

Lecture 4: Backpropagation and Neural Networks (part 1)

Tuesday January 31, 2017

Announcements!

- If you are adversely affected by immigration ban, please talk to me about accommodations
- Send in paper choices by **tonight**
- Should be able to run Jupyter server on Tufts was and network machines now
 - (deep-venv)> pip install --upgrade jupyter
- hw1 deadline in two days — Thurs Feb 2: Don't forget to read the course notes.
- Redo calculation of dL/dW for hinge loss

Python/Numpy of the Day

- `y_pred = scores.argmax(axis=1)`
- `inds = np.random.choice(X.shape[0], batch_size)`
 - randomly select N numbers in a range,
 - useful for subsampling
- `[:, np.newaxis]`
 - reshapes matrices of size (N,) to size (N,1)

Where we are...

$$s = f(x; W) = Wx$$

scores function

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

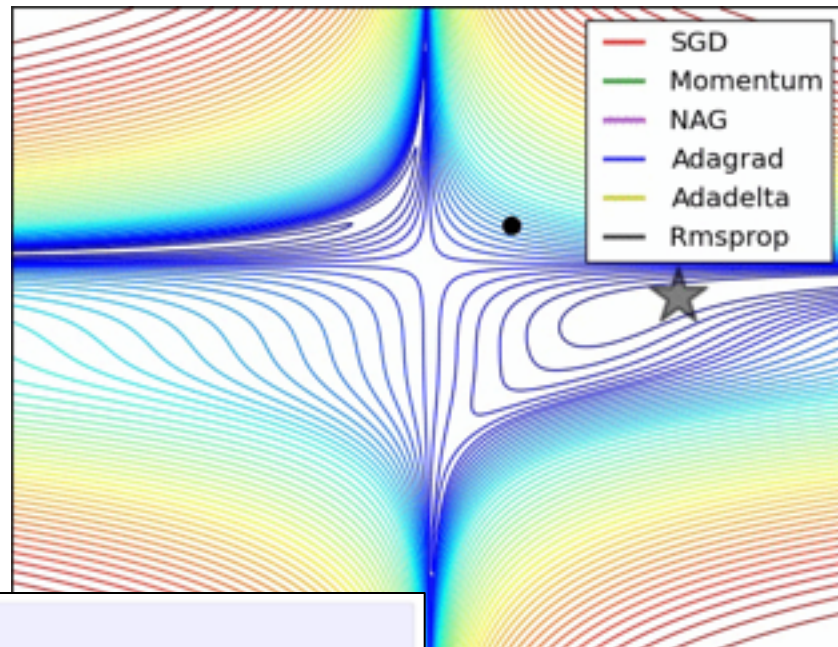
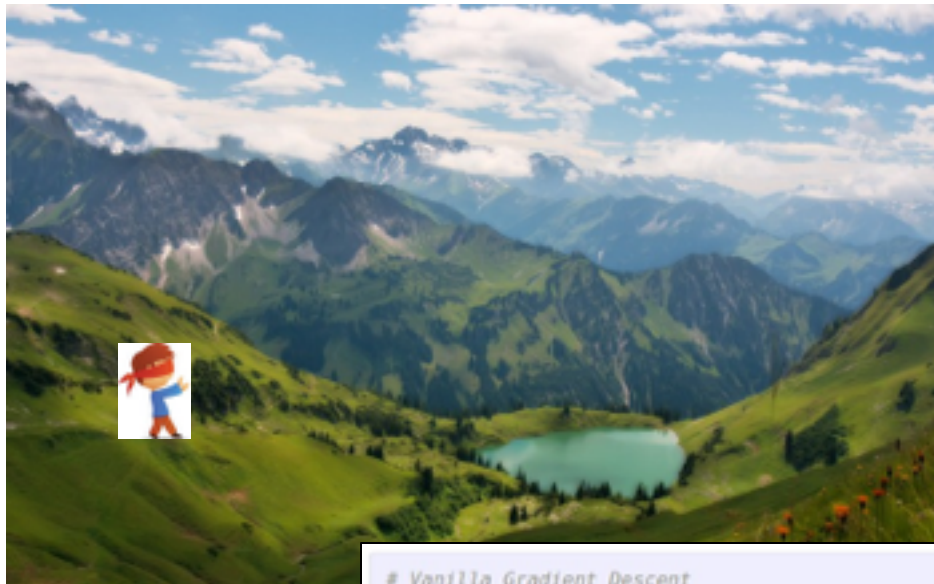
SVM loss

