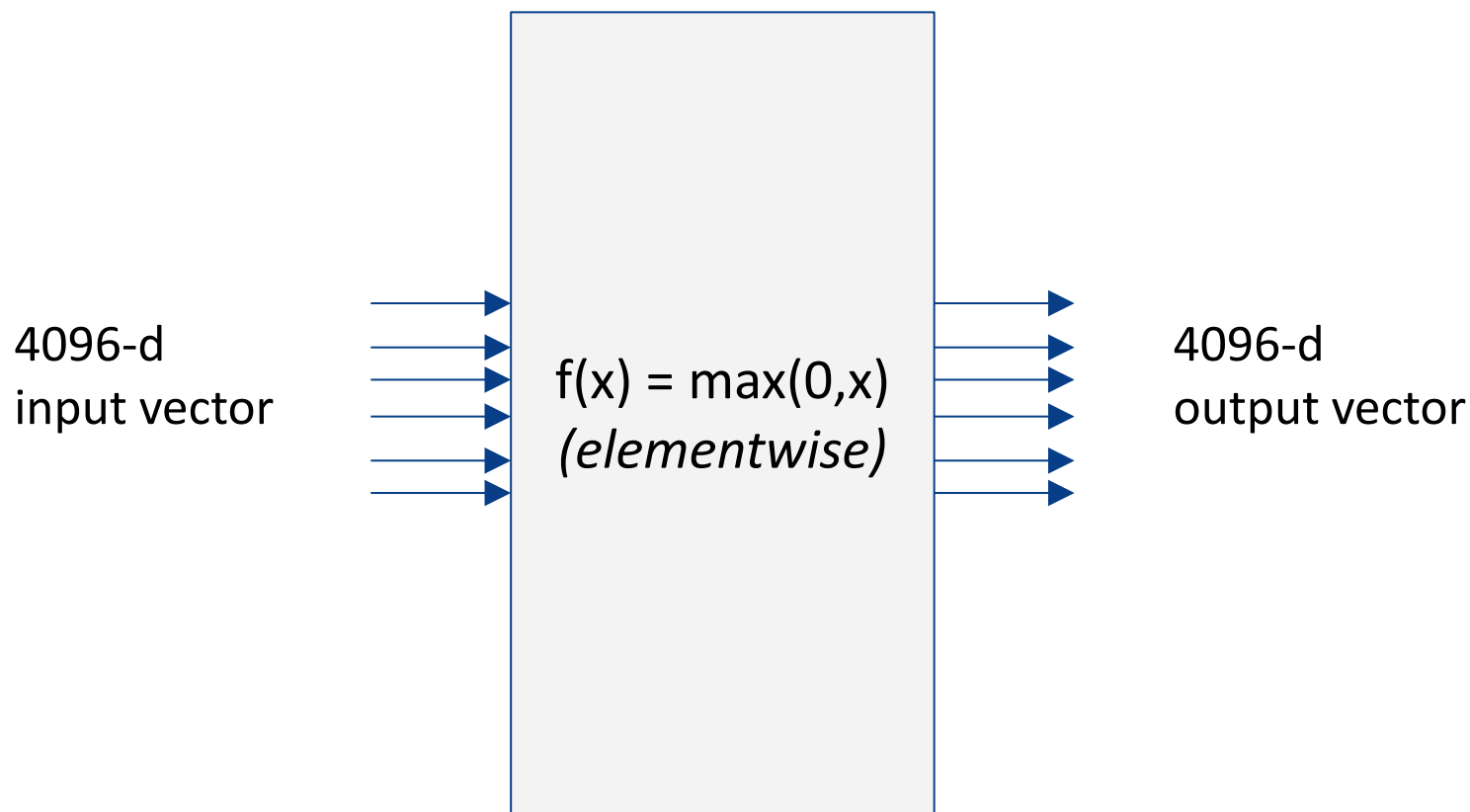reverse-mode differentiation (if you want effect of many things on one thing)

