

# Lecture 4: Backpropagation and Neural Networks part 1

# Administrative

A1 is due Jan 20 (Wednesday). ~150 hours left

Warning: Jan 18 (Monday) is Holiday (no class/office hours)

Also note:

Lectures are non-exhaustive.

Read course notes for completeness.

I'll hold make up office hours on Wed Jan20, 5pm @ Gates 259

## Where we are...

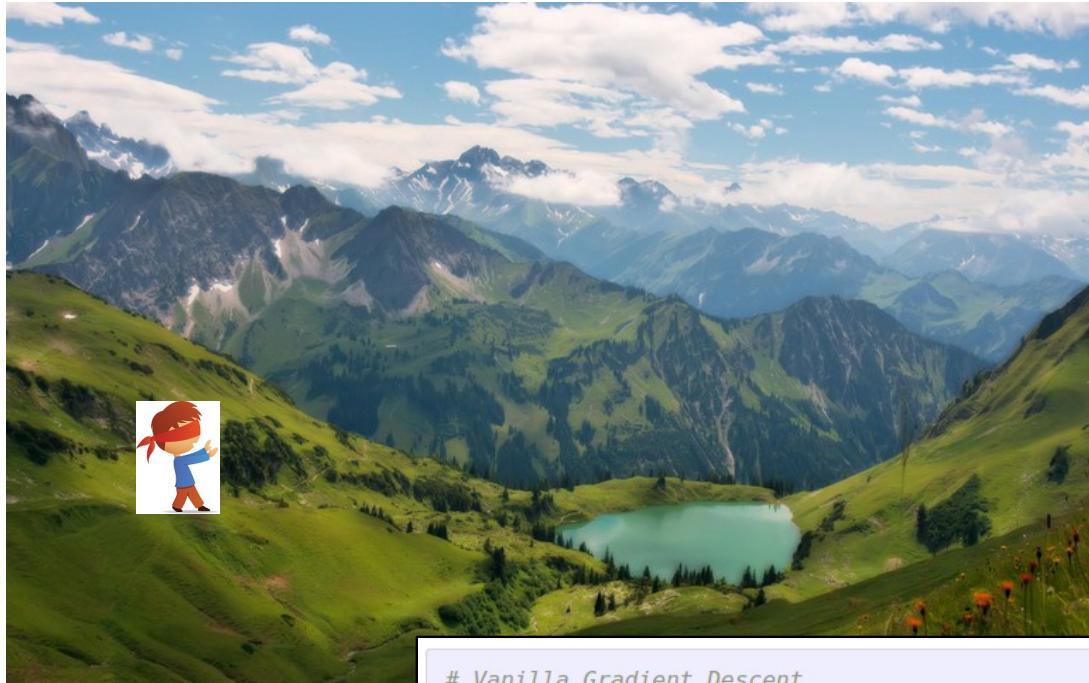
$$s = f(x; W) = Wx \quad \text{scores function}$$

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