$$L = \frac{1}{N} \sum_{i=1}^N L_i + \sum_k W_k^2$$

data loss + regularization

want $\nabla_W L$

Optimization



(image credits
to Alec Radford)

```
# Vanilla Gradient Descent

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

Gradient Descent

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

Numerical gradient: slow :(, approximate :(, easy to write :)

Analytic gradient: fast :), exact :), error-prone :(

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

Hinge Loss Gradient wrt Weights W

margin size, usually 1.0

$$L_i = \sum_{j \neq y_i} [\max(0, w_j^T x_i - w_{y_i}^T x_i + \Delta)]$$

- We want the Jacobian Matrix of all gradients
 - partial derivatives of all output dimensions by all input dimensions

$$\nabla w L = \begin{bmatrix} \nabla w_1 L_1 & \dots & \dots & \nabla w_1 L_N \\ \vdots & \nabla w_j L_i & \ddots & \vdots \\ \nabla w_k L_1 & \dots & \dots & \nabla w_k L_N \end{bmatrix}$$

For all rows of dW where the row corresponds to the GT value for that training instance, i.e. $j = y_i$

$$\nabla_{w_{y_i}} L_i = - \left(\sum_{j \neq y_i} \mathbb{1}(w_j^T x_i - w_{y_i}^T x_i + \Delta > 0) \right) x_i$$

For all rows of dW where $j \neq y_i$

$$\nabla_{w_j} L_i = \mathbb{1}(w_j^T x_i - w_{y_i}^T x_i + \Delta > 0) x_i$$

Softmax Loss Gradient wrt Score S

* note change of subscripts from last slide

$$\begin{array}{l} a_j = w_j^T x_j \\ S_j = \frac{e^{a_j}}{\sum_{k=1}^N e^{a_k}} \quad \forall j \in 1..N \\ \frac{\partial S_i}{\partial a_j} = \frac{\partial \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}}}{\partial a_j} \end{array} \left| \begin{array}{l} \nabla a_j S_i, \text{ when } i = j \\ \frac{\partial \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}}}{\partial a_j} = \frac{e^{a_i} \Sigma - e^{a_j} e^{a_i}}{\Sigma^2} \\ = \frac{e^{a_i}}{\Sigma} \frac{\Sigma - e^{a_j}}{\Sigma} \\ = S_i(1 - S_j) \end{array} \right| \begin{array}{l} \nabla a_j S_i, \text{ when } i \neq j \\ \frac{\partial \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}}}{\partial a_j} = \frac{0 - e^{a_j} e^{a_i}}{\Sigma^2} \\ = -\frac{e^{a_j}}{\Sigma} \frac{e^{a_i}}{\Sigma} \\ = -S_j S_i \end{array}$$

Skipping some steps for space,
please see original notes.

$$\nabla a_j S_i = S_i(\mathbb{1}(i = j) - S_j)$$

Softmax Loss Gradient wrt Score S

$$a_j = w_j^T x_j$$

$$S_j = \frac{e^{a_j}}{\sum_{k=1}^N e^{a_k}} \quad \forall j \in 1..N$$

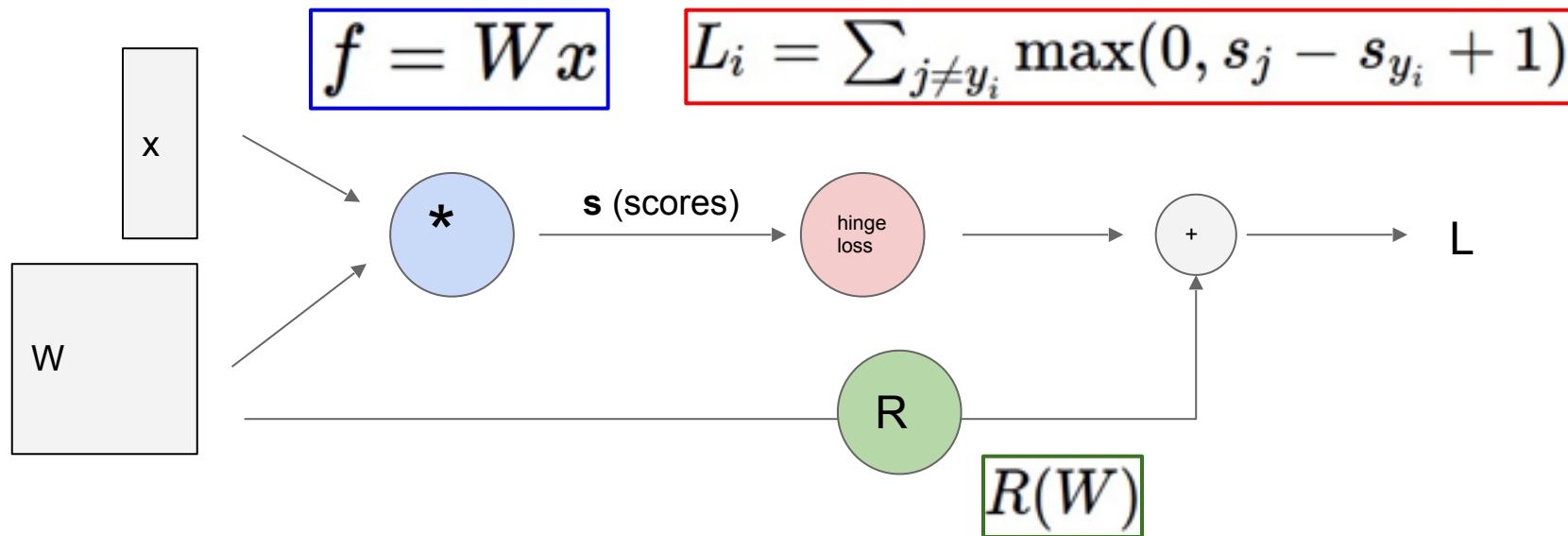
$$\nabla a_j S_i = S_i (\mathbb{1}(i = j) - S_j)$$

$$\nabla S_i L = \frac{\partial}{\partial S_i} - \log(S_i) = S_j - \mathbb{1}(i = j)$$

$$\nabla W_j L = \frac{\partial L}{\partial S_i} * \frac{\partial S_i}{\partial W_j} = (S_j - \mathbb{1}(i = j)) x_i$$

Skipping some steps for space,
please see original notes.

Computational Graph

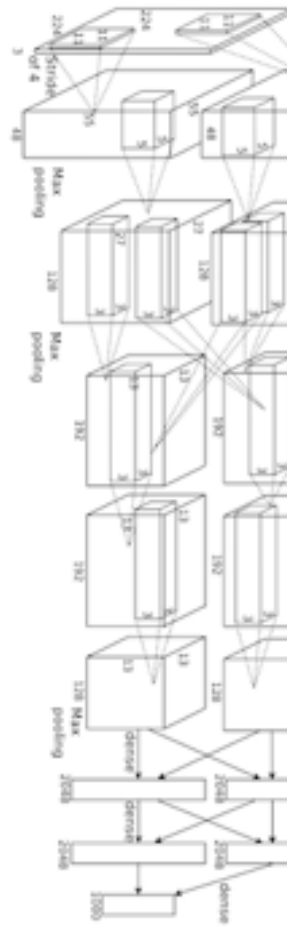


Convolutional Network (AlexNet)

input image

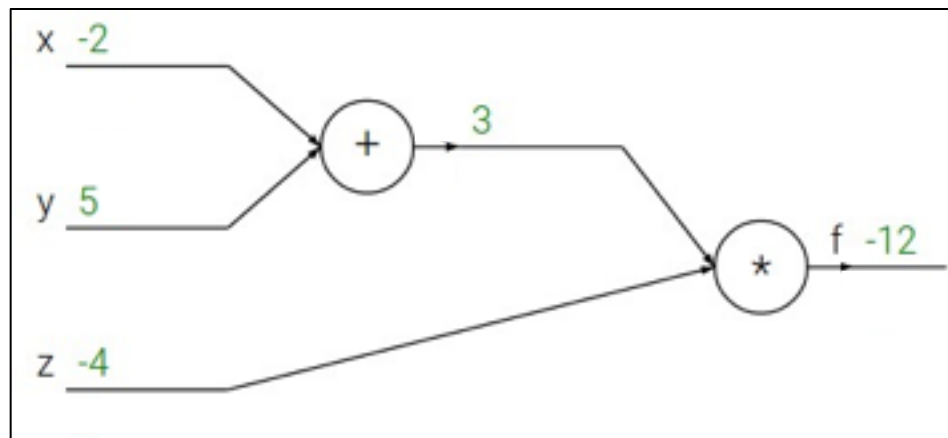
weights

loss



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

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



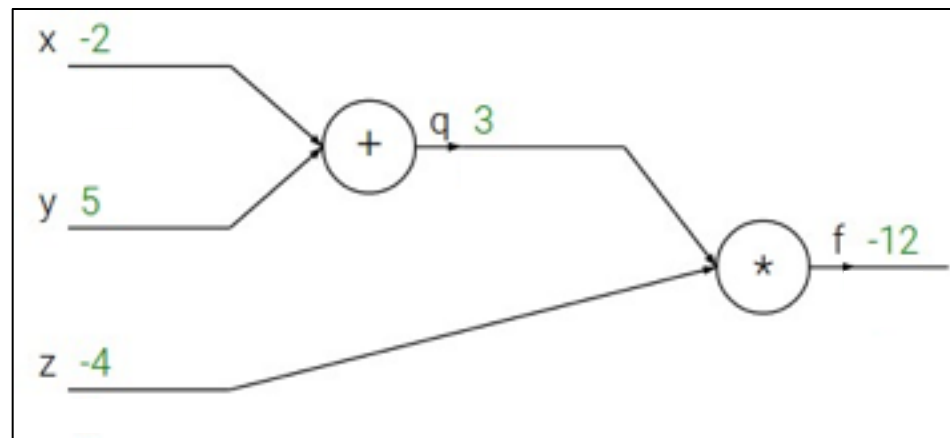
$$f(x, y, z) = (x + y)z$$

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

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

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

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



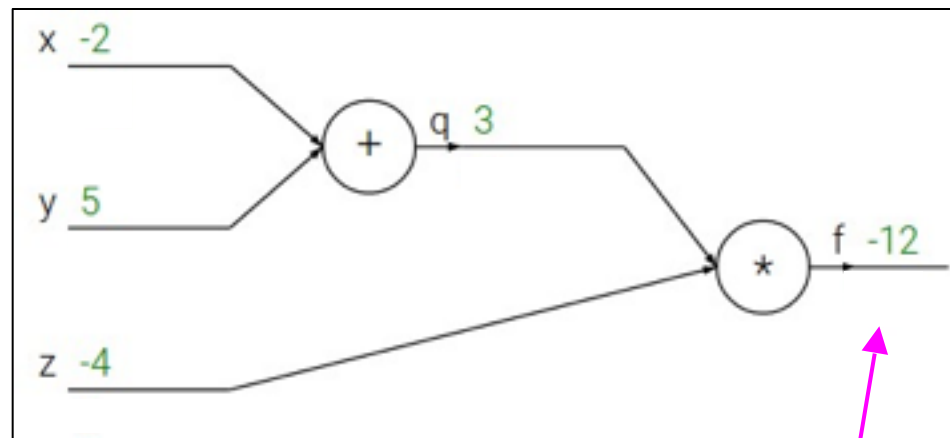
$$f(x, y, z) = (x + y)z$$

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

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

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

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



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

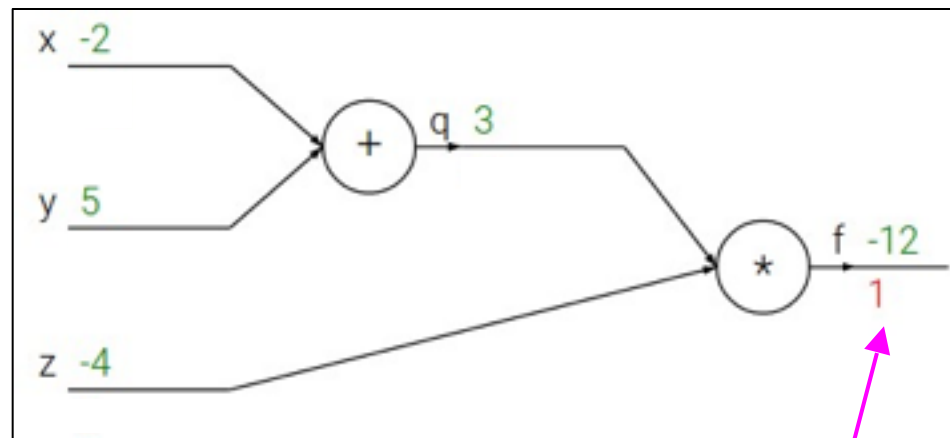
$$f(x, y, z) = (x + y)z$$

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

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

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

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



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

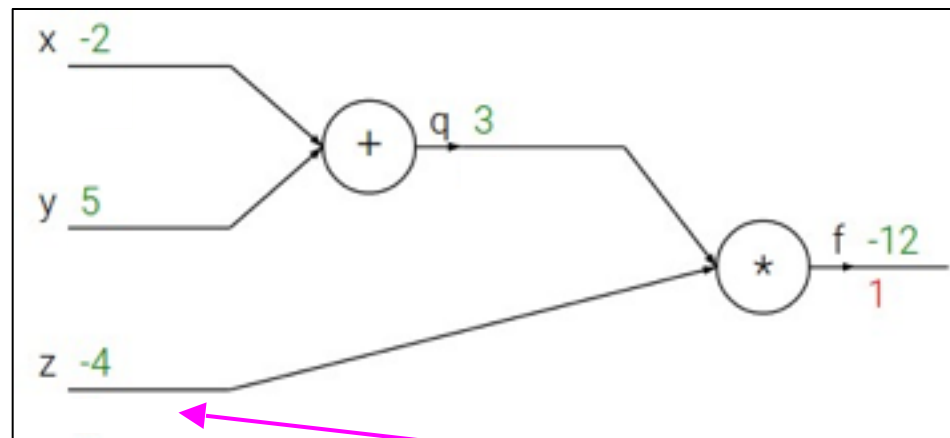
$$f(x, y, z) = (x + y)z$$

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

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