$\frac{\partial y}{\partial x}$ for many different x

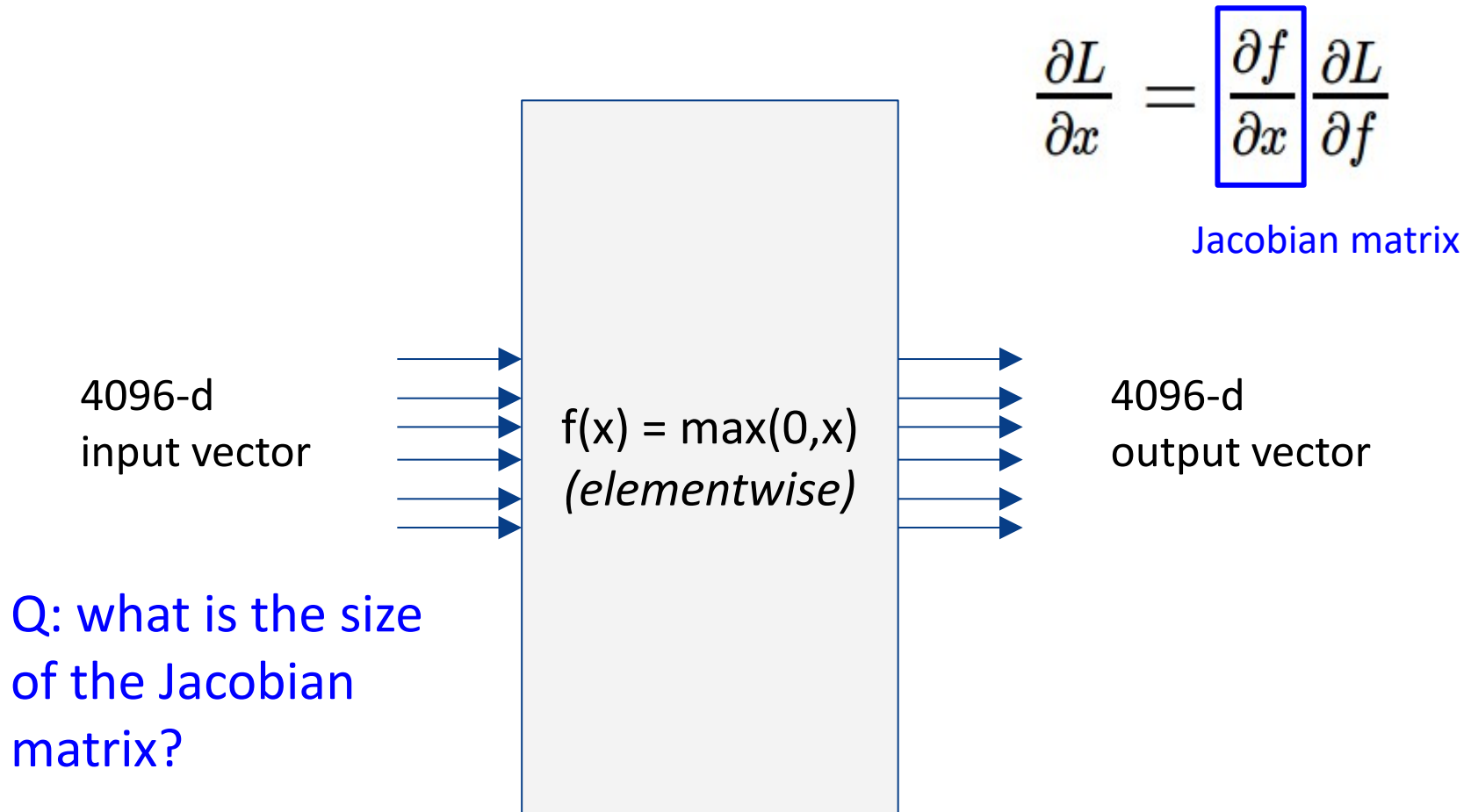More common simply because many nets have a scalar loss function as output.

# Gradients for vector data



(x,y,z are now vectors)

This is now the **Jacobian matrix** (derivative of each element of z w.r.t. each element of x)

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

$x$

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

$$\boxed{\frac{\partial z}{\partial y}}$$ "local gradient"

$y$

$$\boxed{\frac{\partial L}{\partial y}} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial y}$$

$f$

$z$

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

gradients

# Vectorized operations



4096-d
input vector

f(x) = max(0,x)
*(elementwise)*

4096-d
output vector

# Vectorized operations

$$\frac{\partial L}{\partial x} = \boxed{\frac{\partial f}{\partial x}} \frac{\partial L}{\partial f}$$

Jacobian matrix

4096-d
input vector

f(x) = max(0,x)
*(elementwise)*

4096-d
output vector

Q: what is the size of the Jacobian matrix?