$$L = \frac{1}{N} \sum_{i=1}^N L_i + \sum_k W_k^2 \quad \text{data loss + regularization}$$

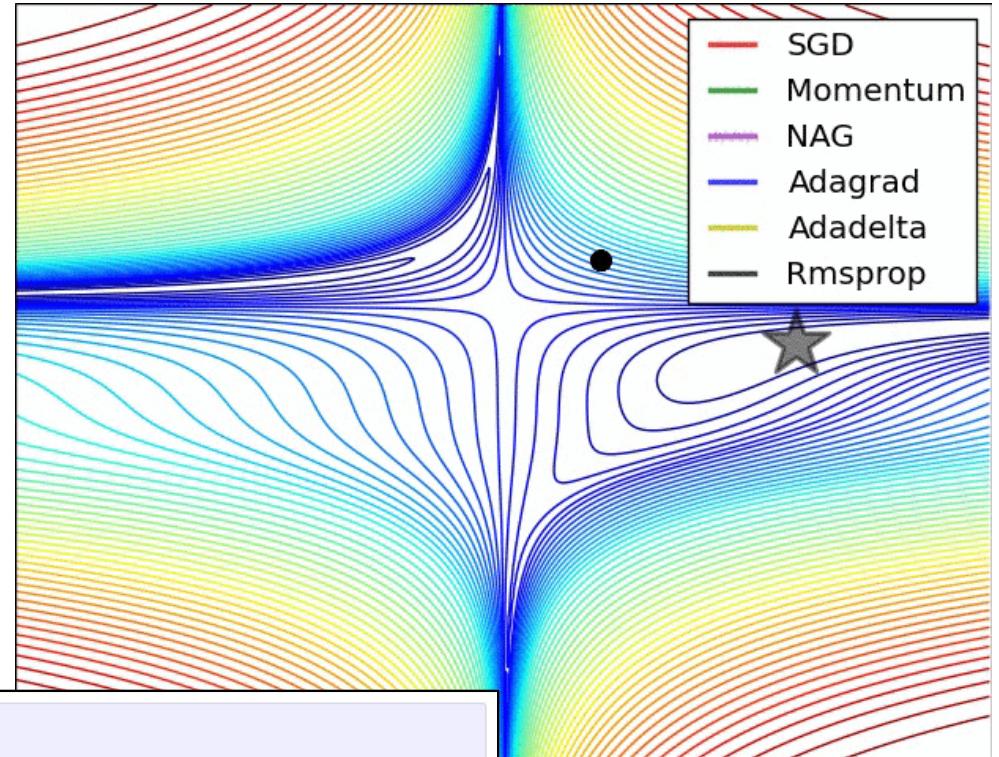
want  $\nabla_W L$

# Optimization



```
# Vanilla Gradient Descent

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



(image credits  
to Alec Radford)

# Gradient Descent

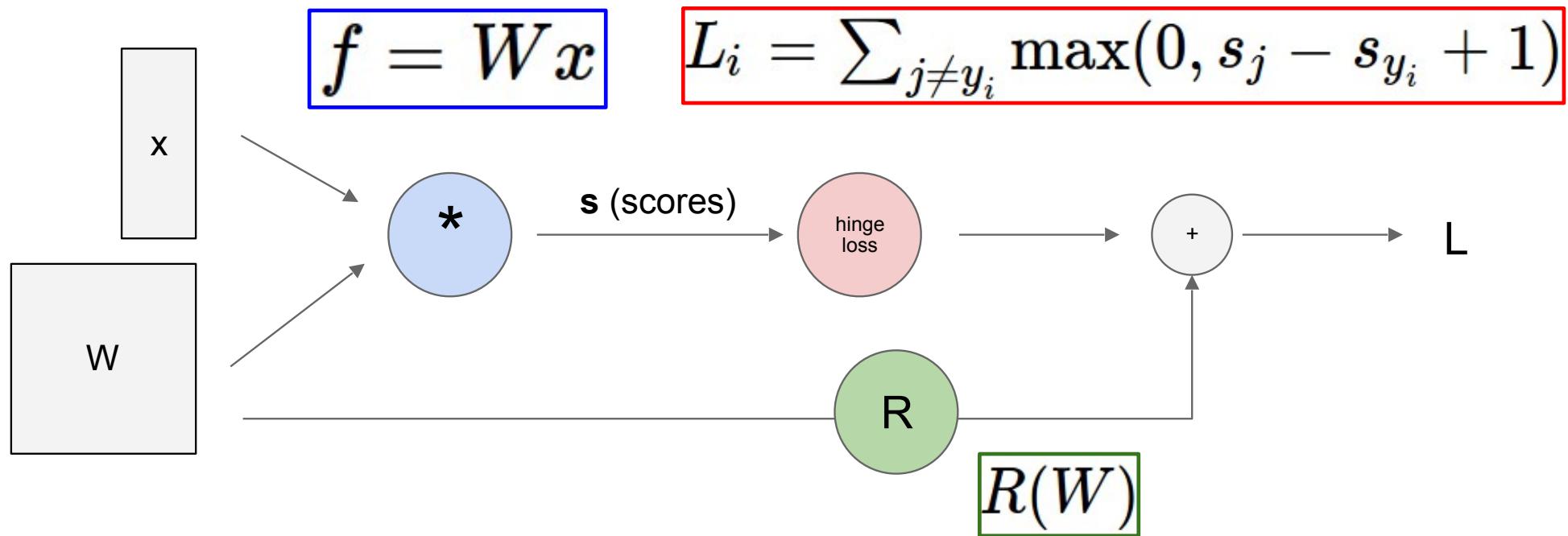
$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h}$$

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

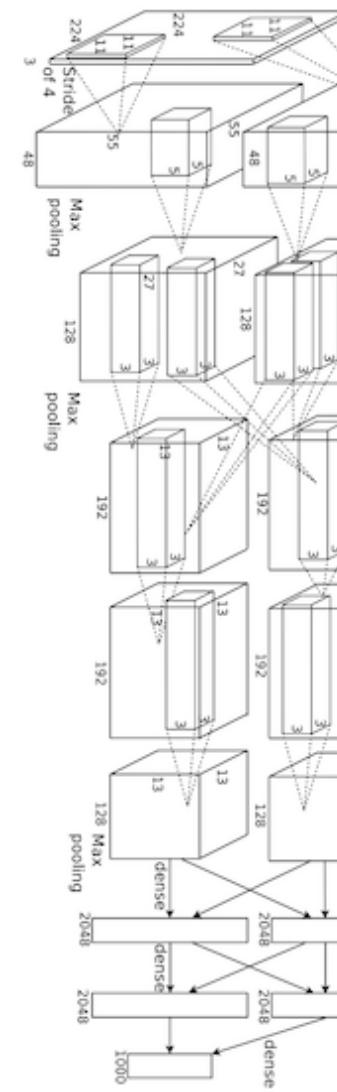
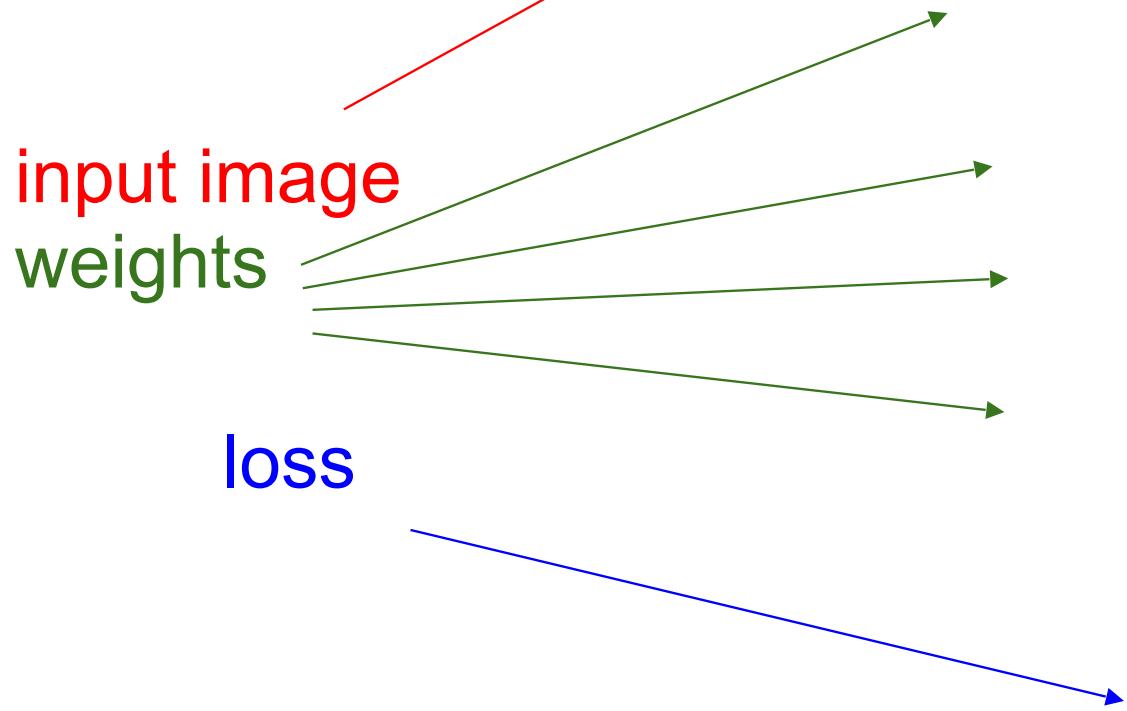