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

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



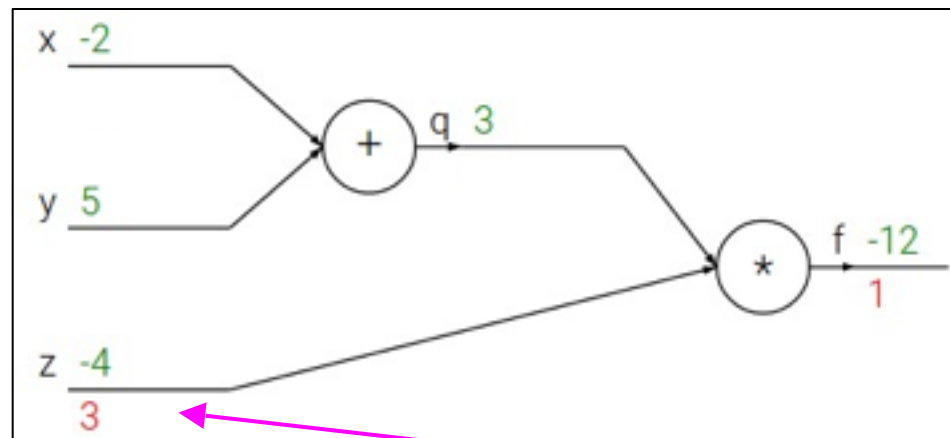
$$f(x, y, z) = (x + y)z$$

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

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

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

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



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

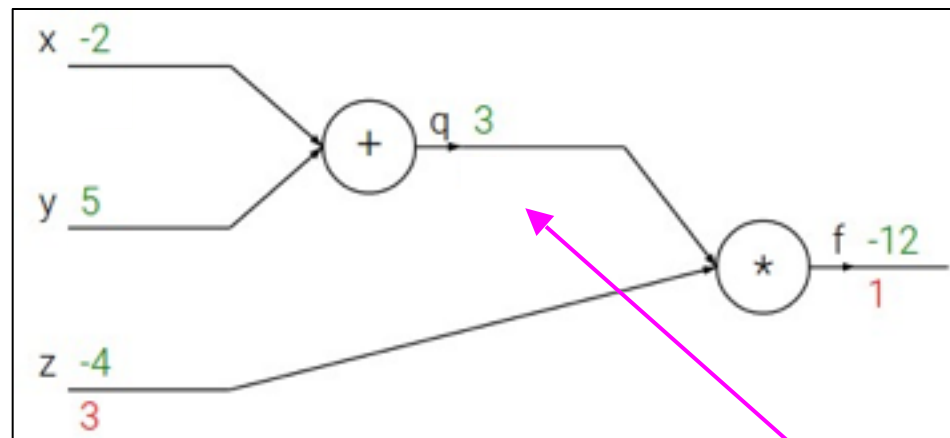
$$f(x, y, z) = (x + y)z$$

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

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

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

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



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

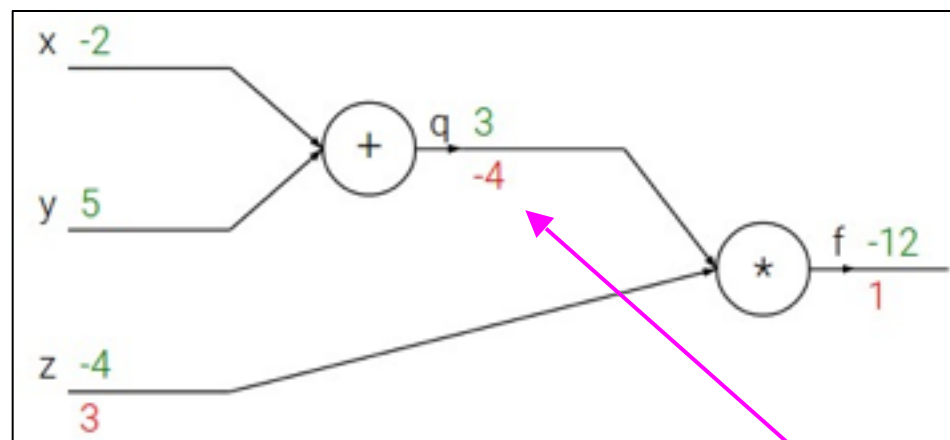
$$f(x, y, z) = (x + y)z$$

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

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

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

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



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

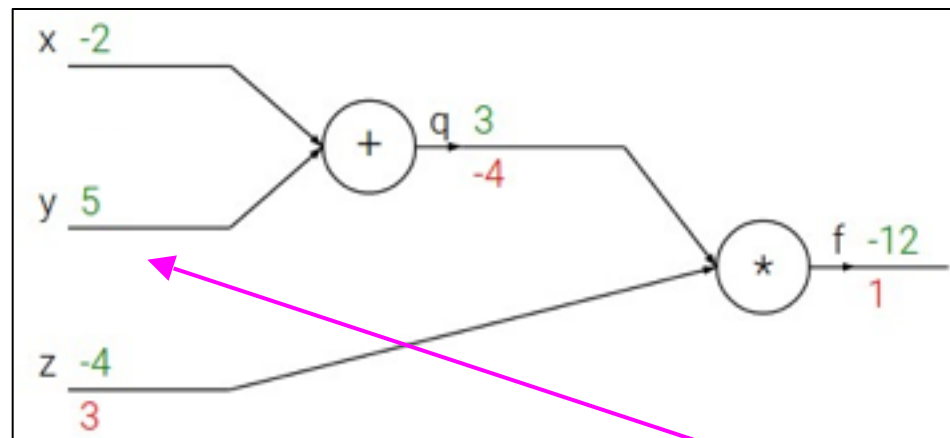
$$f(x, y, z) = (x + y)z$$

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

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

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

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



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

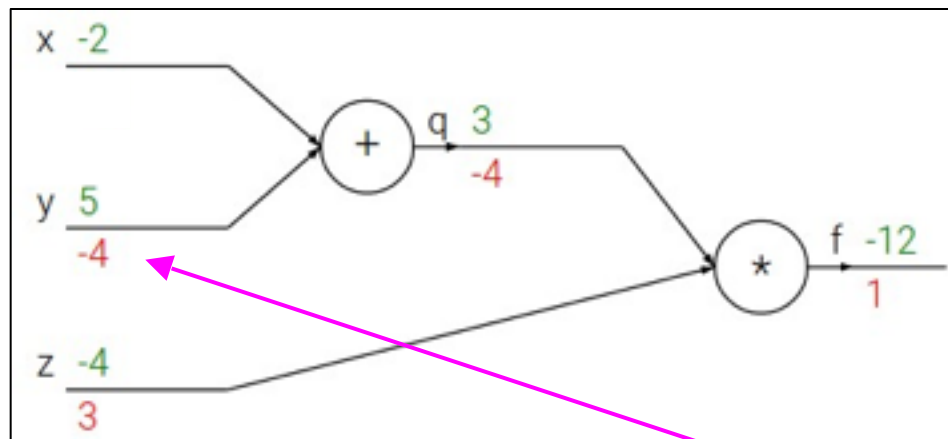
$$f(x, y, z) = (x + y)z$$

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

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

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

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



Chain rule:

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

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

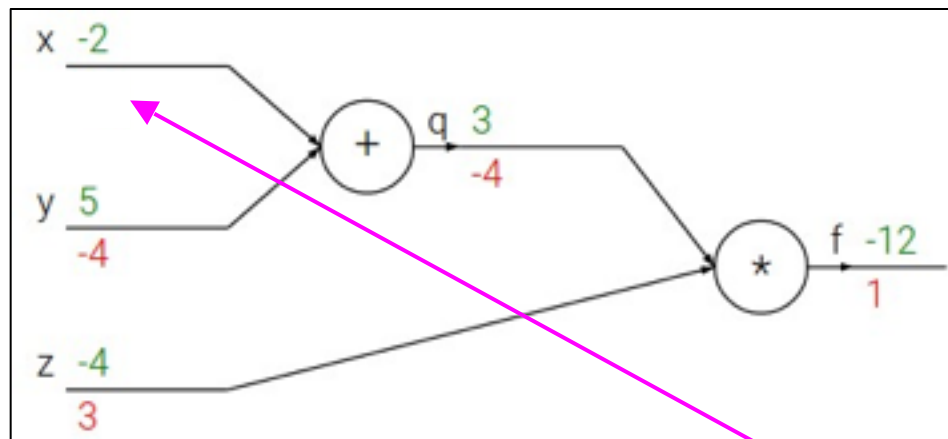
$$f(x, y, z) = (x + y)z$$

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

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

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

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



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

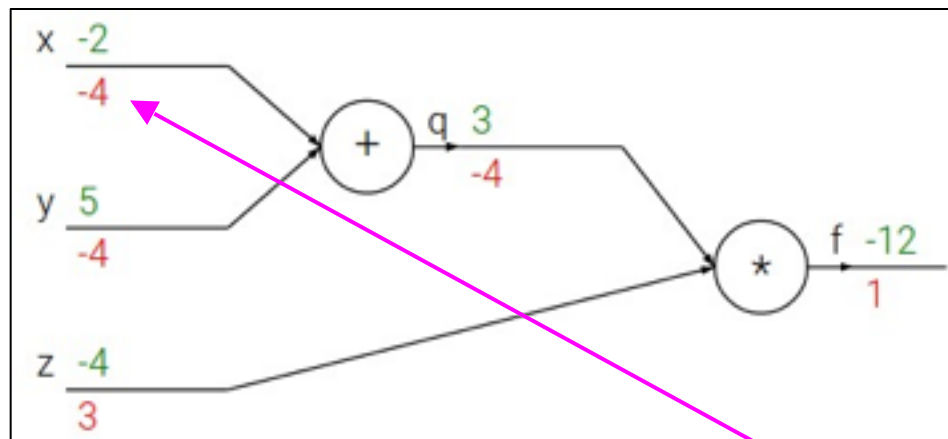
$$f(x, y, z) = (x + y)z$$

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

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

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

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

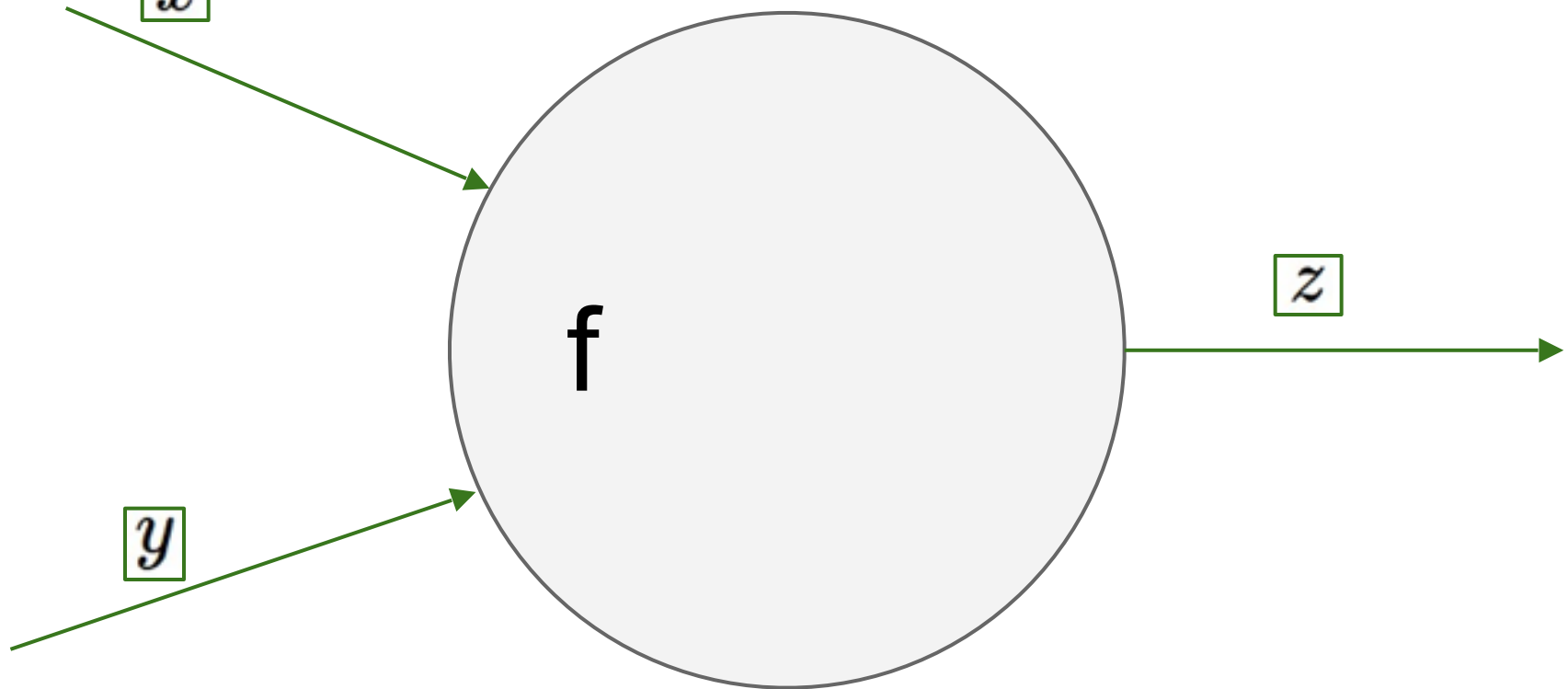


Chain rule:

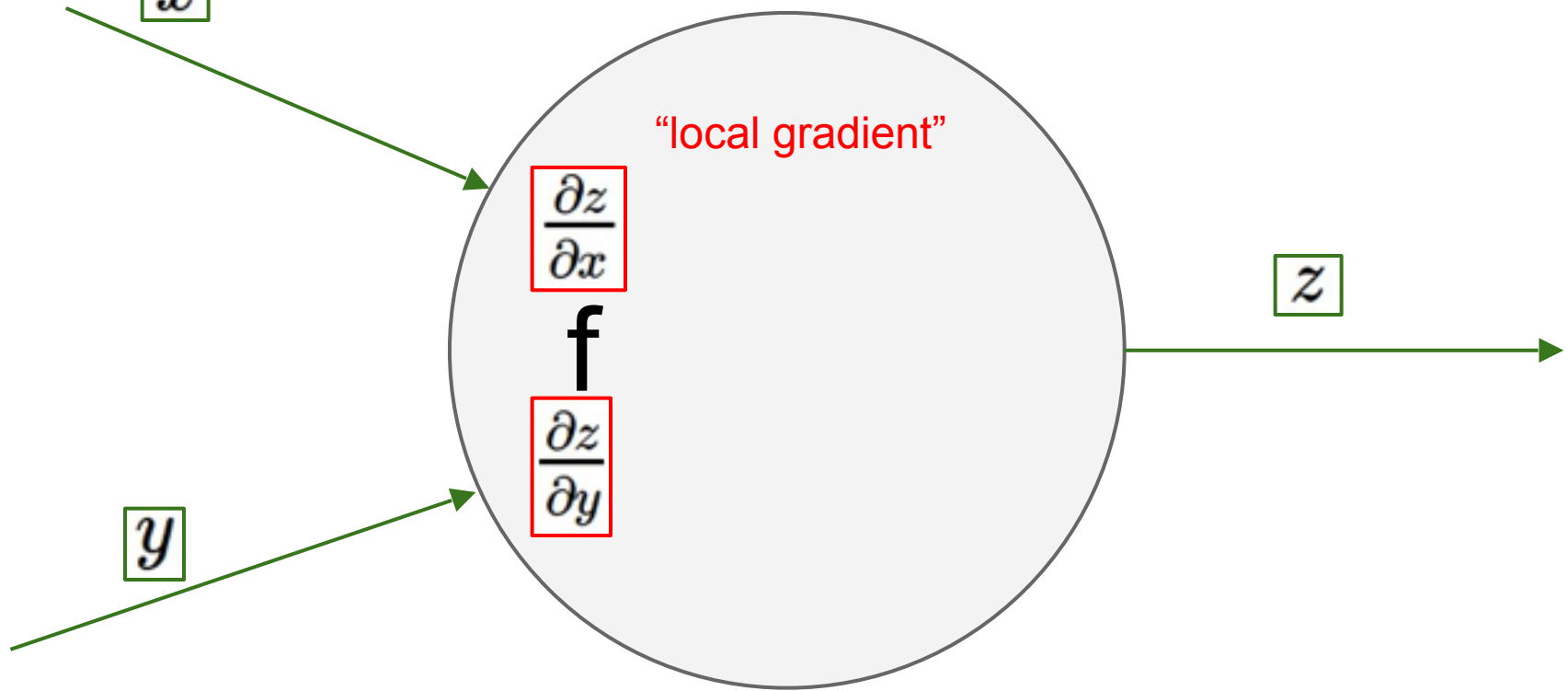
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

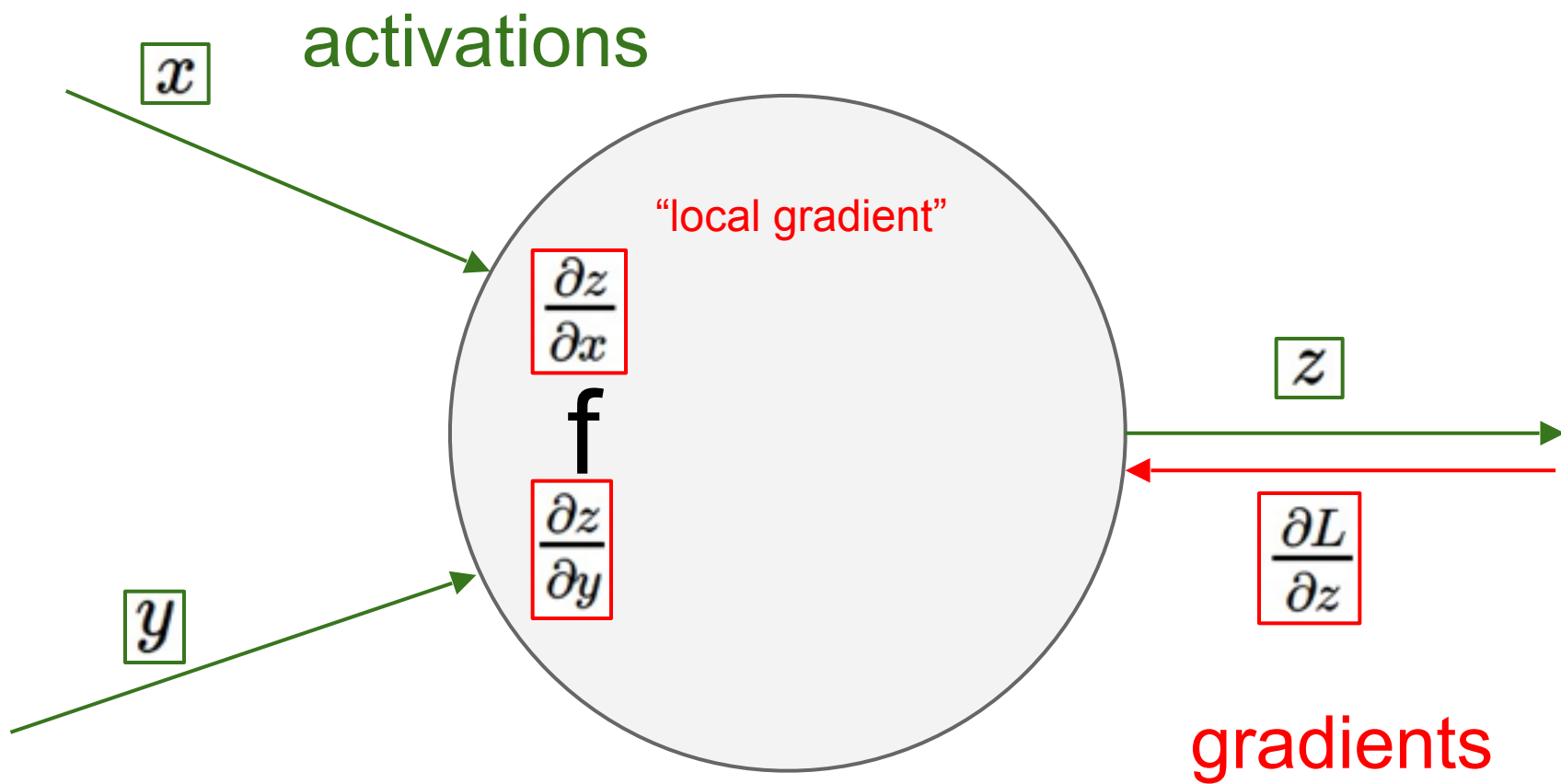
$$\frac{\partial f}{\partial x}$$

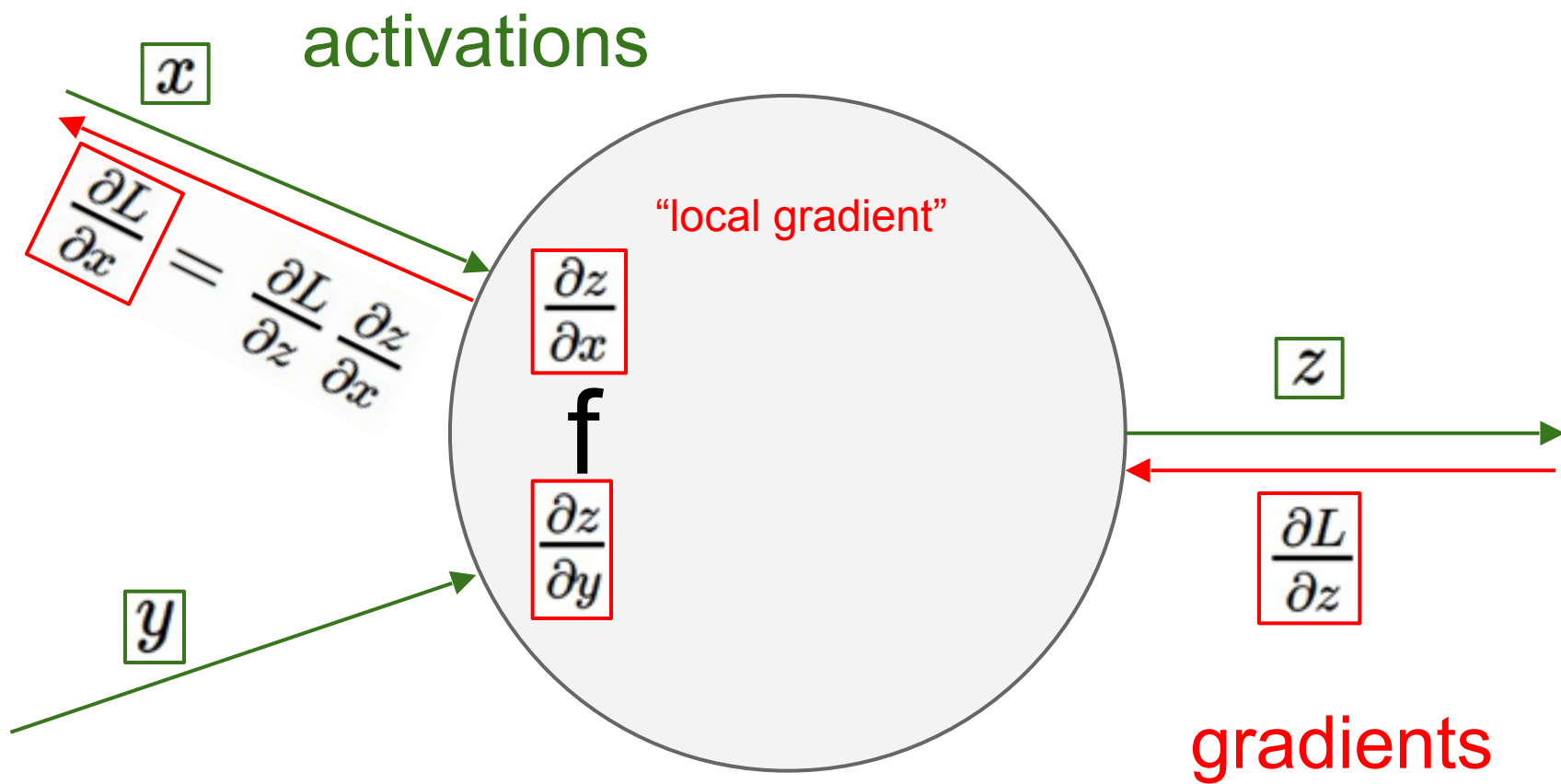
activations



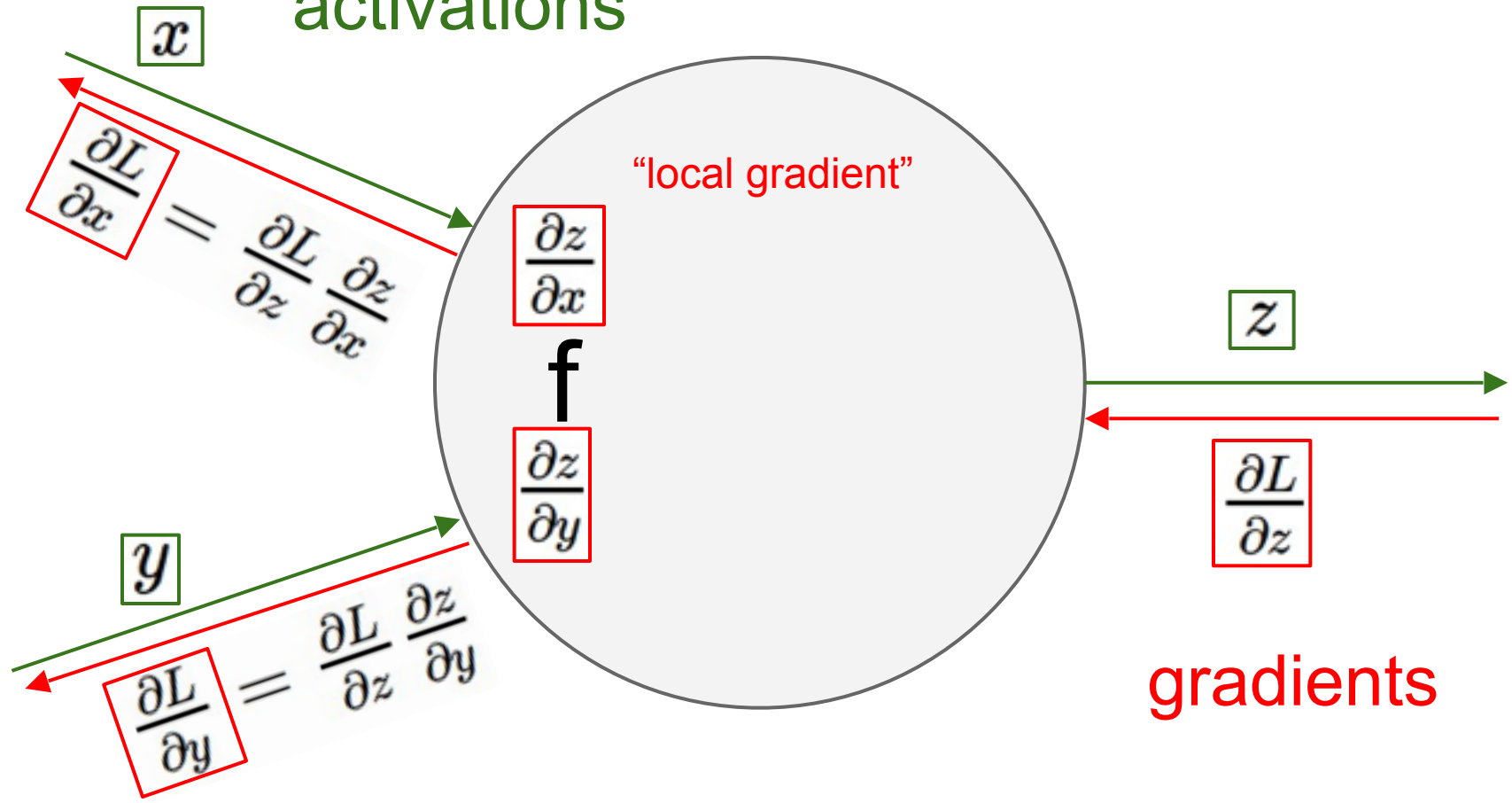
activations







activations



activations

x

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

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

“local gradient”

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

f

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

y

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

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

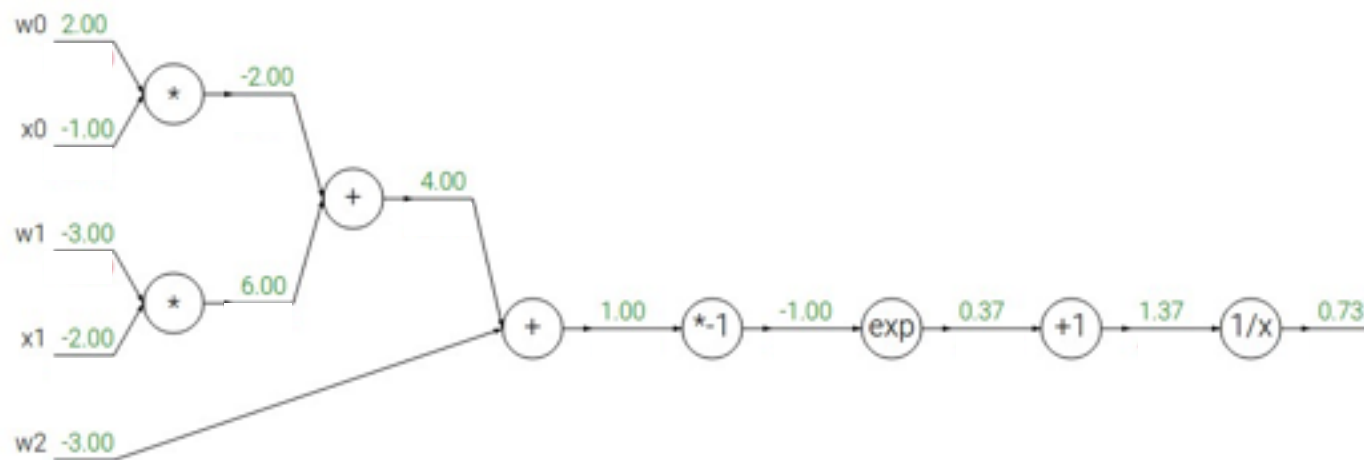
z

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

gradients

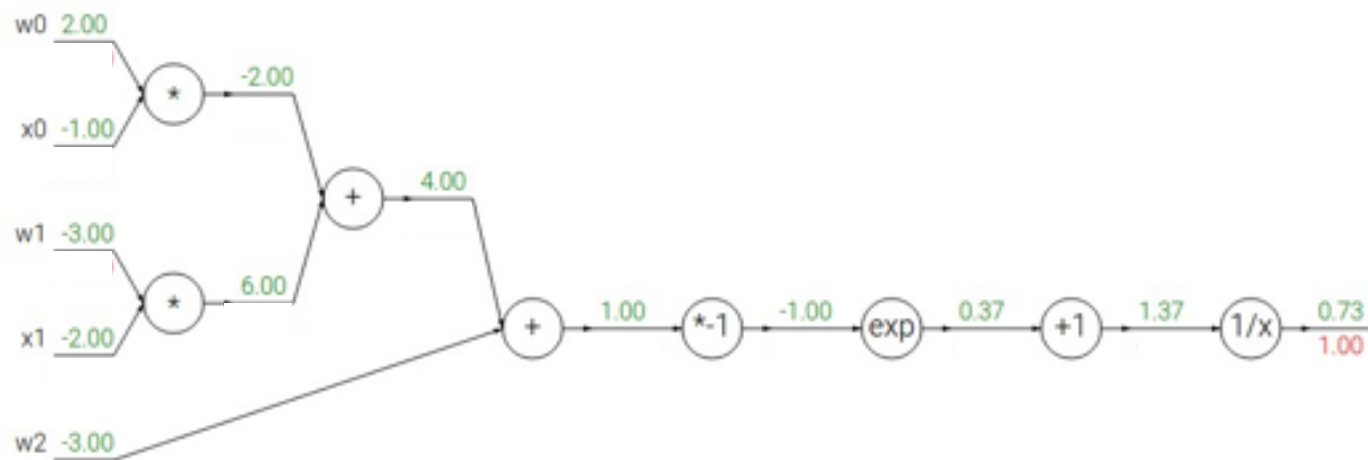
Another example:

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



Another example:

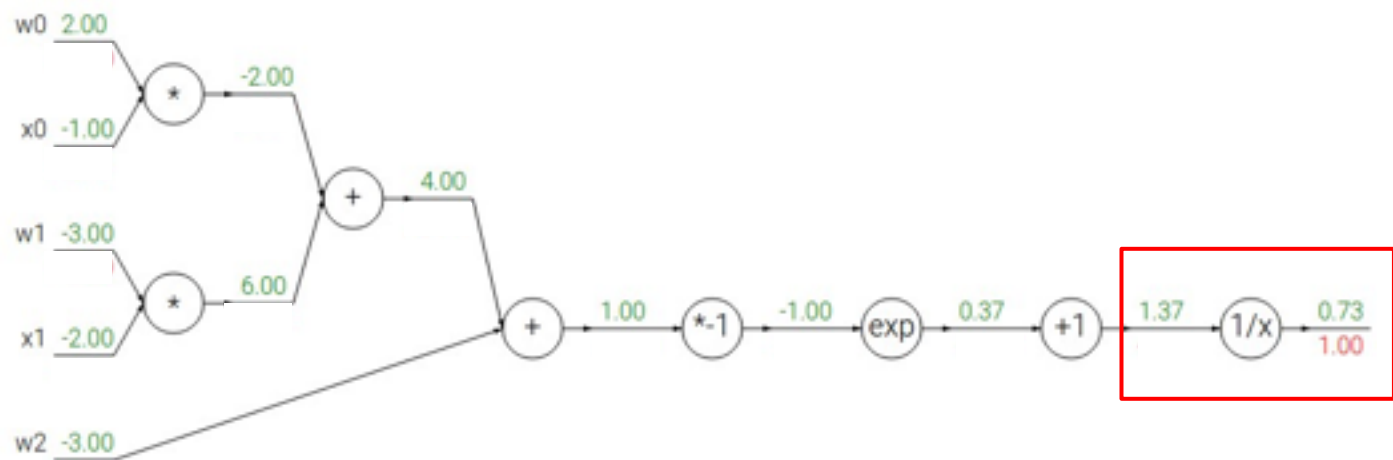
$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$



$f(x) = e^x$	\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:

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



$$f(x) = e^x$$

→

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

→

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

→

$$\frac{df}{dx} = -1/x^2$$

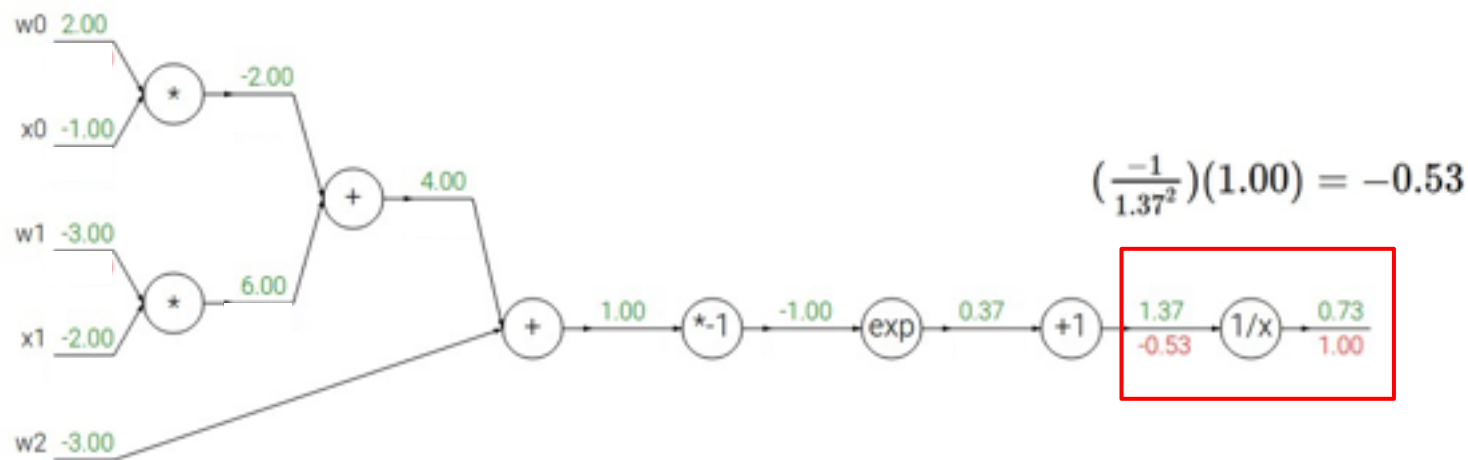
$$f_c(x) = c + x$$

→

$$\frac{df}{dx} = 1$$

Another example:

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



$$f(x) = e^x$$

→

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

→