$$\frac{\partial L}{\partial x} = \boxed{\frac{\partial f}{\partial x}}\frac{\partial L}{\partial f}$$

Jacobian matrix

4096-d
input vector

f(x) = max(0,x)
*(elementwise)*

4096-d
output vector

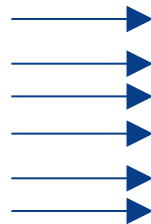Q: what is the size of the Jacobian matrix?
[4096 x 4096!]
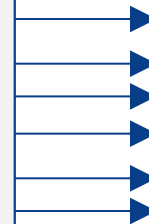
Q2: what does it look like?

# Vectorized operations

in practice we process an entire minibatch (e.g. 100) of examples at one time:

100 4096-d input vectors

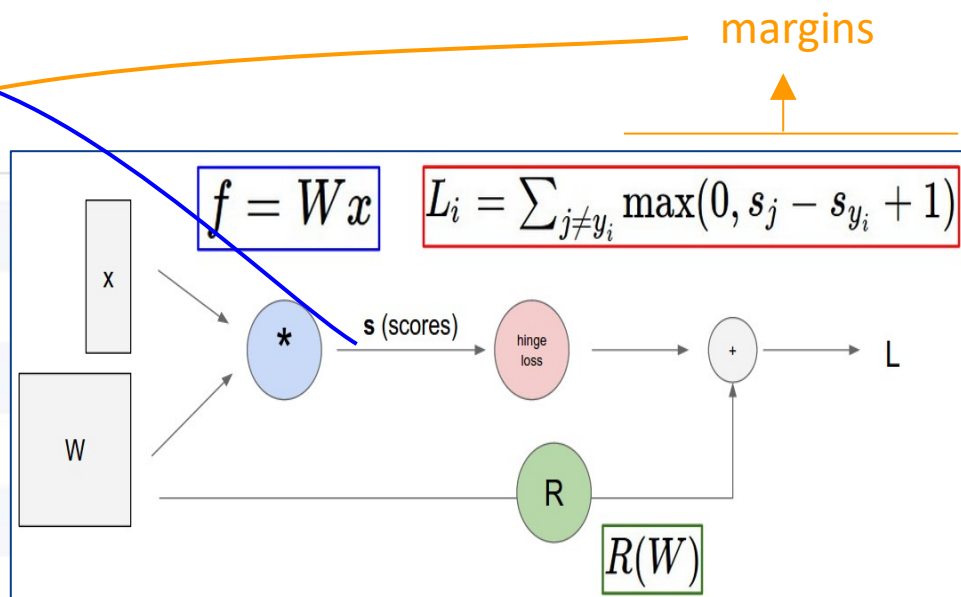$f(x) = max(0,x)$
*(elementwise)*

100 4096-d output vectors

i.e. Jacobian would technically be a [409,600 x 409,600] matrix :\

Why don't we compute it that way?

margins

E.g. for the SVM:

```
# receive W (weights), X (data)
# forward pass (we have 8 lines)
scores = #...
margins = #...
data_loss = #...
reg_loss = #...
loss = data_loss + reg_loss
# backward pass (we have 5 lines)
dmargins = # ... (optionally, we go direct to dscores)
dscores = #...
dW = #...
```

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

x

W

*

s (scores)

hinge loss

+

L

R

$R(W)$

# Summary so far

- neural nets will be very large: no hope of writing down gradient formula by hand for all parameters

- **backpropagation** = recursive application of the chain rule along a computational graph to compute the gradients of all inputs/parameters/intermediates

- implementations maintain a graph structure, where the nodes implement the **forward**() / **backward**() API.

- **forward**: compute result of an operation and save any intermediates needed for gradient computation in memory

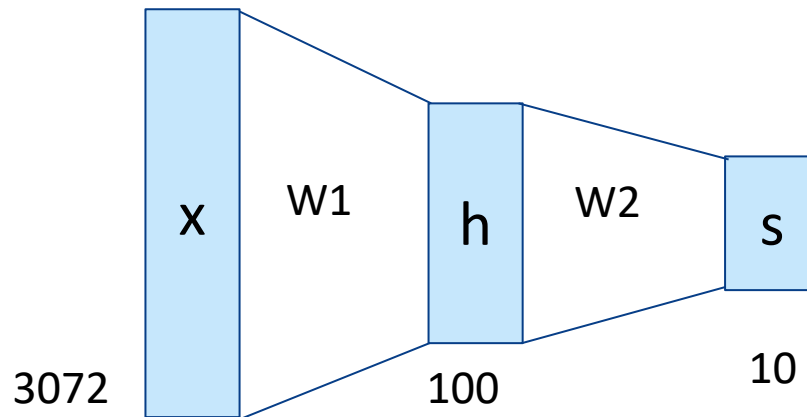- **backward**: apply the chain rule to compute the gradient of the loss function with respect to the inputs.