**Analytic gradient:** fast :), exact :), error-prone :(

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

# Computational Graph



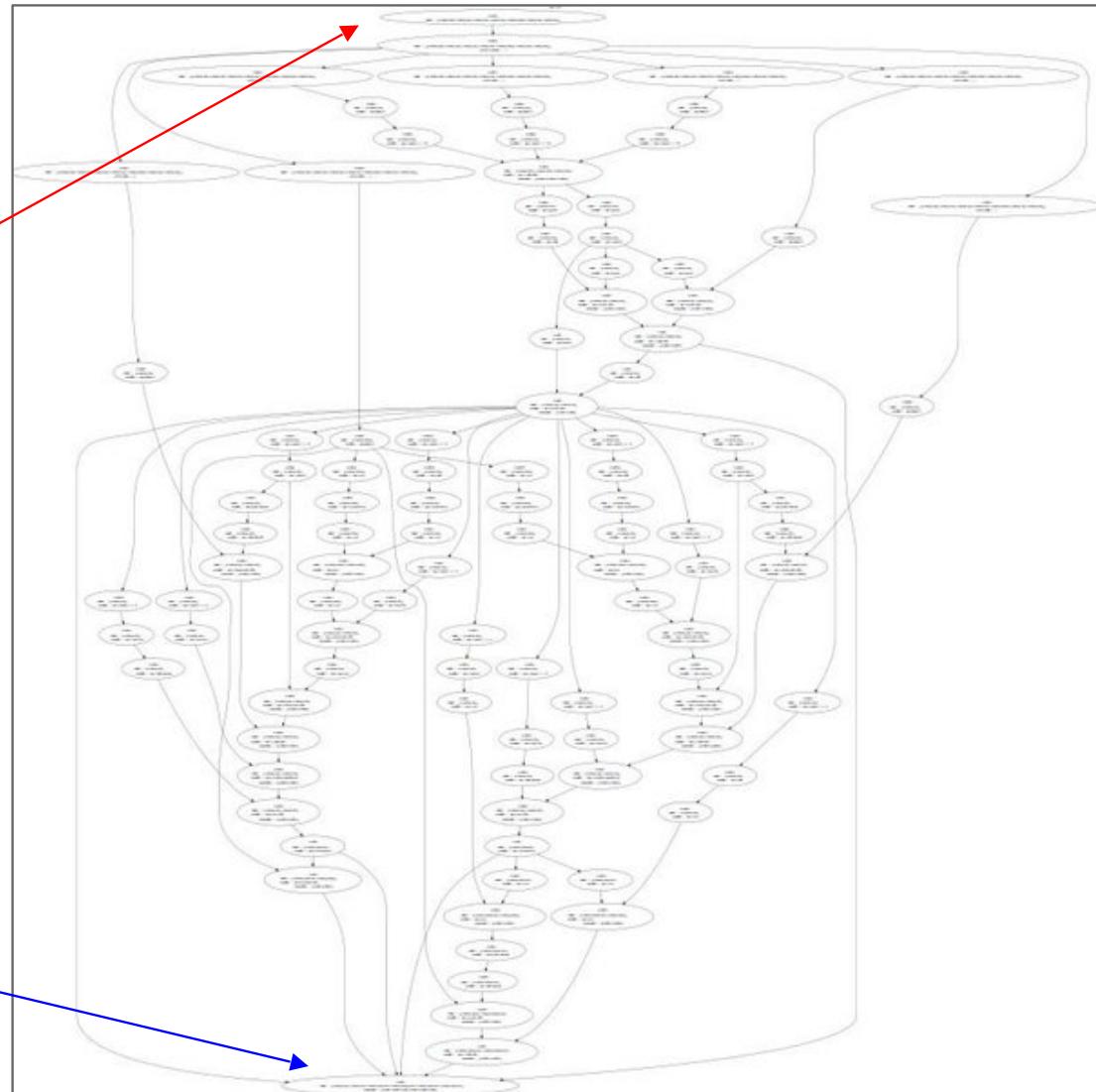
# Convolutional Network (AlexNet)



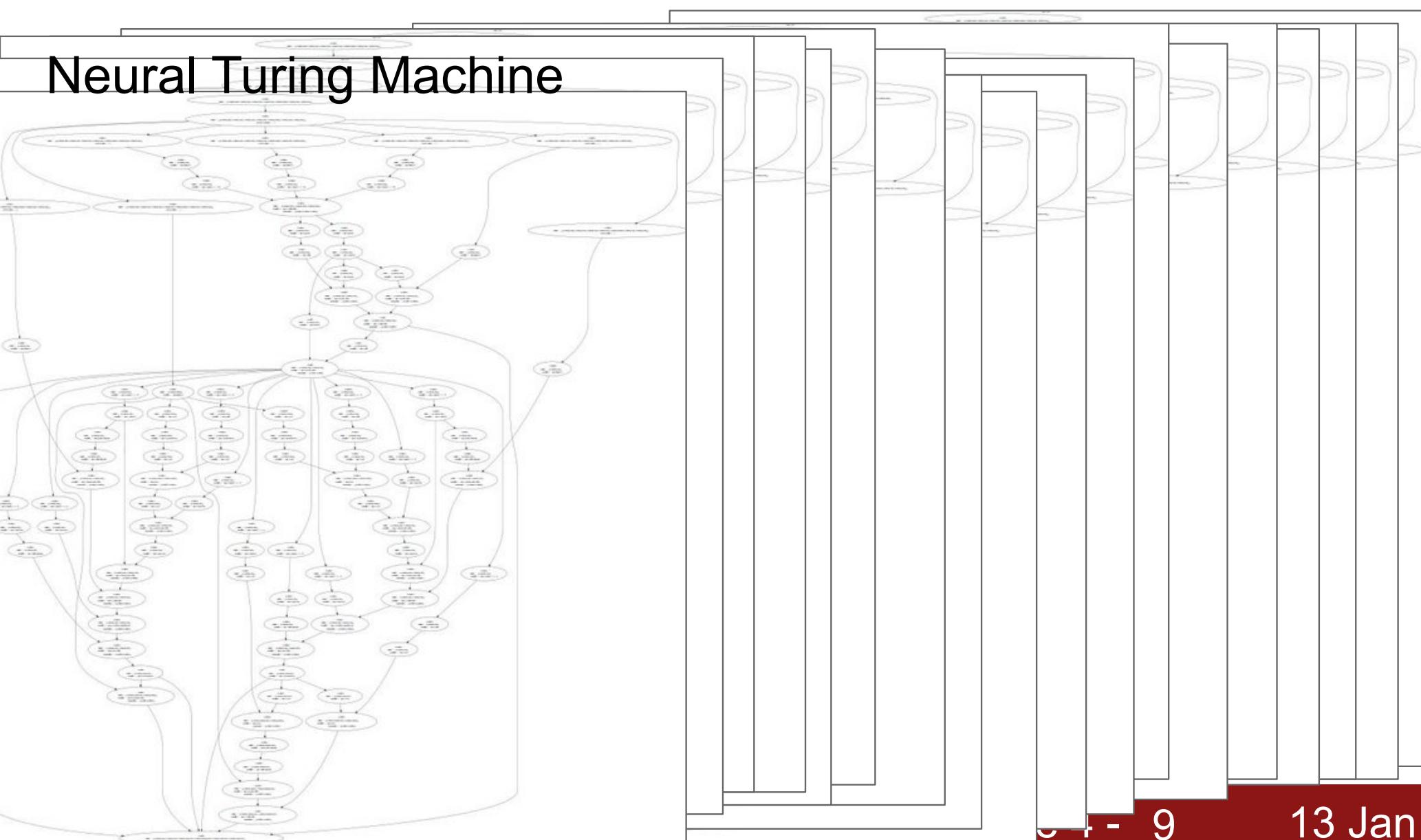
# Neural Turing Machine

input tape

loss

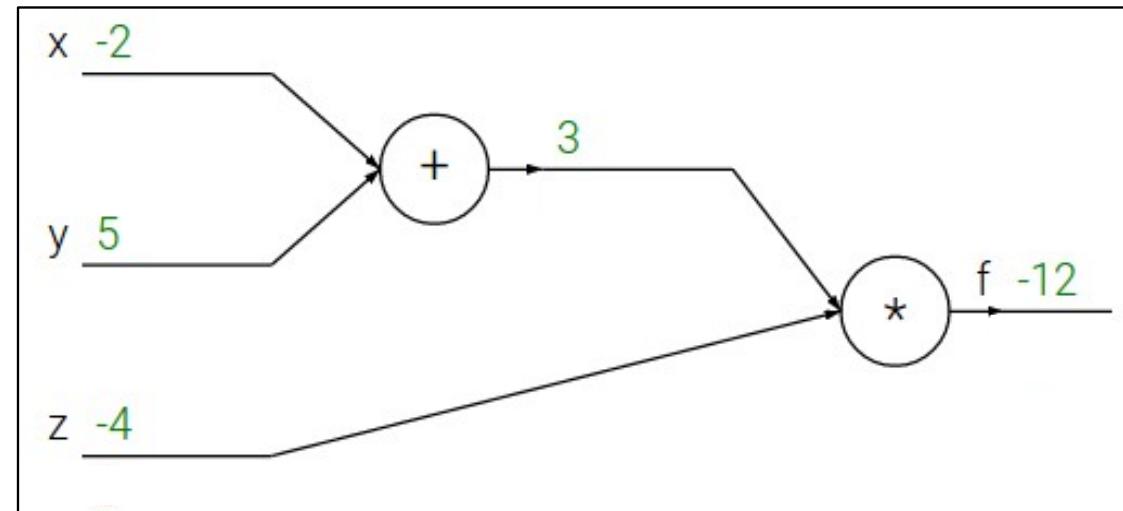


# Neural Turing Machine



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

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

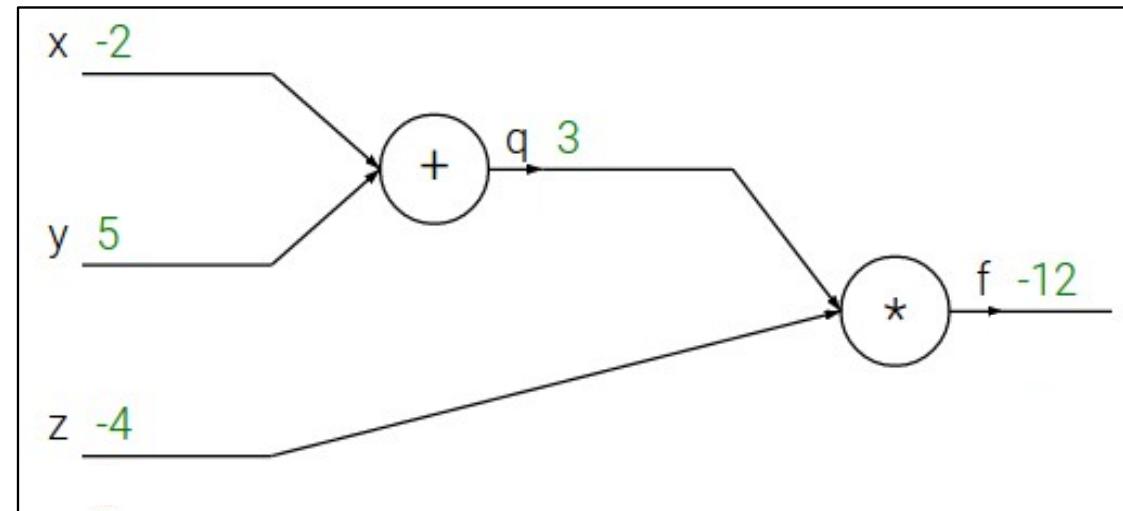