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

→

$$\frac{df}{dx} = -1/x^2$$

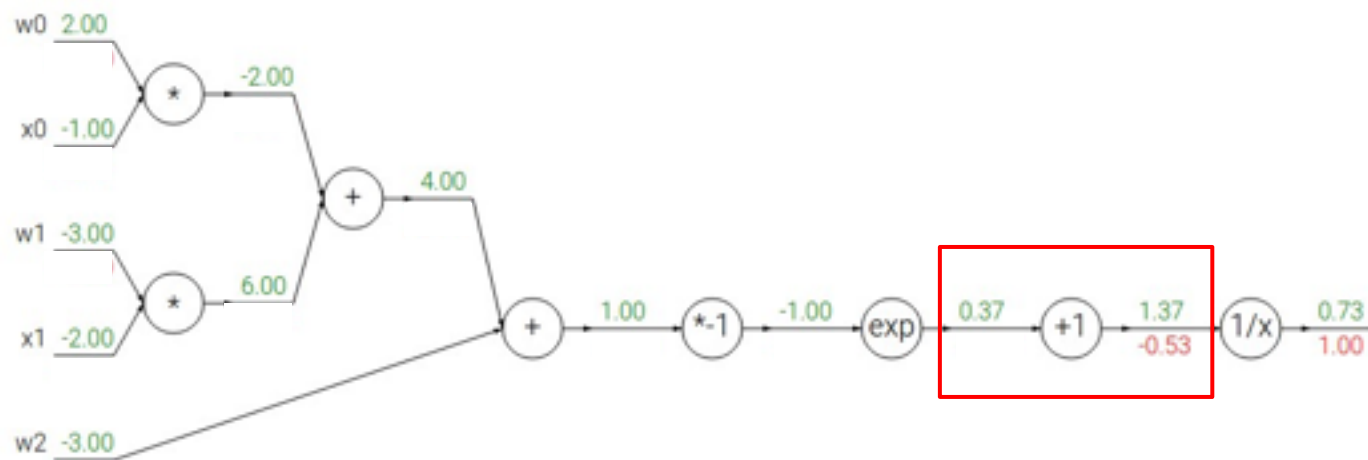
$$f_c(x) = c + x$$

→

$$\frac{df}{dx} = 1$$

Another example:

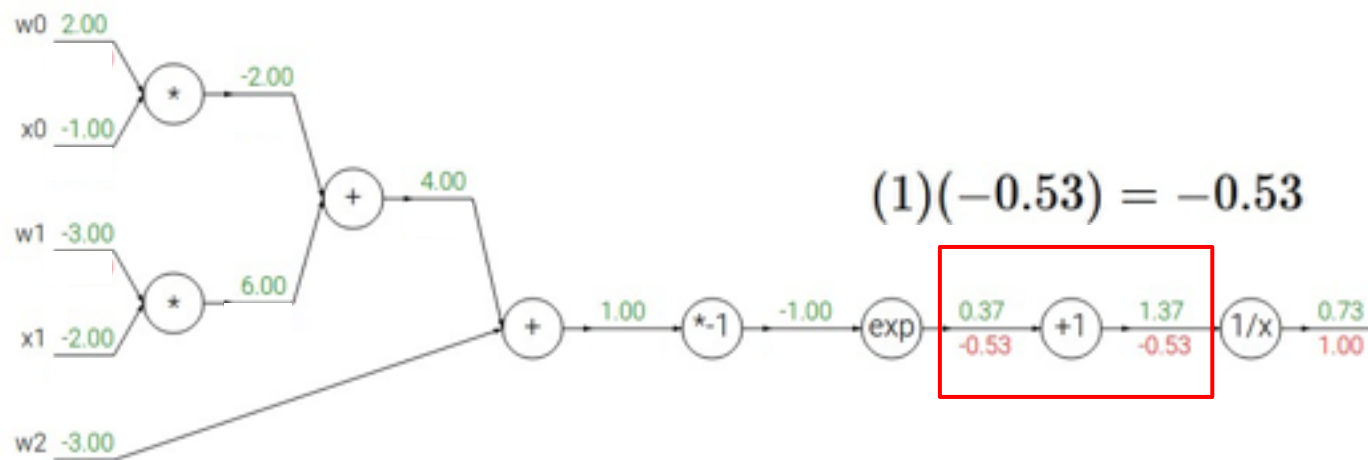
$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$



$f(x) = e^x$	\rightarrow	$\frac{df}{dx} = e^x$		$f(x) = \frac{1}{x}$	\rightarrow	$\frac{df}{dx} = -1/x^2$
$f_a(x) = ax$	\rightarrow	$\frac{df}{dx} = a$		<div style="border: 1px solid red; padding: 5px;">$f_c(x) = c + x$</div>	\rightarrow	<div style="border: 1px solid red; padding: 5px;">$\frac{df}{dx} = 1$</div>

Another example:

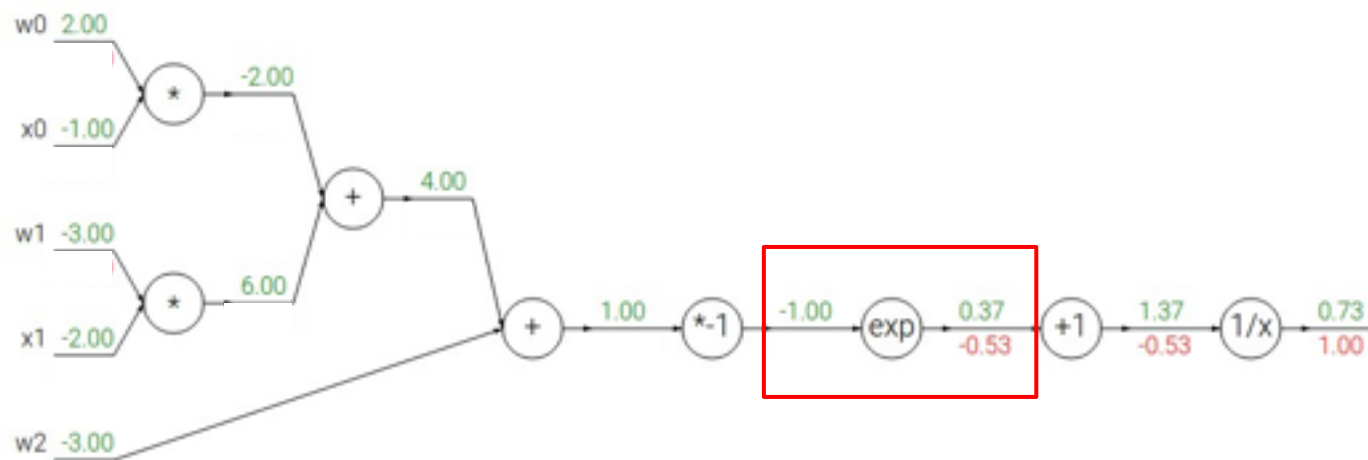
$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$



$f(x) = e^x$	\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:

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



$$f(x) = e^x$$

→

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

→

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

→

$$\frac{df}{dx} = -1/x^2$$

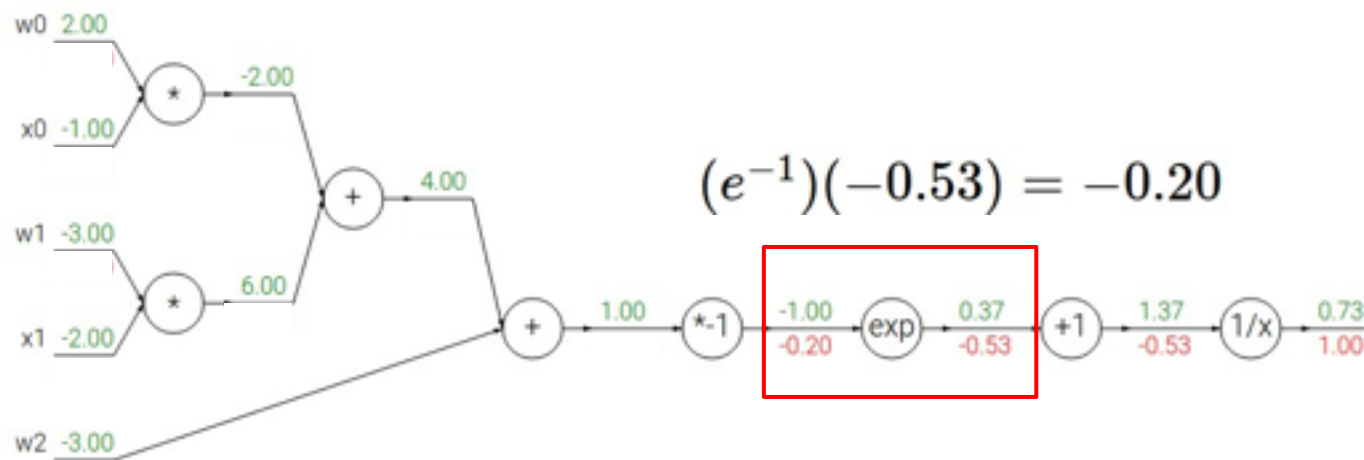
$$f_c(x) = c + x$$

→

$$\frac{df}{dx} = 1$$

Another example:

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



$$f(x) = e^x$$

→

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

→

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

→

$$\frac{df}{dx} = -1/x^2$$

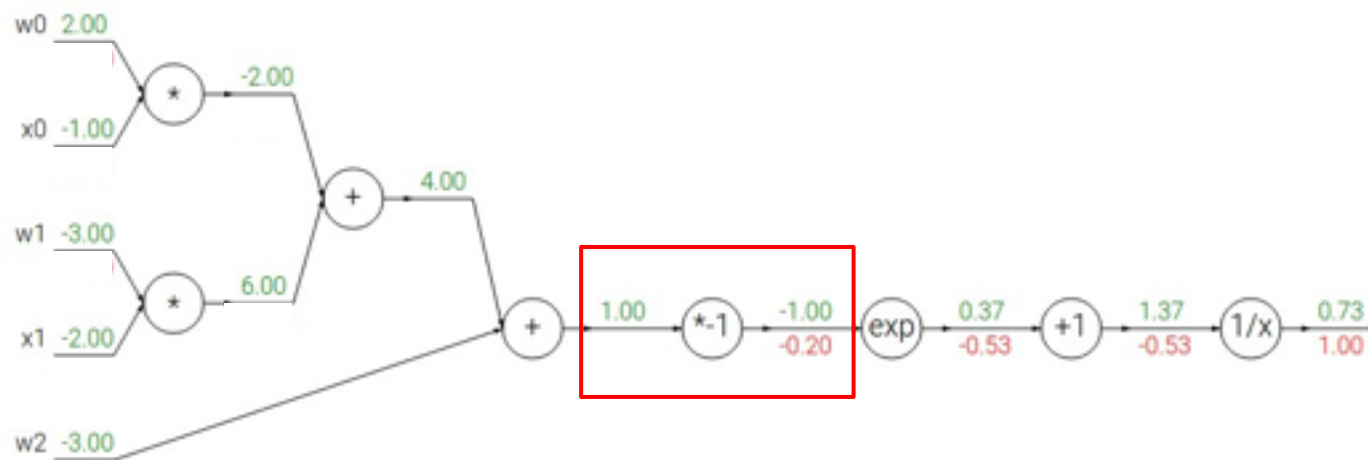
$$f_c(x) = c + x$$

→

$$\frac{df}{dx} = 1$$

Another example:

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



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

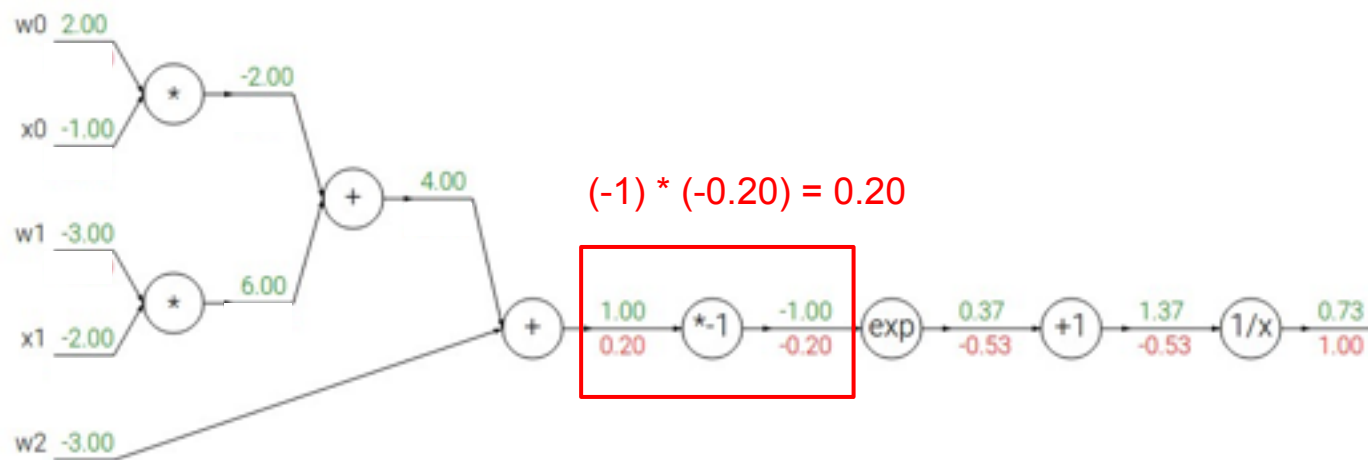
$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

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

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

Another example:

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



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

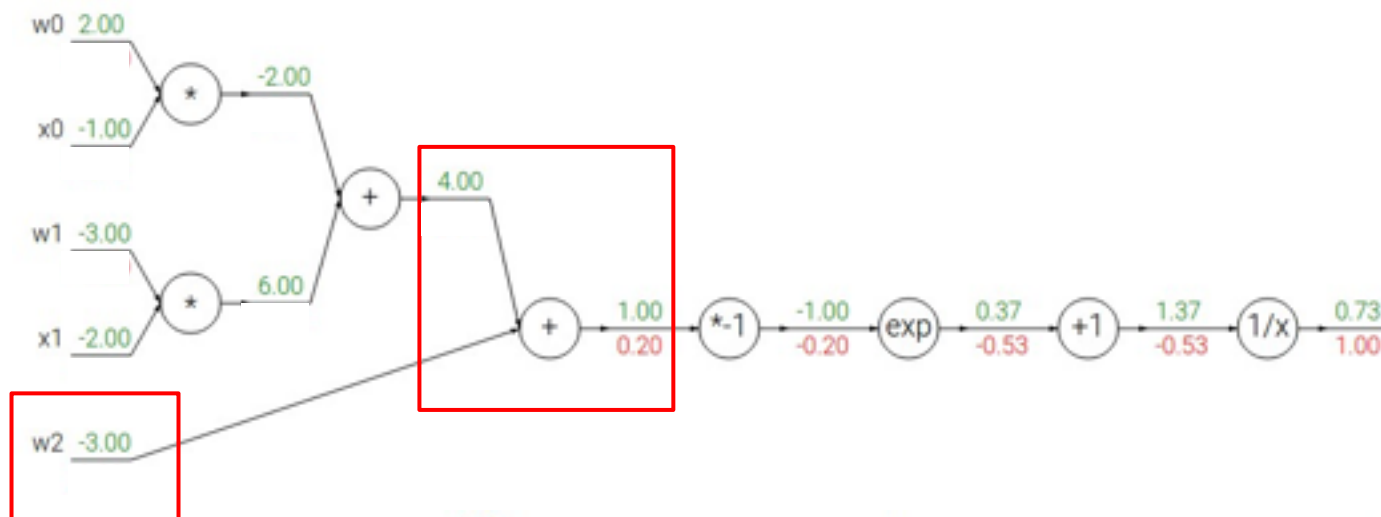
$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

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

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

Another example:

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



$$f(x) = e^x$$

→

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

→

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

→

$$\frac{df}{dx} = -1/x^2$$

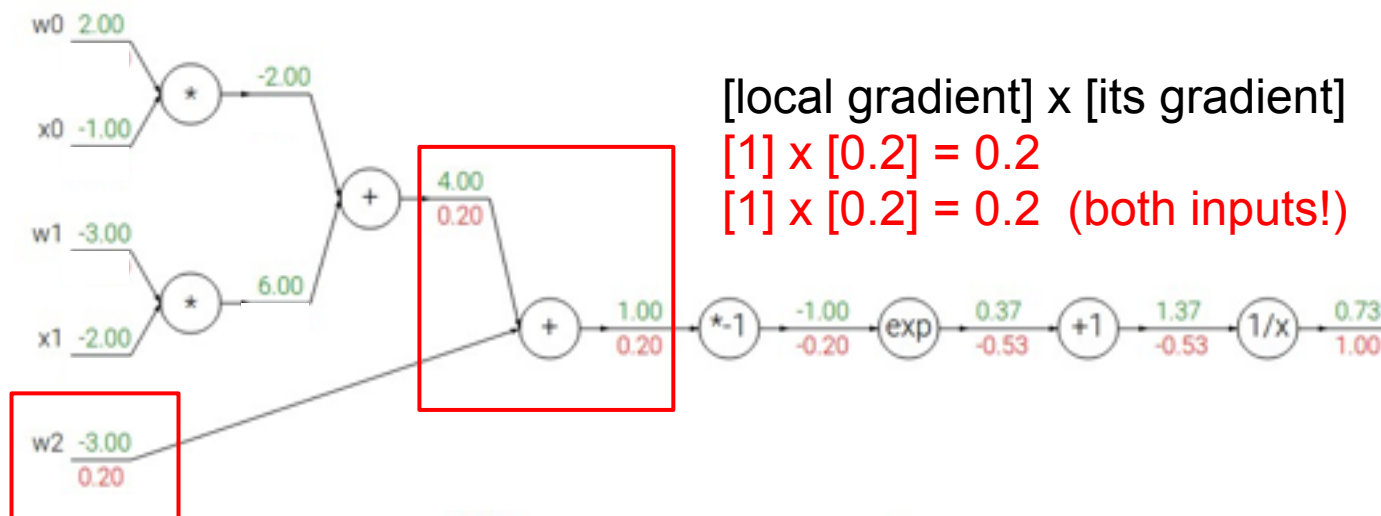
$$f_c(x) = c + x$$

→

$$\frac{df}{dx} = 1$$

Another example:

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



$$f(x) = e^x$$

→

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

→

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

→

$$\frac{df}{dx} = -1/x^2$$

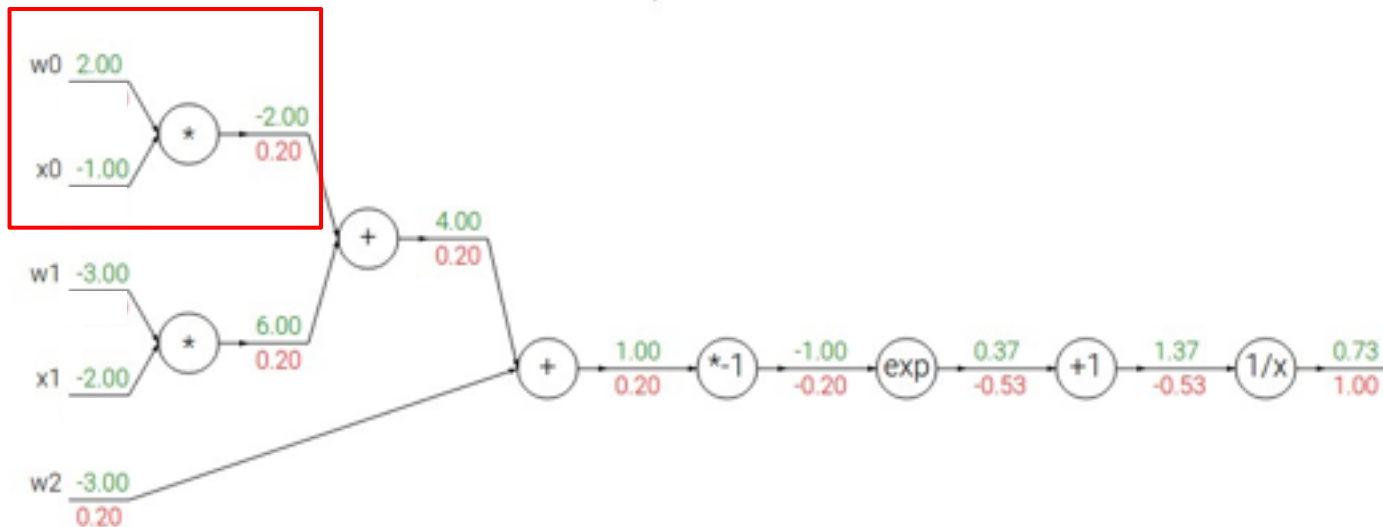
$$f_c(x) = c + x$$

→

$$\frac{df}{dx} = 1$$

Another example:

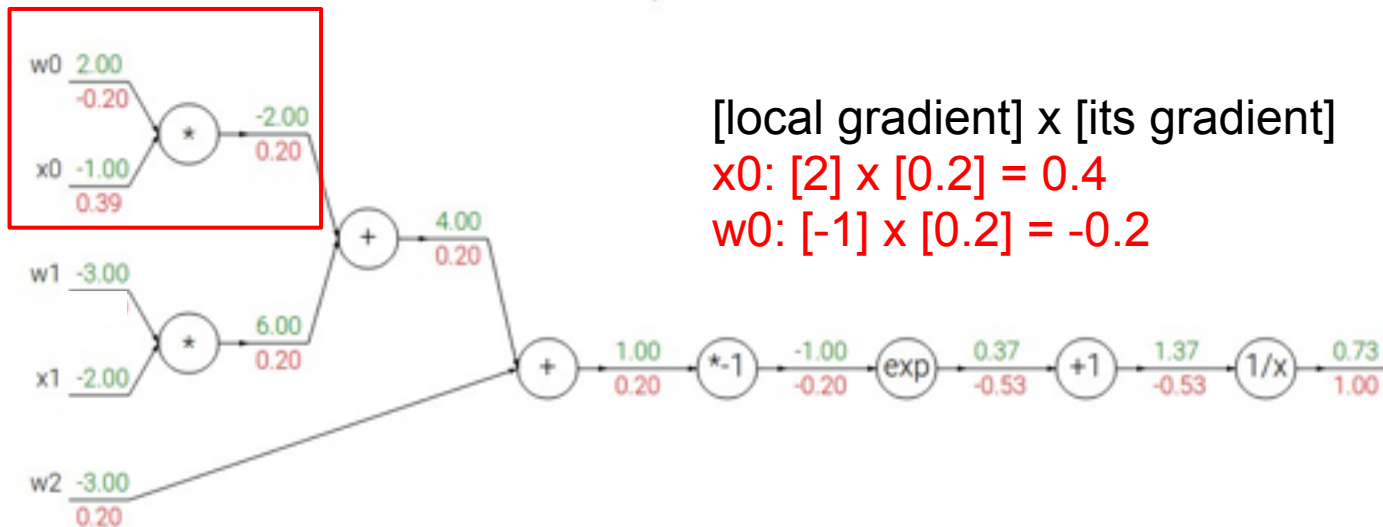
$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$



$f(x) = e^x$	\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:

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



[local gradient] x [its gradient]

$$x_0: [2] \times [0.2] = 0.4$$

$$w_0: [-1] \times [0.2] = -0.2$$

$$f(x) = e^x$$

→

$$\frac{df}{dx} = e^x$$

$$f(x) = \frac{1}{x}$$

→

$$\frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax$$

→

$$\frac{df}{dx} = a$$

$$f_c(x) = c + x$$

→

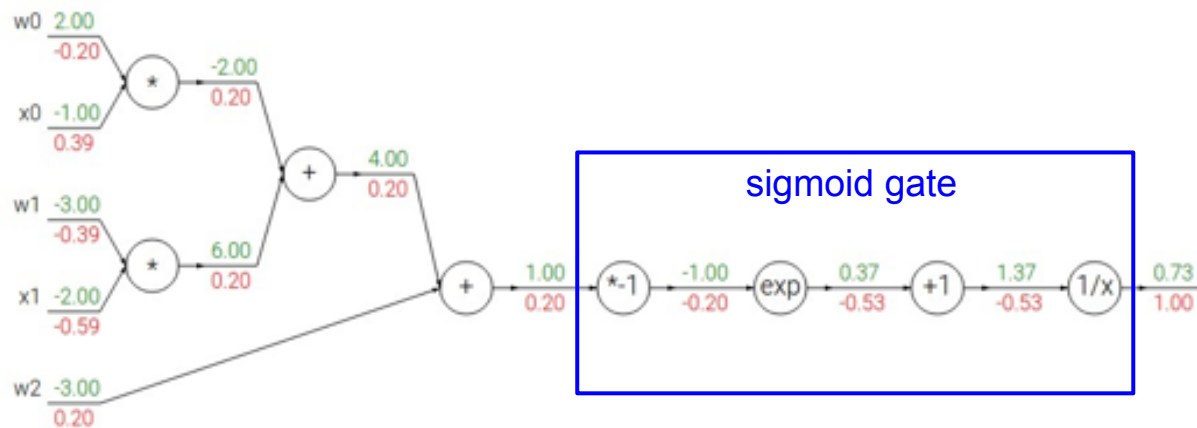
$$\frac{df}{dx} = 1$$

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2 x_2)}}$$

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

sigmoid function

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

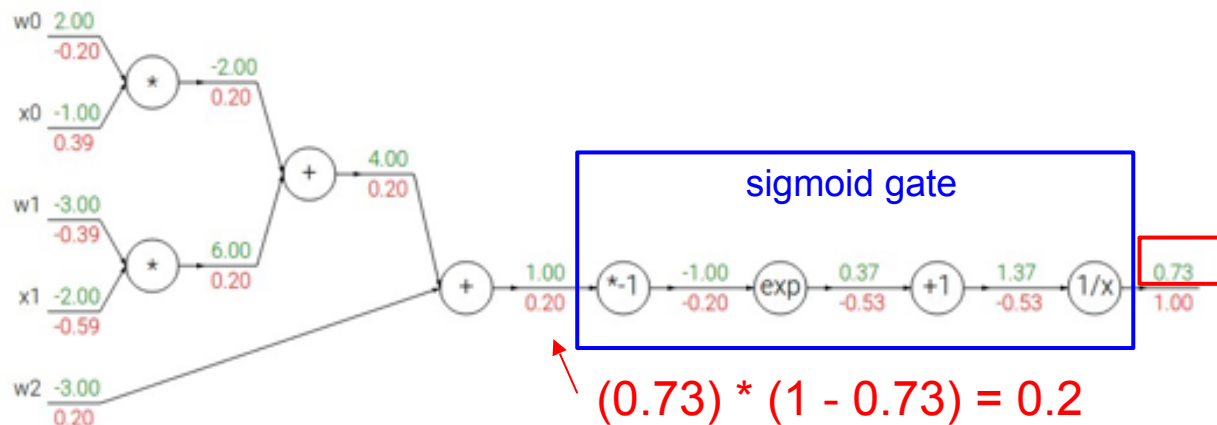


$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2 x_2)}}$$

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

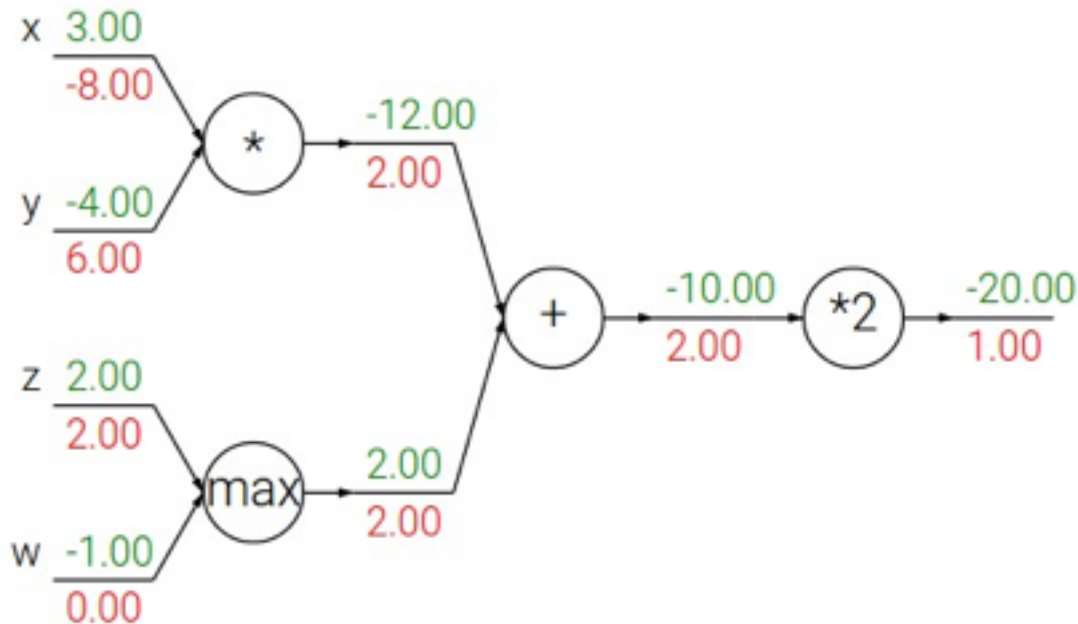
sigmoid function

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

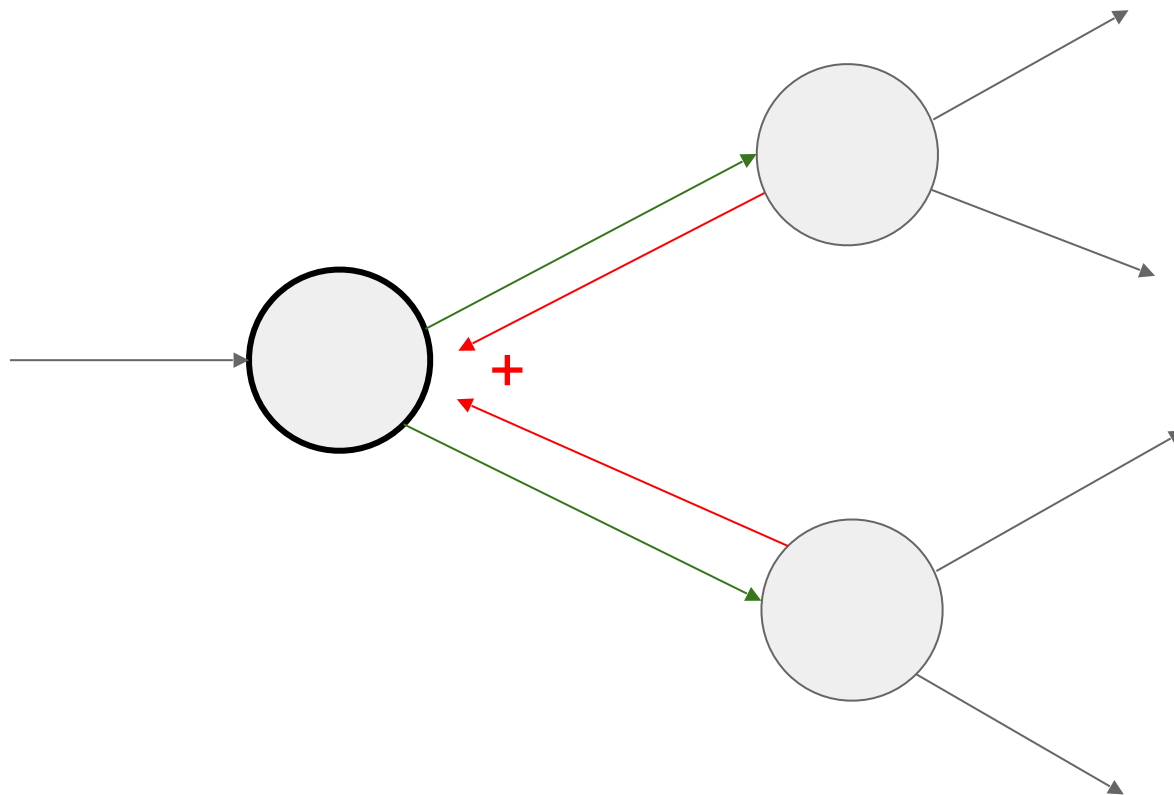


Patterns in backward flow

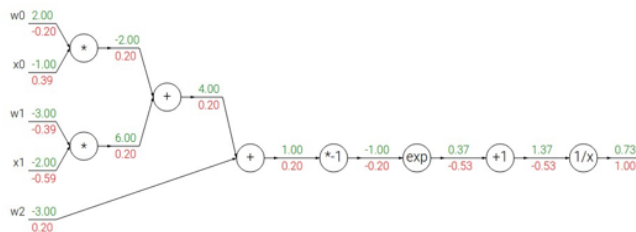
add gate: gradient distributor
max gate: gradient router
mul gate: gradient... “switcher”?



Gradients add at branches



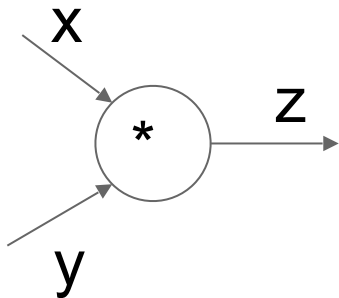
Implementation: forward/backward API



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
```