# Neural Network

2-layer Neural Network

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

3-layer Neural Network:

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



x   W1   h   W2   s

3072      100      10

# Full implementation of training a 2-layer Neural Network needs ~11 lines:

```python
01.  X = np.array([ [0,0,1],[0,1,1],[1,0,1],[1,1,1] ])
02.  y = np.array([[0,1,1,0]]).T
03.  syn0 = 2*np.random.random((3,4)) - 1
04.  syn1 = 2*np.random.random((4,1)) - 1
05.  for j in xrange(60000):
06.      l1 = 1/(1+np.exp(-(np.dot(X,syn0))))
07.      l2 = 1/(1+np.exp(-(np.dot(l1,syn1))))
08.      l2_delta = (y - l2)*(l2*(1-l2))
09.      l1_delta = l2_delta.dot(syn1.T) * (l1 * (1-l1))
10.      syn1 += l1.T.dot(l2_delta)
11.      syn0 += X.T.dot(l1_delta)
```

from @iamtrask, http://iamtrask.github.io/2015/07/12/basic-python-network/

Based on cs231n by Fei-Fei Li & Andrej Karpathy & Justin Johnson

Stage your forward/backward computation!

```
# receive W1,W2,b1,b2 (weights/biases), X (data)

# forward pass:

h1 = #... function of X,W1,b1

scores = #... function of h1,W2,b2

loss = #... (several lines of code to evaluate Softmax loss)

# backward pass:

dscores = #...

dh1,dW2,db2 = #...

dW1,db1 = #...
```
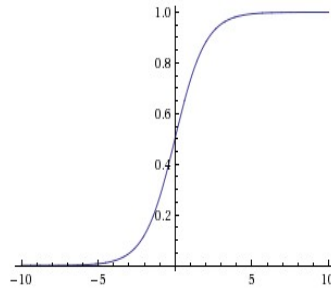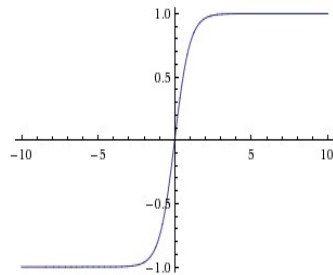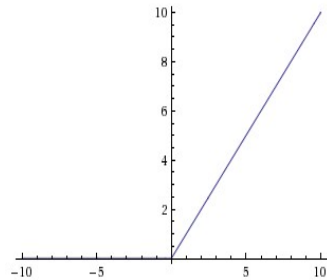
# Activation Functions

**Sigmoid**

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

**Leaky ReLU**
max(0.1x, x)

**Maxout**

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

**tanh**   tanh(x)

**ELU**

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha\,(\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$

**ReLU**   max(0,x)

# Neural Networks: Architectures



"2-layer Neural Net", or
"1-hidden-layer Neural Net"