$$f(x, y, z) = (x + y)z$$

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

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

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



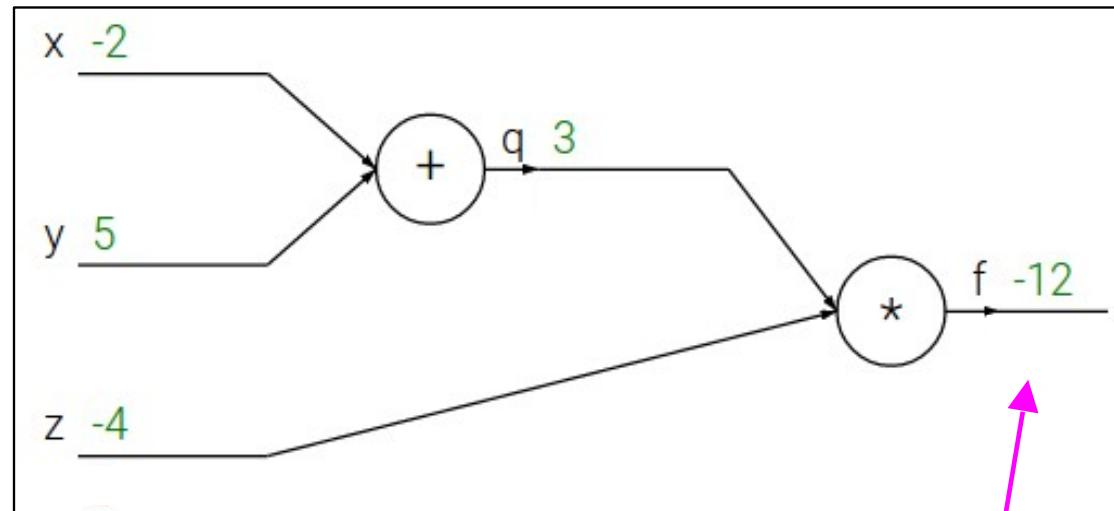
Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

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

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

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

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