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]
```

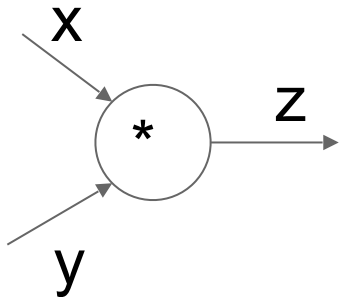
$$\frac{\partial L}{\partial z}$$

An arrow points from this box to the 'dz' parameter in the backward method signature of the code block above.

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

An arrow points from this box to the 'dx' element in the return list of the backward method in the code block above.

Implementation: forward/backward API



```
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]
```

(x,y,z are scalars)



[illegible][illegible]

ModuleTable.lua	cancel unused variable and useless expression
Module.lua	Revert "Don't re-flatten parameters if they are already flattened"
Mul.lua	removing the requirement for providing size in nn.Mul
MulConstant.lua	Ignore updateGradInput if self.gradInput is nil
MultiCriterion.lua	asserts in MultiCriterion and ParallelCriterion add
MultiLabelMarginCriterion.lua	initial revamp of torch7 tree
MultiMarginCriterion.lua	multimargin supports p=2
Narrow.lua	typeAs in Narrow not done in place.
NarrowTable.lua	NarrowTable
Normalize.lua	Remove brnn and baddbrnn from Normalize, because they allocate memory, ...
PRReLU.lua	Buffers for PReLU cuda implementation.
Padding.lua	fixed broken nn.Padding: input was returned in backprop
PoincareDistance.lua	Merge pull request #532 from xweng/master
Parallel.lua	fix a bug in conditional expression
ParallelCriterion.lua	asserts in MultiCriterion and ParallelCriterion add
ParallelTable.lua	Parallel optimization. ParallelTable inherits Container. unit tests
Power.lua	Use UNIX line endings
README.md	doc readthedocs
RRReLU.lua	Add randomized leaky rectified linear unit (RRReLU)
ReLU.lua	adds in-place ReLU and fixes a potential divide-by-zero in nn.Sqrt
Replicate.lua	Replicate batchSize
Reshape.lua	Added more informative pretty-printing.
Select.lua	initial revamp of torch7 tree
SelectTable.lua	nn.Module preserve type sharing semantics (#187); add nn.Module apply
Sequential.lua	fixing Sequential.remove corner case

Example: Torch Layers

[illegible][illegible]

* Original slides borrowed from Andrej Karpathy and Li Fei-Fei, Stanford cs231n

comp150dl

Example: Torch MulConstant

$$f(X) = aX$$