**"Fully-connected" layers**

"3-layer Neural Net", or
"2-hidden-layer Neural Net"

```python
class Neuron:
  # ...
  def neuron_tick(inputs):
    """ assume inputs and weights are 1-D numpy arrays and bias is a number """
    cell_body_sum = np.sum(inputs * self.weights) + self.bias
    firing_rate = 1.0 / (1.0 + math.exp(-cell_body_sum)) # sigmoid activation function
    return firing_rate
```
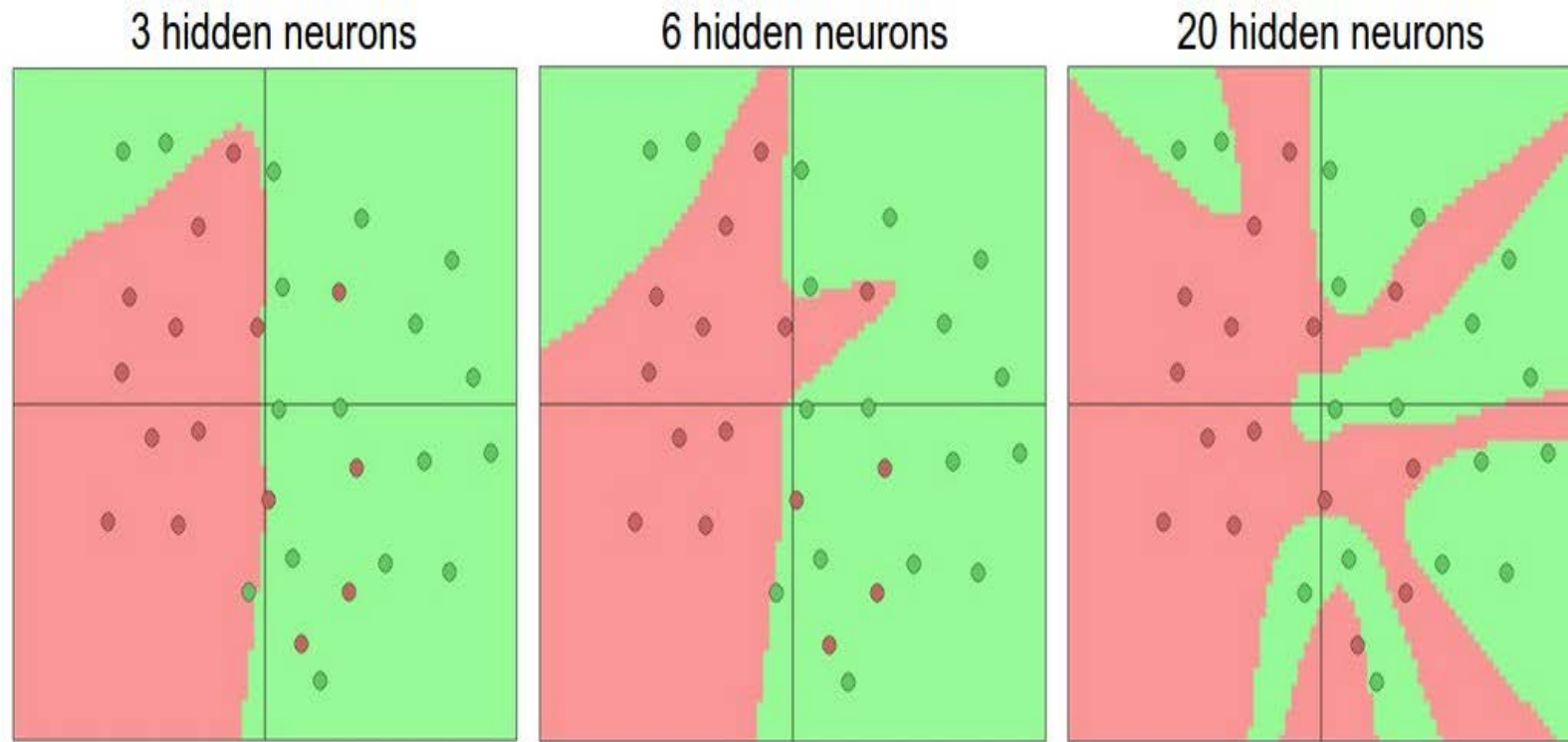
We can efficiently evaluate an entire layer of neurons.

# Example Feed-forward computation of a Neural Network



```
# forward-pass of a 3-layer neural network:
f = lambda x: 1.0/(1.0 + np.exp(-x)) # activation function (use sigmoid)
x = np.random.randn(3, 1) # random input vector of three numbers (3x1)
h1 = f(np.dot(W1, x) + b1) # calculate first hidden layer activations (4x1)
h2 = f(np.dot(W2, h1) + b2) # calculate second hidden layer activations (4x1)
out = np.dot(W3, h2) + b3 # output neuron (1x1)
```

3 hidden neurons    6 hidden neurons    20 hidden neurons

more neurons = more capacity

# Summary

- we arrange neurons into fully-connected layers

- the abstraction of a **layer** has the nice property that it allows us to use efficient vectorized code (e.g. matrix multiplies)

- neural networks are not really *neural*

- neural networks: bigger = better (but might have to regularize more strongly)