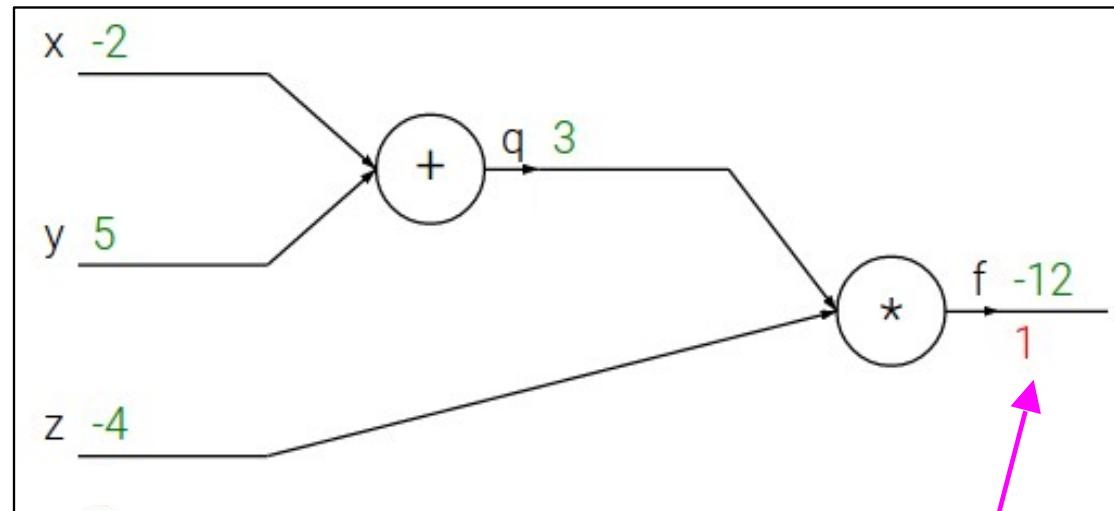
Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

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

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

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

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



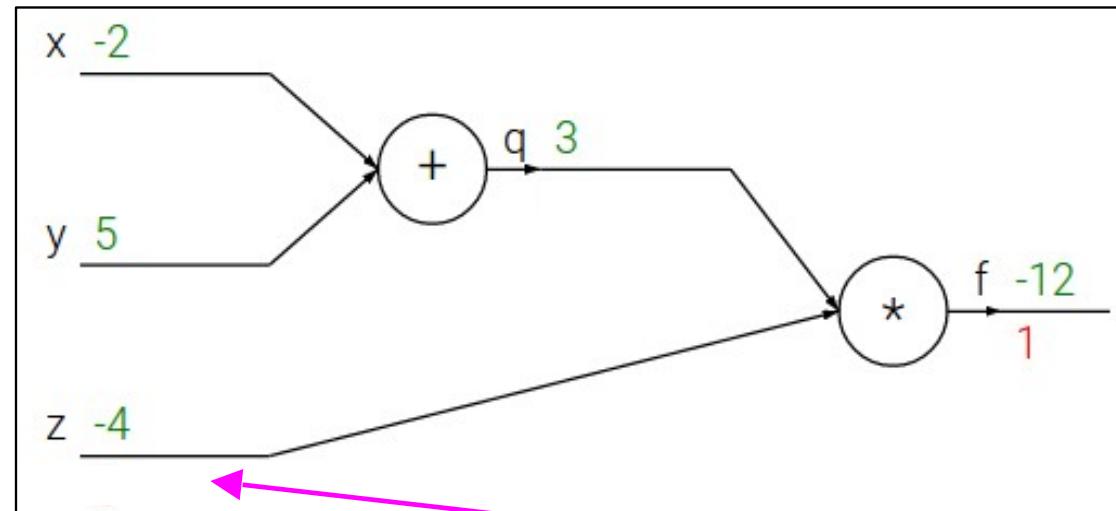
Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

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

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