```
1 local MulConstant, parent = torch.class('nn.MulConstant', 'nn.Module')
2
3 function MulConstant:__init(constant_scalar, ip)
4     parent.__init(self)
5     assert(type(constant_scalar) == 'number', 'input is not scalar!')
6     self.constant_scalar = constant_scalar
7
8     -- default for inplace is false
9     self.inplace = ip or false
10    if (ip and type(ip) ~= 'boolean') then
11        error('in-place flag must be boolean')
12    end
13 end
```

initialization

```
14
15 function MulConstant:updateOutput(input)
16     if self.inplace then
17         input:mul(self.constant_scalar)
18         self.output = input
19     else
20         self.output:resizeAs(input)
21         self.output:copy(input)
22         self.output:mul(self.constant_scalar)
23     end
24     return self.output
25 end
```

forward()

```
26
27 function MulConstant:updateGradInput(input, gradOutput)
28     if self.gradInput then
29         if self.inplace then
30             gradOutput:mul(self.constant_scalar)
31             self.gradInput = gradOutput
32             -- restore previous input value
33             input:div(self.constant_scalar)
34         else
35             self.gradInput:resizeAs(gradOutput)
36             self.gradInput:copy(gradOutput)
37             self.gradInput:mul(self.constant_scalar)
38         end
39         return self.gradInput
40     end
41 end
```

backward()

Example: Caffe Layers

[illegible]

* Original slides borrowed from Andrej Karpathy and Li Fei-Fei, Stanford cs231n

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 #ifndef CPU_ONLY
41 STUB_GPU(SigmoidLayer);
42 #endif
43
44 INSTANTIATE_CLASS(SigmoidLayer);
45
46 } // namespace caffe
```

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

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

```

1 #include <vector>
2 #include <random>
3 #include <iterator>
4 #include "caffe/layers/shuffle_layer.hpp"
5 #include "caffe/util/math_functions.hpp"
6
7 namespace caffe {
8
9 template <typename Dtype>
10 void Shuffle(Dtype* bottom_data, Dtype* top_data, const int item_size,
11             const bool forward, const int* shuffle_order, const int count) {
12     // data_shape is expected to be the shape (count, M) of the blob
13     // data in bottom_data and top_data
14     for (int i = 0; i < count; ++i) {
15         for (int j = 0; j < item_size; ++j) {
16             if (forward) {
17                 top_data[i*item_size+j] = bottom_data[i*item_size+shuffle_order[j]];
18             } else {
19                 bottom_data[i*item_size+shuffle_order[j]] = top_data[i*item_size+j];
20             }
21         }
22     }
23 }
24
25 template <typename Dtype>
26 void ShuffleLayer<Dtype>::LayerSetUp(const vector<Blob<Dtype>*>& bottom,
27                                     const vector<Blob<Dtype>*>& top) {
28     // Check there is only one bottom layer
29     CHECK_EQ(bottom.size(), 1);
30     // TODO: extend functionality to > 2-D blobs, but for now only 2D works
31     // CHECK_EQ(bottom[0]->shape().size(), 2);
32
33     shuffle_seed_ = rand();
34     // calculate count of each item in batch
35     batch_item_size_ = bottom[0]->count(1, bottom[0]->CanonicalAxisIndex(-1));
36
37     vector<int> shuffle_order;
38     // Make a vector of ordered inds
39     for (int i=0; i<bottom[0]->shape(1); i++) shuffle_order.push_back(i);
40     // Mersenne twister initialized with input seed
41     std::mt19937 gen(shuffle_seed_);
42     std::shuffle(shuffle_order.begin(), shuffle_order.end(), gen);
43
44     // copy randomized shuffle order to layer member variable shuffle_order_
45     shuffle_order_.Reshape(shuffle_order.size(), 1, 1, 1);
46
47     for (int i = 0; i < shuffle_order.size(); i++) {
48         shuffle_order_.mutable_cpu_data()[shuffle_order_.offset(i)] = shuffle_order[i];
49     }
50 }

```

```

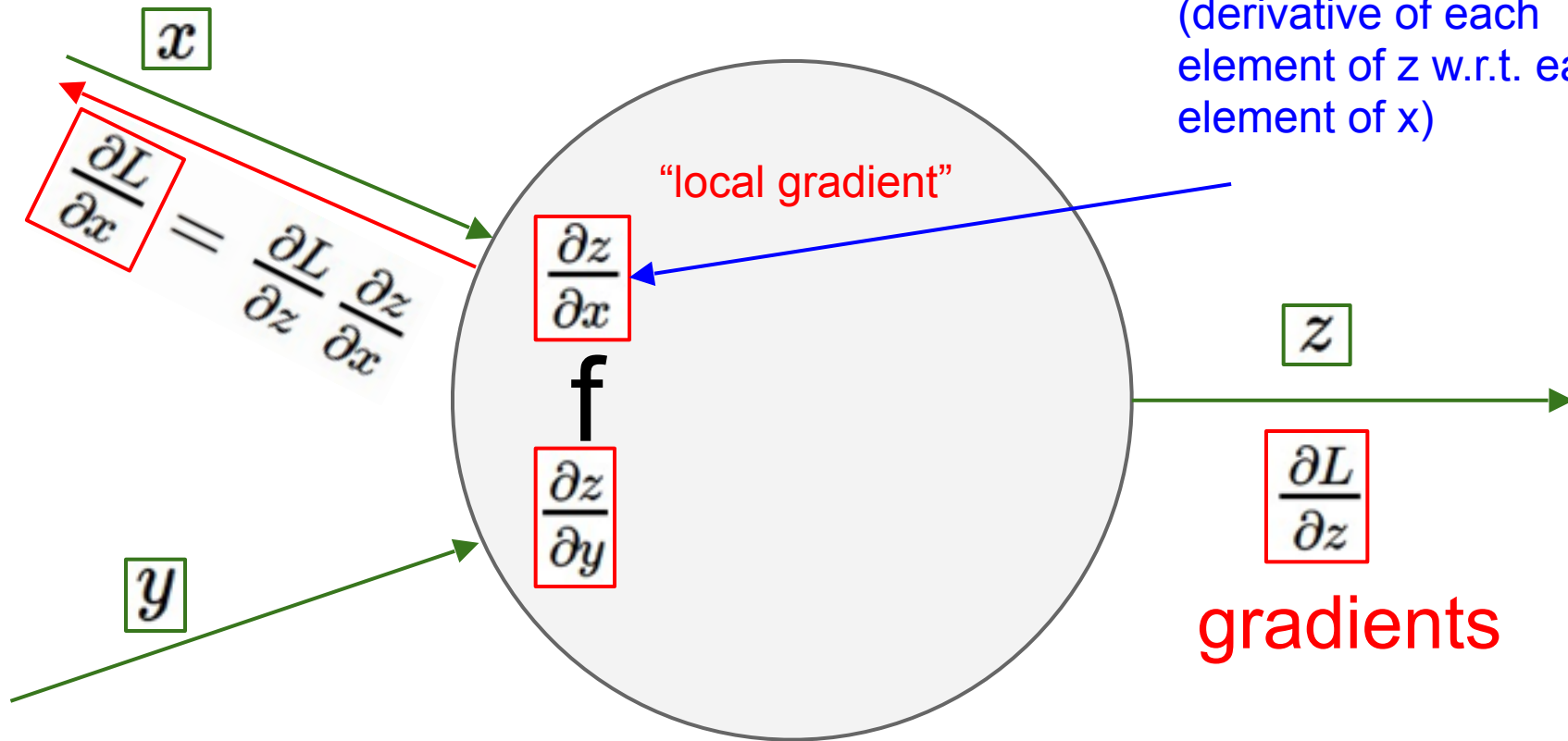
62 template <typename Dtype>
63 void ShuffleLayer<Dtype>::Forward_cpu(const vector<Blob<Dtype>*>& bottom,
64                                     const vector<Blob<Dtype>*>& top) {
65     const int count = top[0]->num();
66     std::cout << bottom[0]->shape(0) << std::endl;
67     std::cout << bottom[0]->shape(1) << std::endl;
68     vector<int> orig_shape = bottom[0]->shape();
69     vector<int> new_shape; new_shape.push_back(count); new_shape.push_back(batch_item_size_);
70     bottom[0]->Reshape(new_shape);
71     top[0]->Reshape(new_shape);
72     Dtype* bottom_data = bottom[0]->mutable_cpu_data();
73     Dtype* top_data = top[0]->mutable_cpu_data();
74     const int* shuffle_order = shuffle_order_.cpu_data();
75
76     bool forward = true;
77     Shuffle(bottom_data, top_data, batch_item_size_, forward, shuffle_order,
78            count);
79     top[0]->Reshape(orig_shape);
80     bottom[0]->Reshape(orig_shape);
81 }
82
83 template <typename Dtype>
84 void ShuffleLayer<Dtype>::Backward_cpu(const vector<Blob<Dtype>*>& top,
85                                       const vector<bool>& propagate_down, const vector<Blob<Dtype>*>& bottom) {
86     const int count = top[0]->num();
87     vector<int> orig_shape = bottom[0]->shape();
88     vector<int> new_shape; new_shape.push_back(count); new_shape.push_back(batch_item_size_);
89     bottom[0]->Reshape(new_shape);
90     top[0]->Reshape(new_shape);
91     Dtype* bottom_data = bottom[0]->mutable_cpu_data();
92     Dtype* top_data = top[0]->mutable_cpu_data();
93     const int* shuffle_order = shuffle_order_.cpu_data();
94
95     bool forward = false;
96     Shuffle(bottom_data, top_data, batch_item_size_, forward, shuffle_order,
97            count);
98     bottom[0]->Reshape(orig_shape);
99     top[0]->Reshape(orig_shape);
100 }
101
102 #ifdef CPU_ONLY
103 STUB_GPU(ShuffleLayer);
104 #endif
105
106 INSTANTIATE_CLASS(ShuffleLayer);
107 REGISTER_LAYER_CLASS(Shuffle);
108
109 } // namespace caffe

```

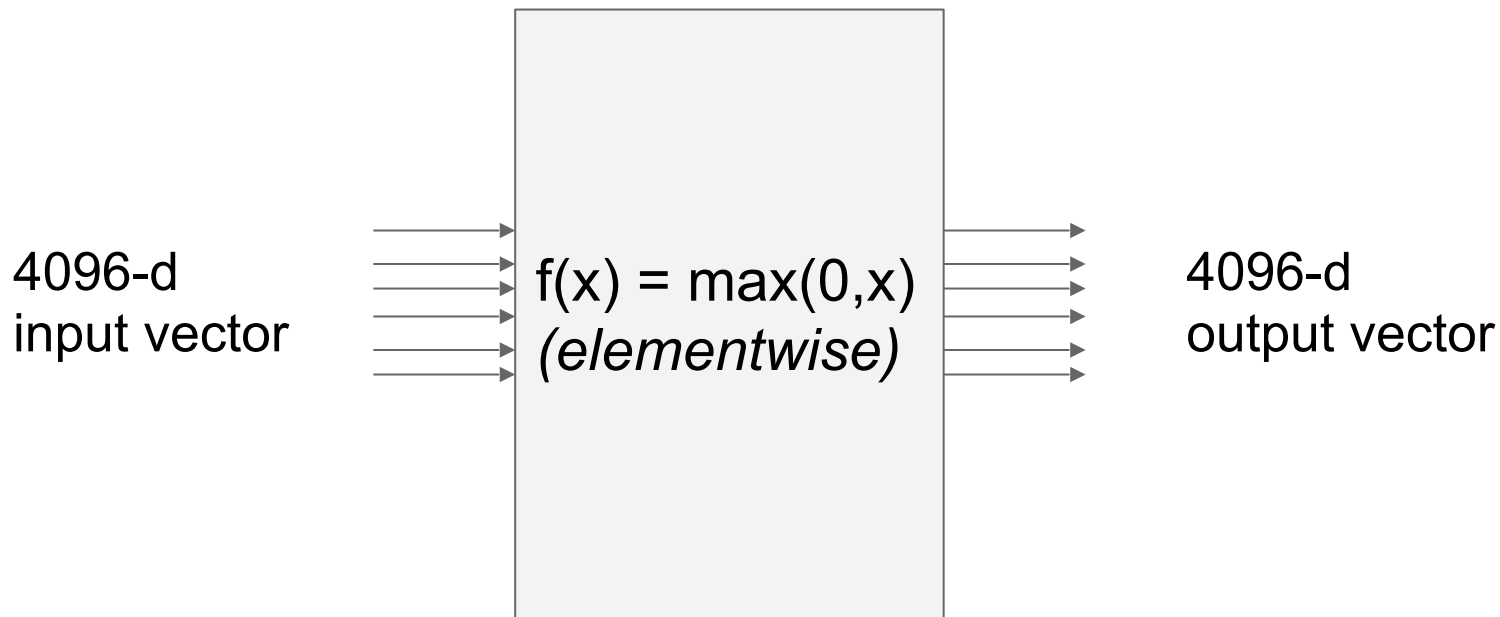
Gradients for vectorized code

(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)



Vectorized operations

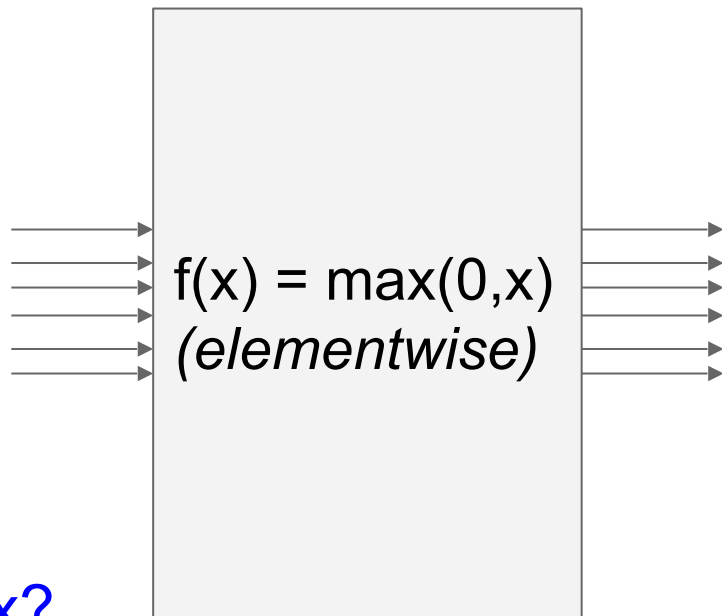


Vectorized operations

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

Jacobian matrix

4096-d
input vector



4096-d
output vector

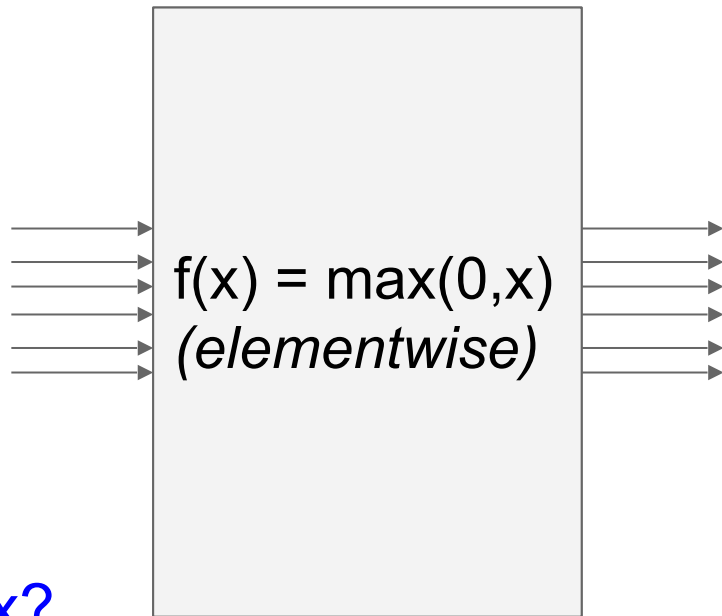
Q: what is the
size of the
Jacobian matrix?

Vectorized operations

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

Jacobian matrix

4096-d
input vector



4096-d
output vector

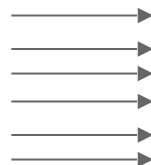
Q: what is the
size of the
Jacobian matrix?
[4096 x 4096!]

Q2: what does it
look like?

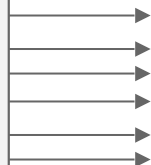
Vectorized operations

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

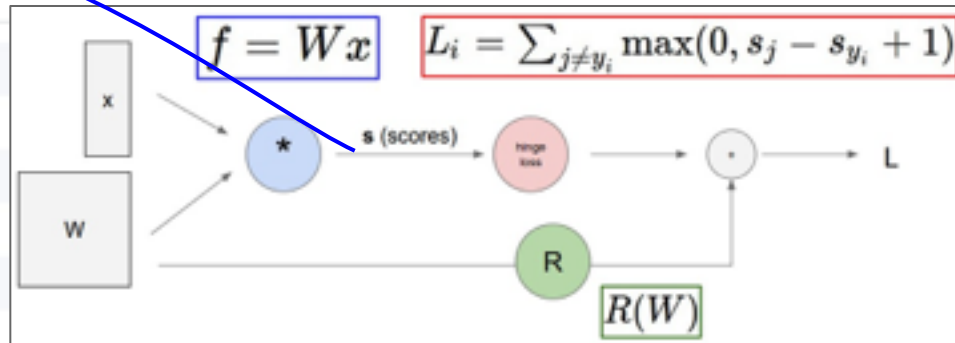
i.e. Jacobian would technically be a [409,600 x 409,600] matrix :\\

Assignment: Writing SVM/Softmax

Stage your forward/backward computation!

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 = #...
```



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()**.
- **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 so far:

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

Neural Network so far:

(**Before**) Linear score function:

$$f = Wx$$

(**Now**) 2-layer Neural Network

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

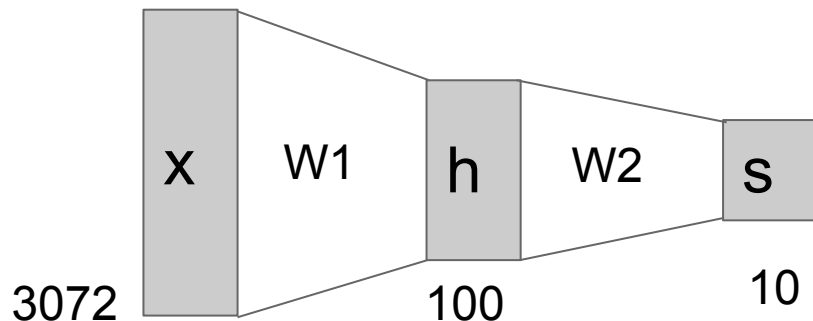
Neural Network so far:

(**Before**) Linear score function:

$$f = Wx$$

(**Now**) 2-layer Neural Network

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



Neural Network so far:

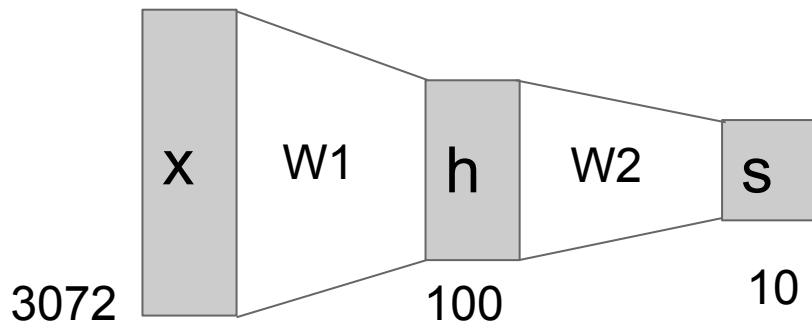


(**Before**) Linear score function:

$$f = Wx$$

(**Now**) 2-layer Neural Network

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



Neural Network so far:

(**Before**) Linear score function:

$$f = Wx$$

(**Now**) 2-layer Neural Network
or 3-layer Neural Network

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

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

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

```
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)
```

Forward pass

Backward
pass

backprop
of derivative

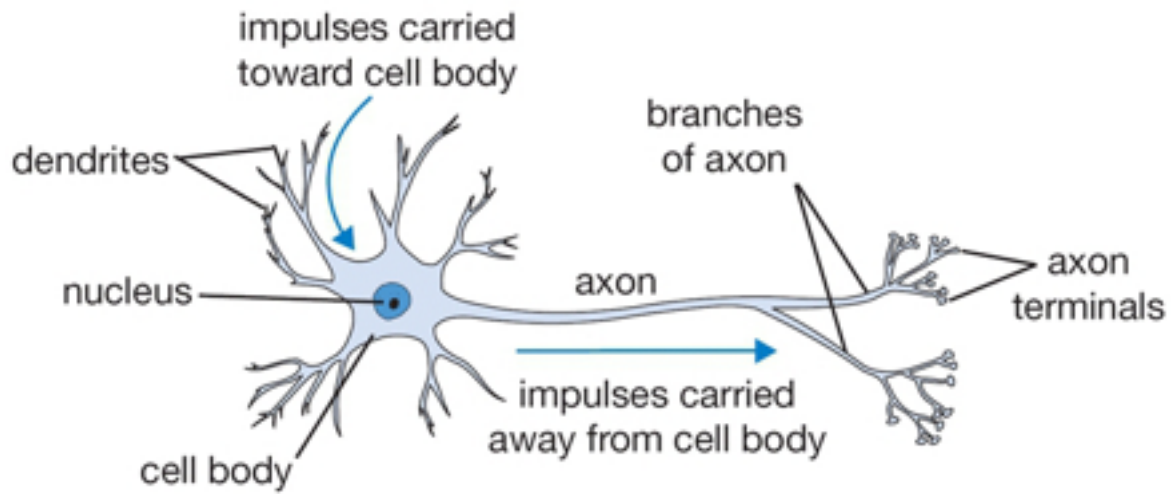
from @iamtrask, <http://iamtrask.github.io/2015/07/12/basic-python-network/>

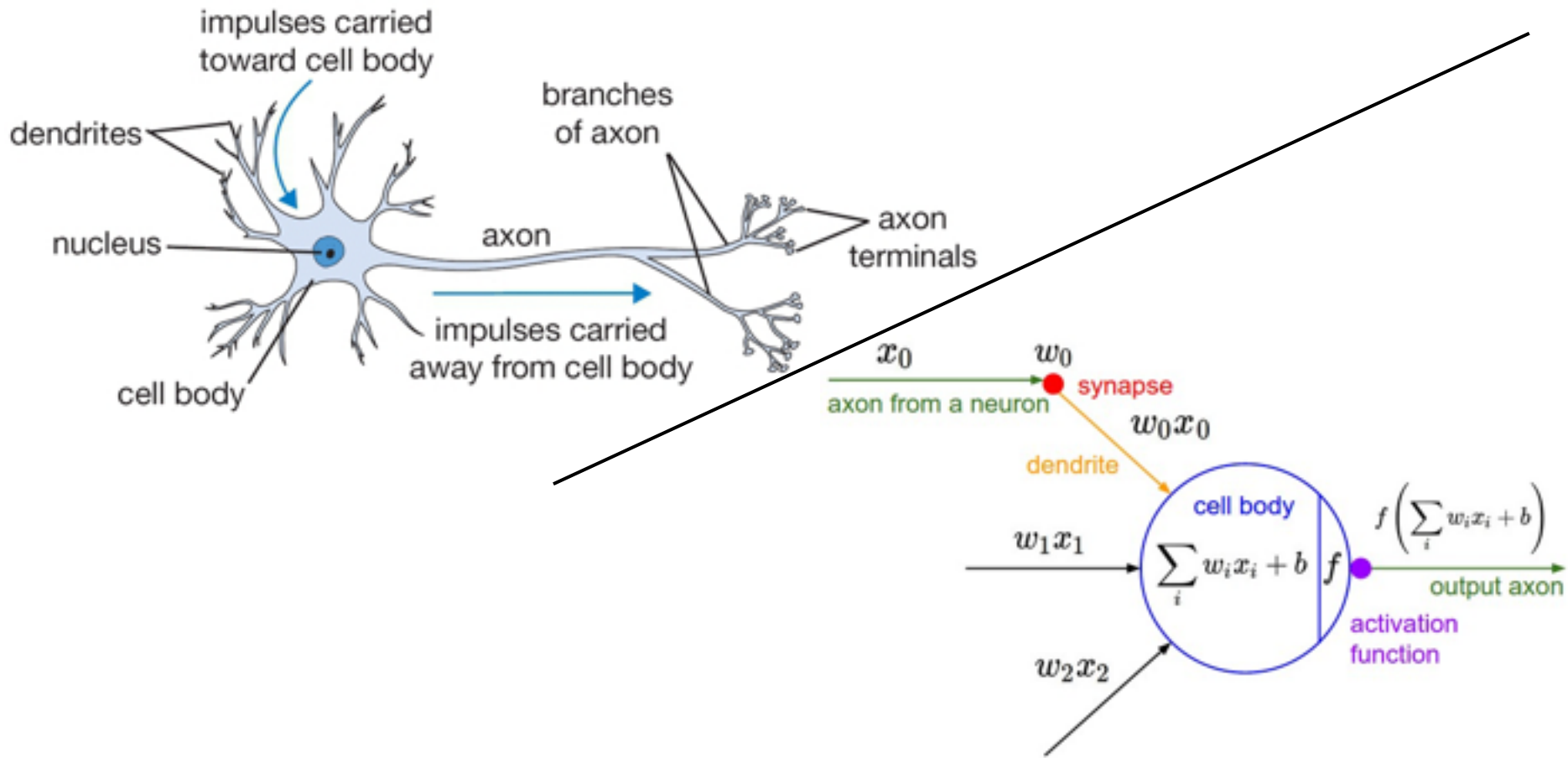
Assignment: Writing 2layer Net

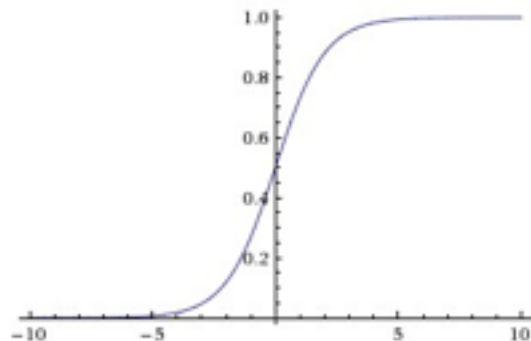
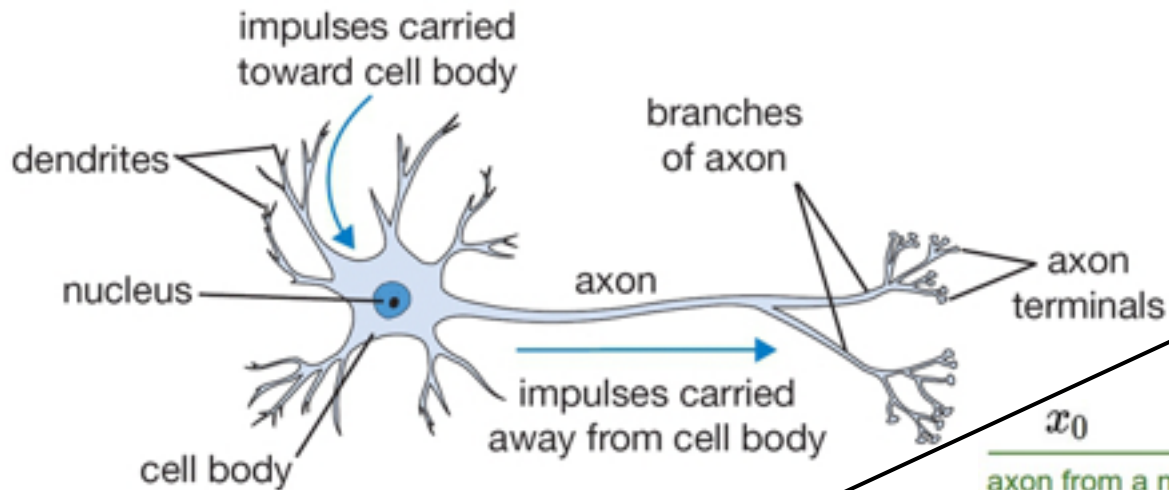
Stage your forward/backward computation!