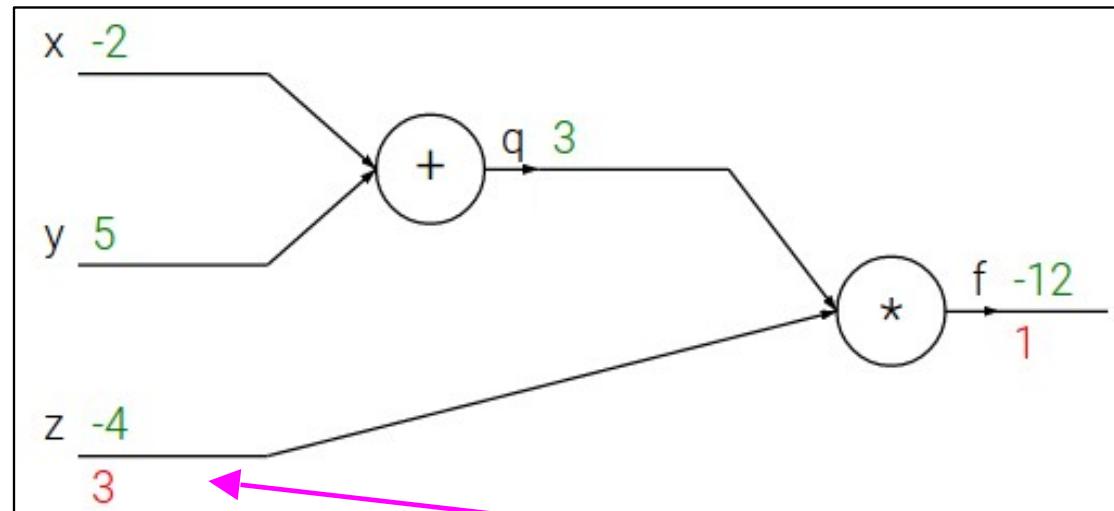
$$\frac{\partial f}{\partial z}$$

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

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

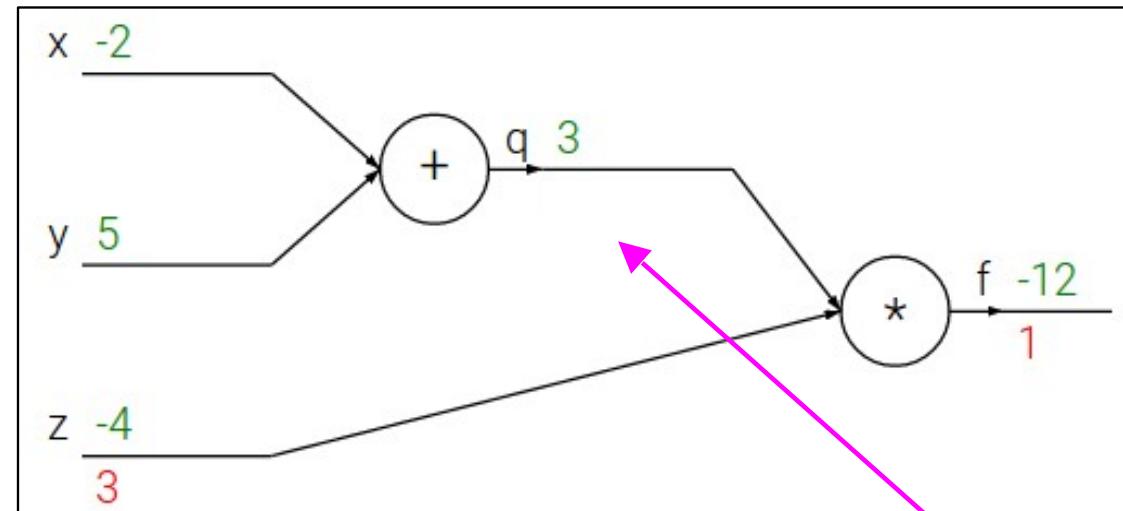
$$\frac{\partial f}{\partial z}$$

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

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

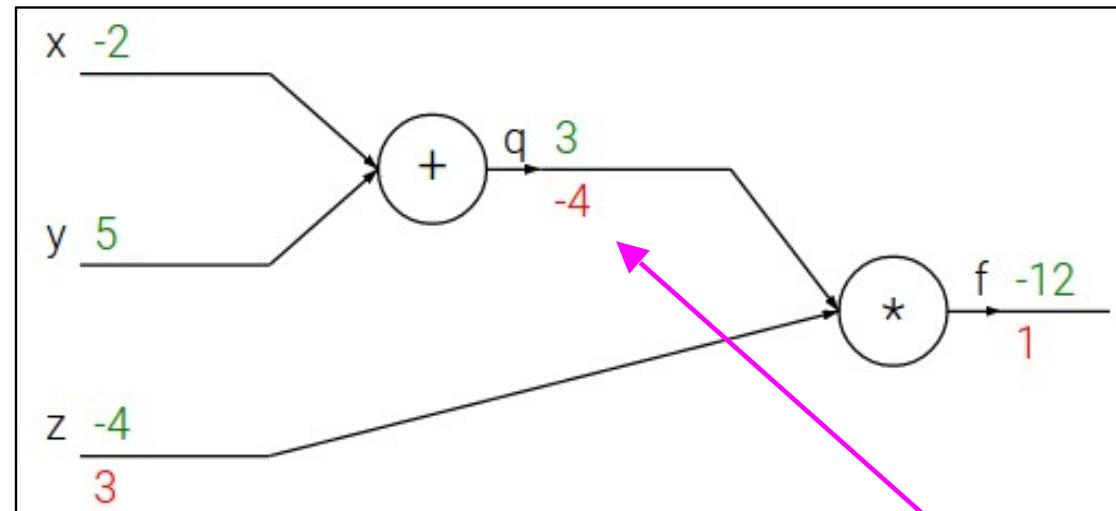
$$\frac{\partial f}{\partial q}$$

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