```
# 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 = #...
```



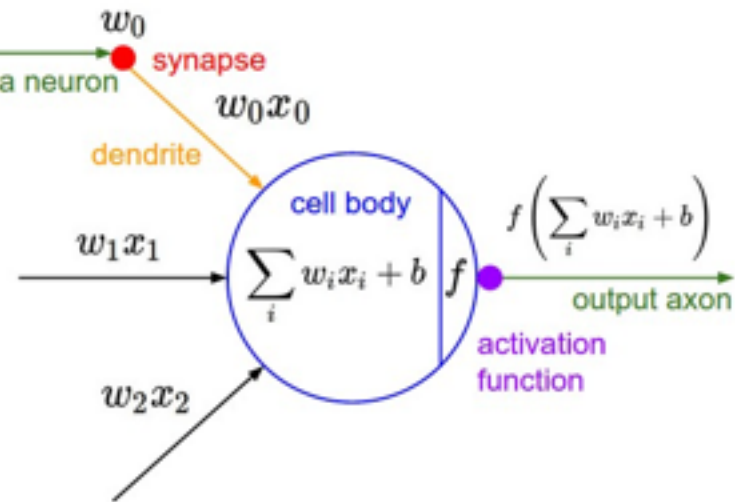


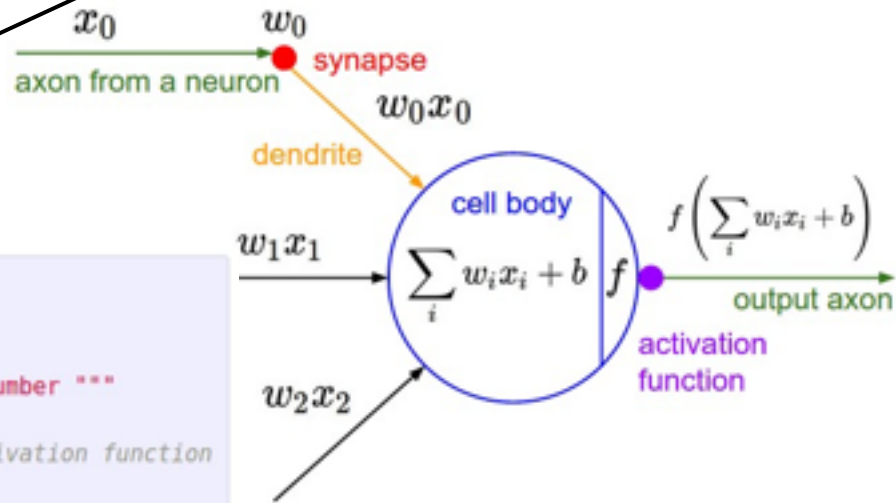
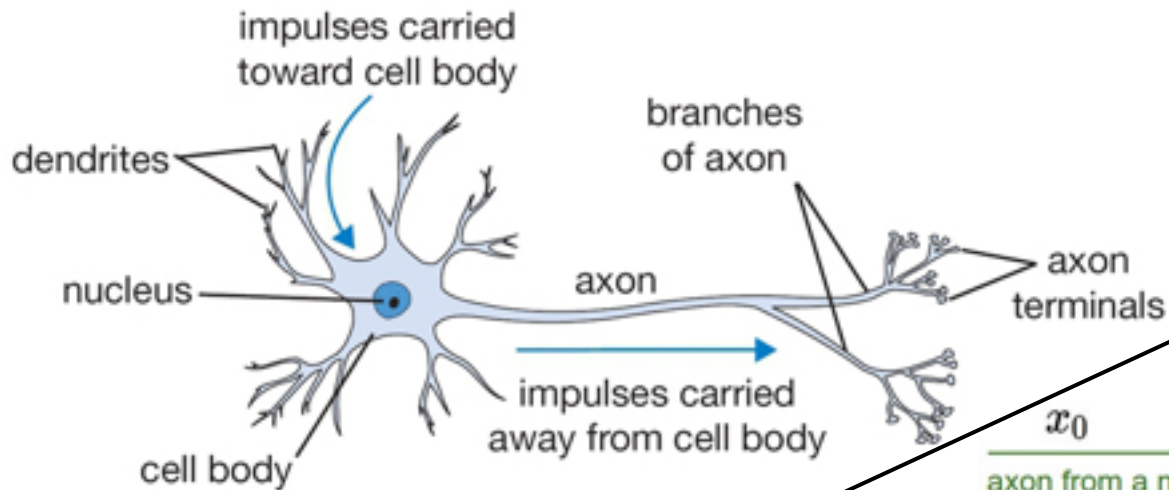




**sigmoid activation
function**

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



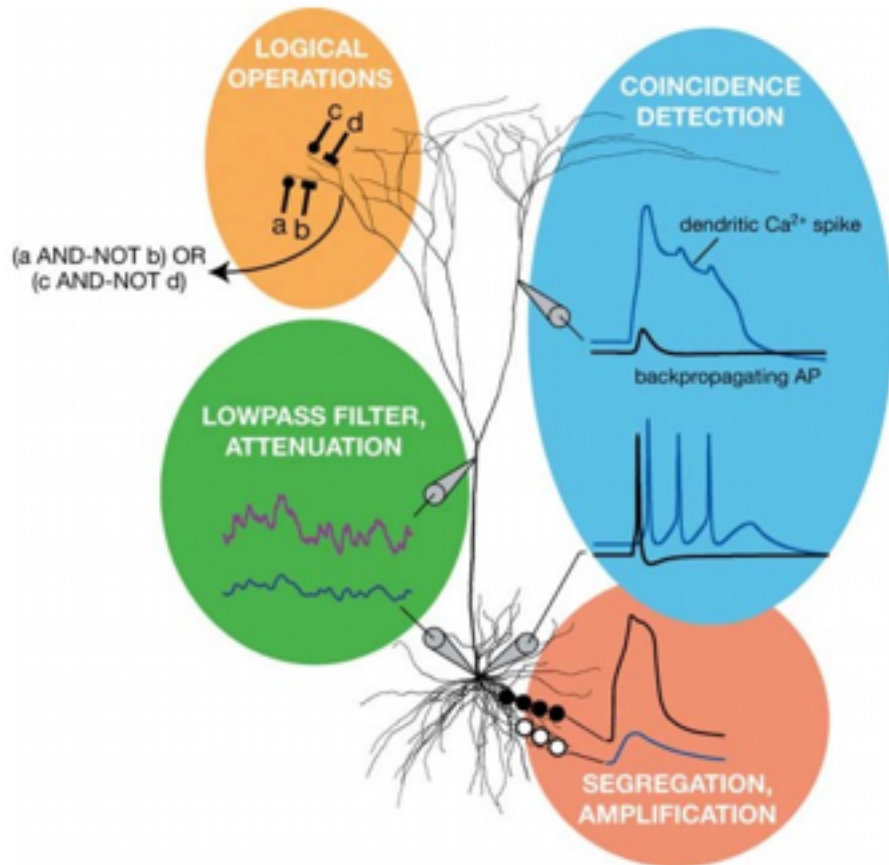


```
class Neuron:
    # ...
    def neuron_tick(inputs):
        """ assume inputs and weights are 1-D numpy arrays and bias is a number """
        cell_body_sum = np.sum(inputs * self.weights) + self.bias
        firing_rate = 1.0 / (1.0 + math.exp(-cell_body_sum)) # sigmoid activation function
  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


Be very careful with your Brain analogies:

Biological Neurons:

- Many different types
- Dendrites can perform complex non-linear computations
- Synapses are not a single weight but a complex non-linear dynamical system
- Rate code may not be adequate

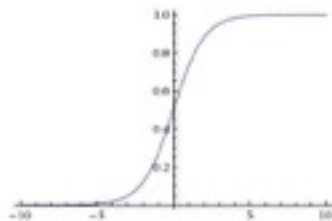


[Dendritic Computation. London and Hausser]

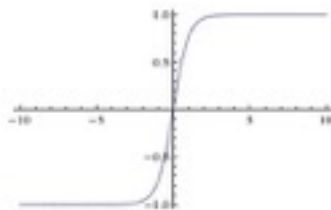
Activation Functions

Sigmoid

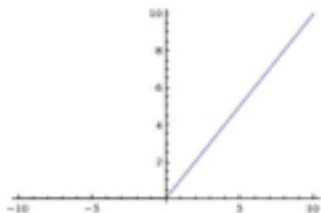
$$\sigma(x) = 1/(1 + e^{-x})$$



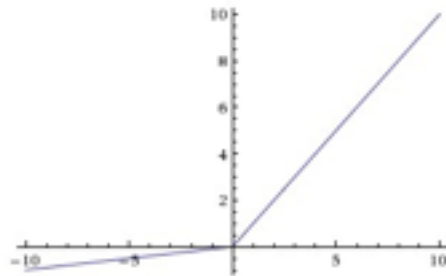
tanh tanh(x)



ReLU max(0,x)



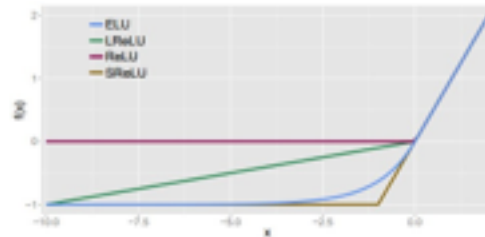
Leaky ReLU $\max(0.1x, x)$



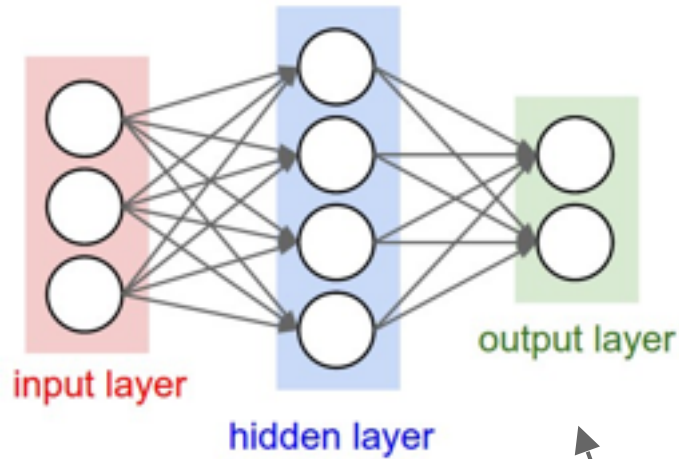
Maxout $\max(w_1^T x + b_1, w_2^T x + b_2)$

ELU

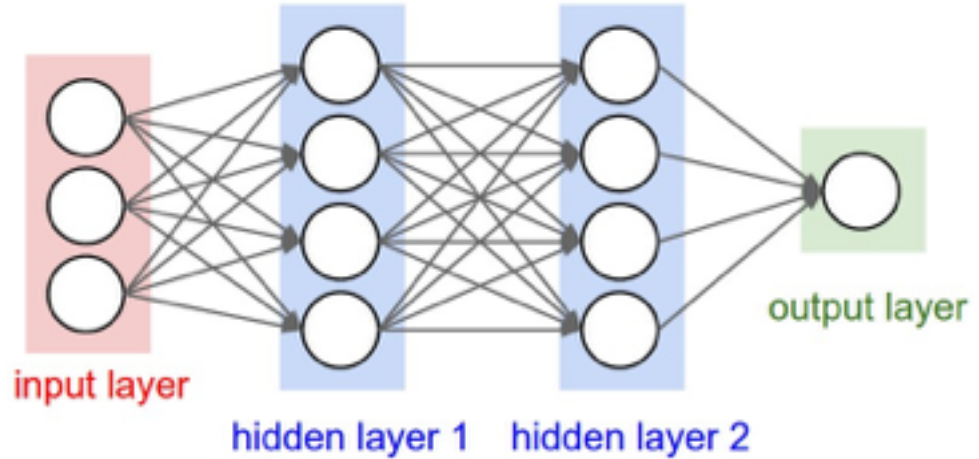
$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$



Neural Networks: Architectures



“2-layer Neural Net”, or
“1-hidden-layer Neural Net”



“3-layer Neural Net”, or
“2-hidden-layer Neural Net”

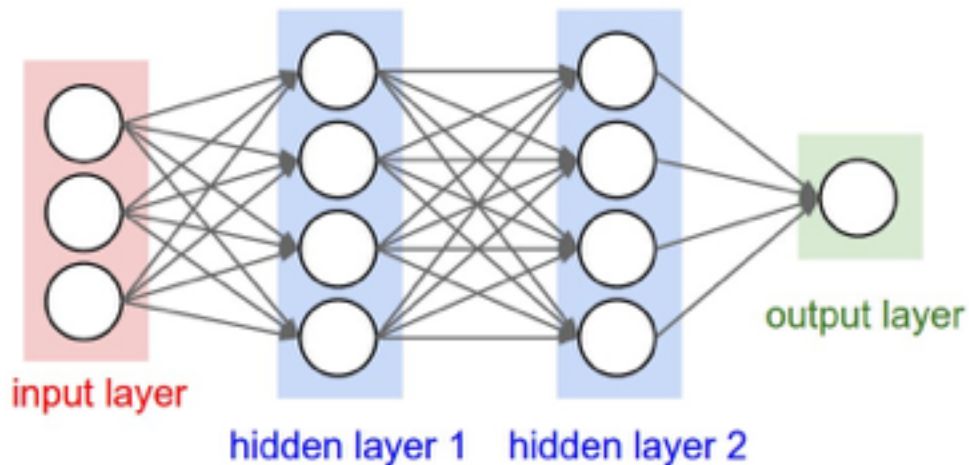
“Fully-connected” layers

Example Feed-forward computation of a Neural Network

```
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



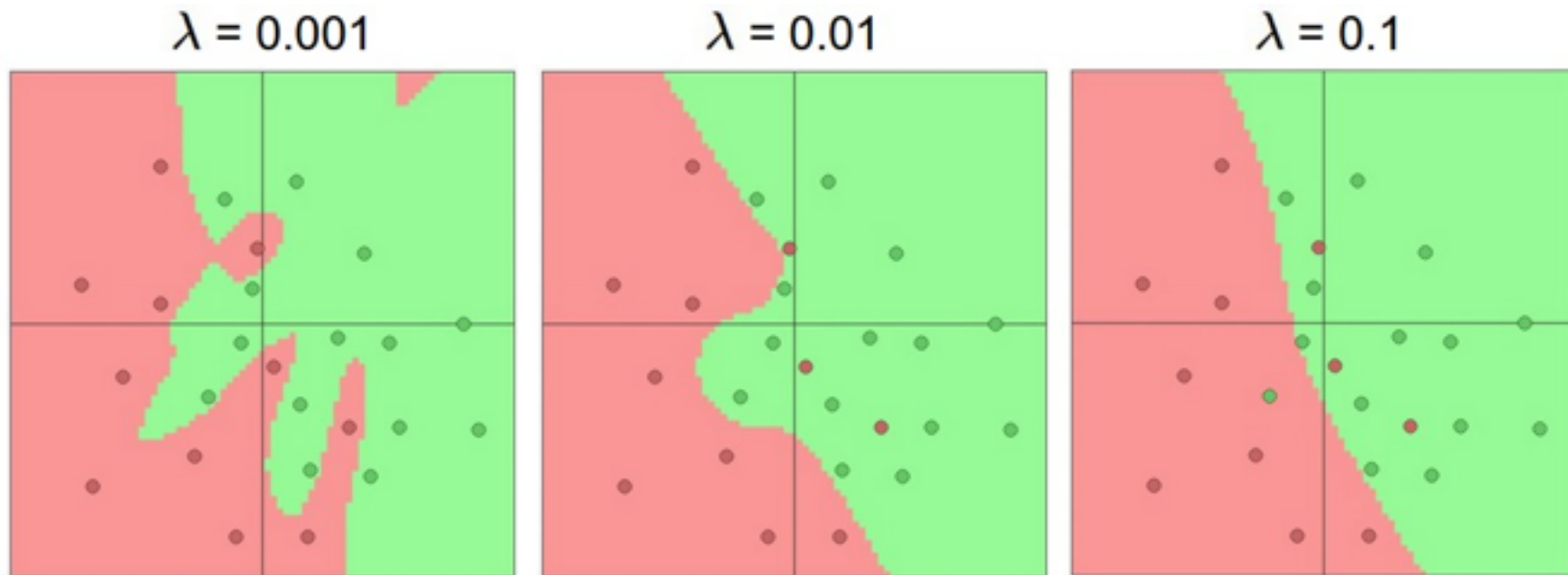
```
# forward-pass of a 3-layer neural network:  
f = lambda x: 1.0/(1.0 + np.exp(-x)) # activation function (use sigmoid)  
x = np.random.randn(3, 1) # random input vector of three numbers (3x1)  
h1 = f(np.dot(W1, x) + b1) # calculate first hidden layer activations (4x1)  
h2 = f(np.dot(W2, h1) + b2)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

Setting the number of layers and their sizes



↑
more neurons = more capacity

Do not use size of neural network as a regularizer. Use stronger regularization instead:



(you can play with this demo over at ConvNetJS: <http://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>)

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)

Next Lecture:

More than you ever wanted to know about Neural Networks and how to train them.