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

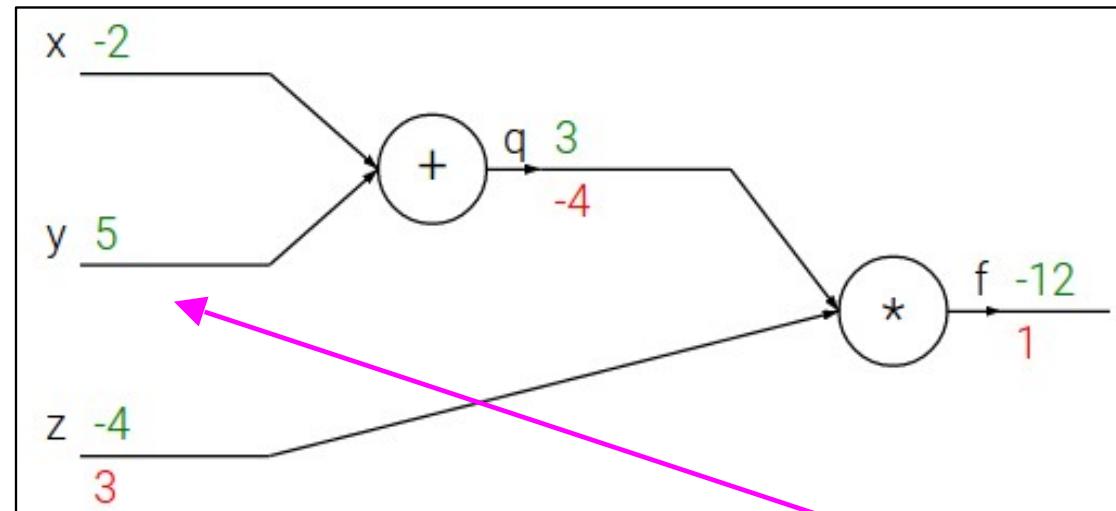
$$\frac{\partial f}{\partial q}$$

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

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

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

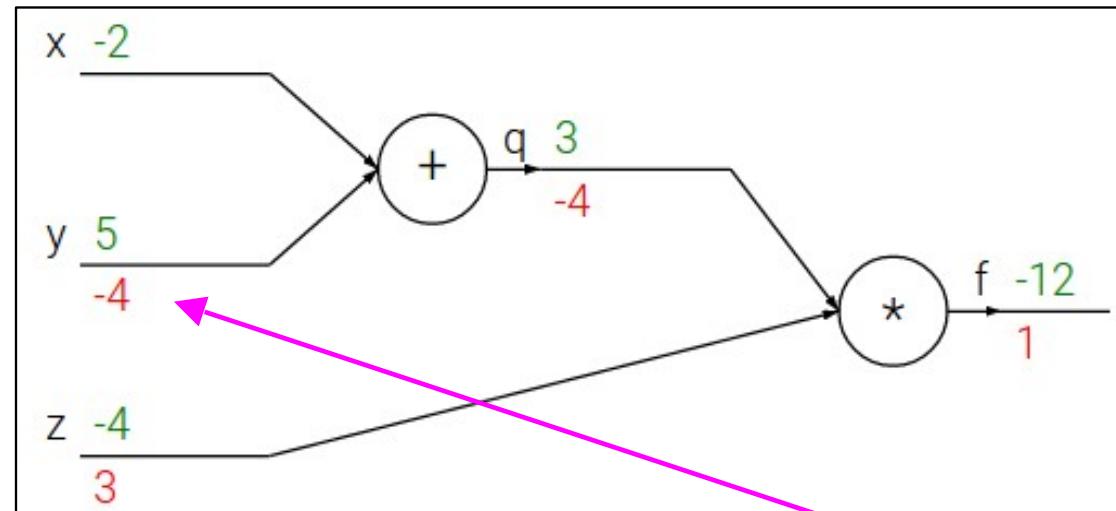
$$f(x, y, z) = (x + y)z$$

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:

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

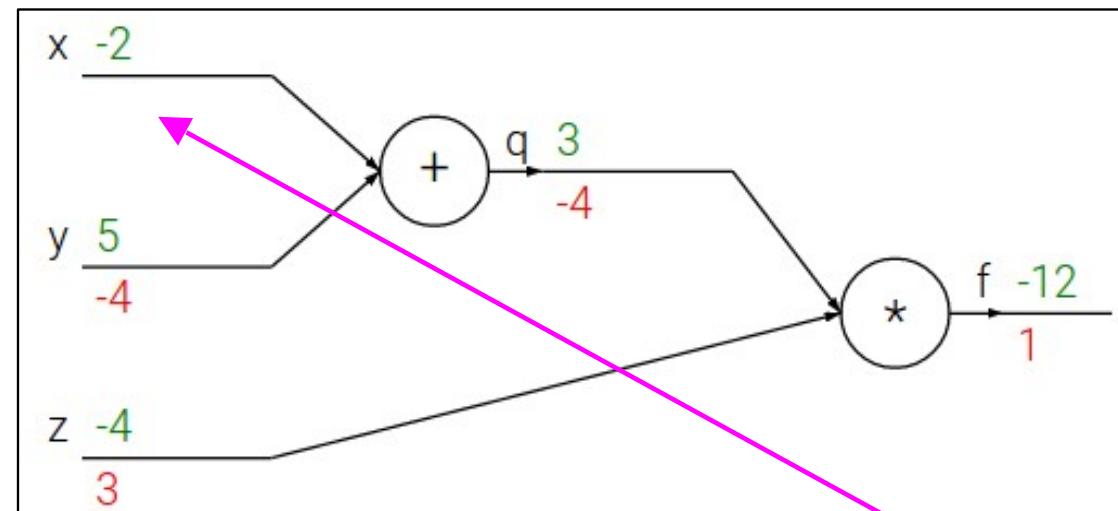
$$\boxed{\frac{\partial f}{\partial y}}$$

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

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

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

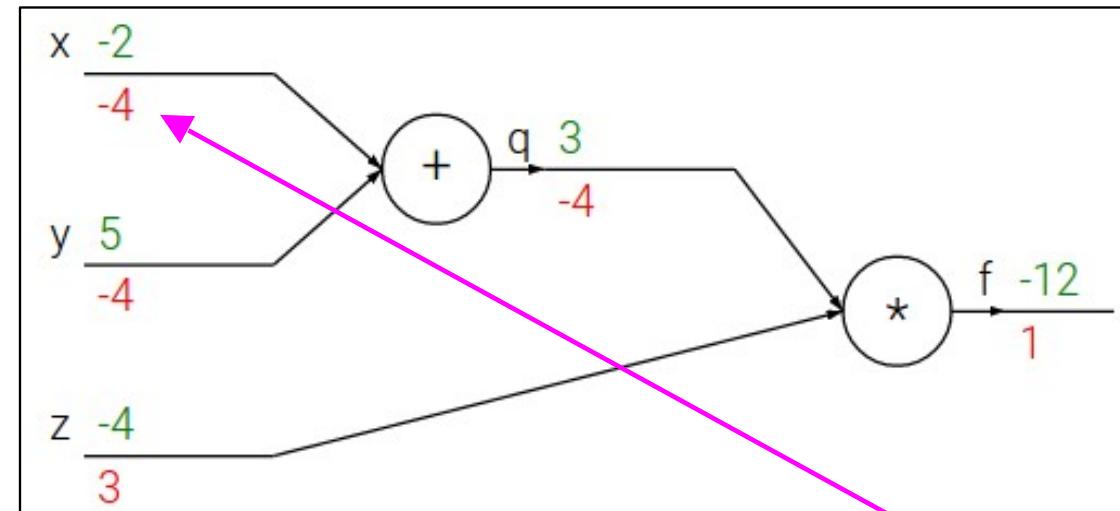
$$f(x, y, z) = (x + y)z$$

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

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

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

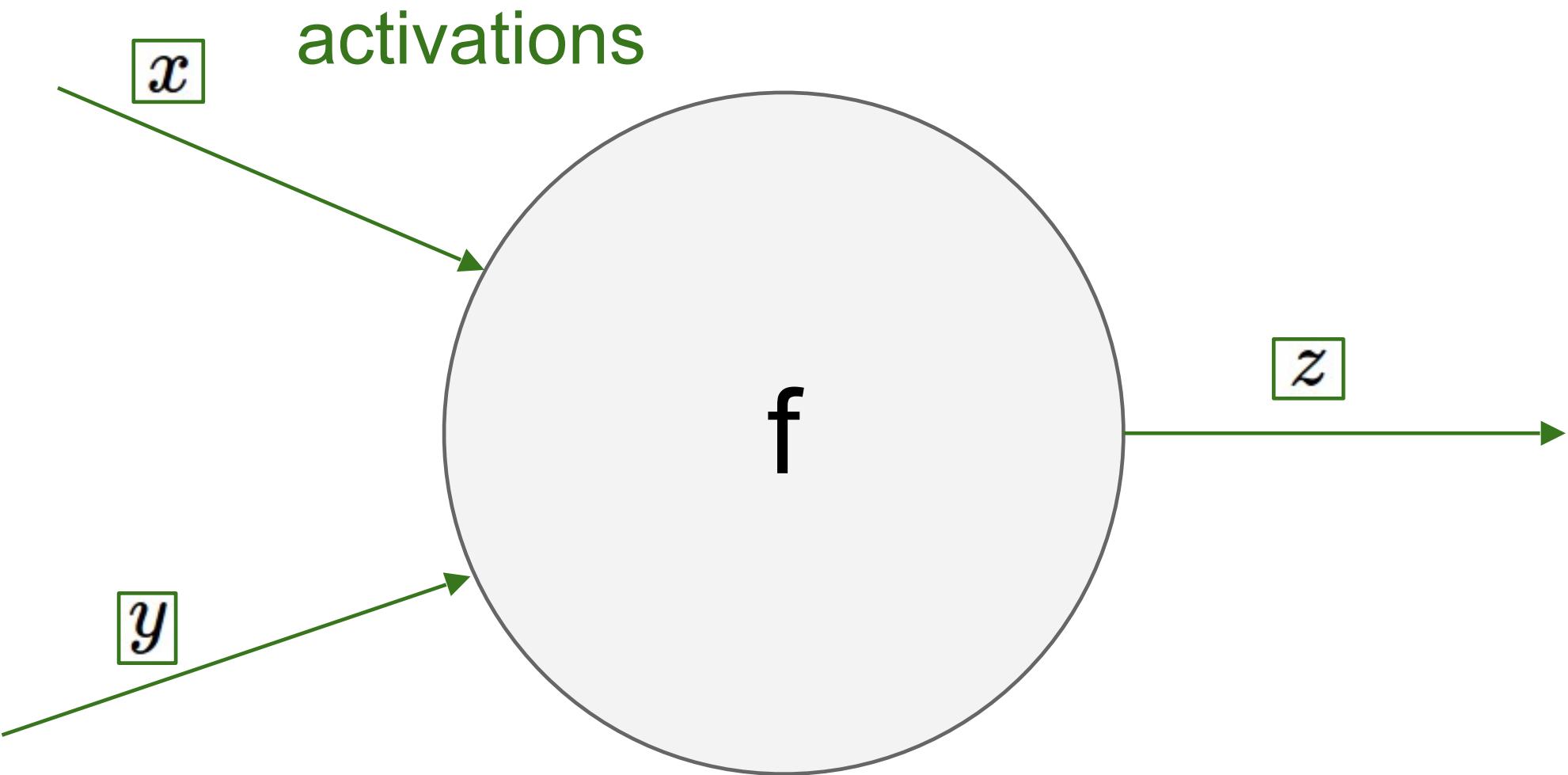
Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

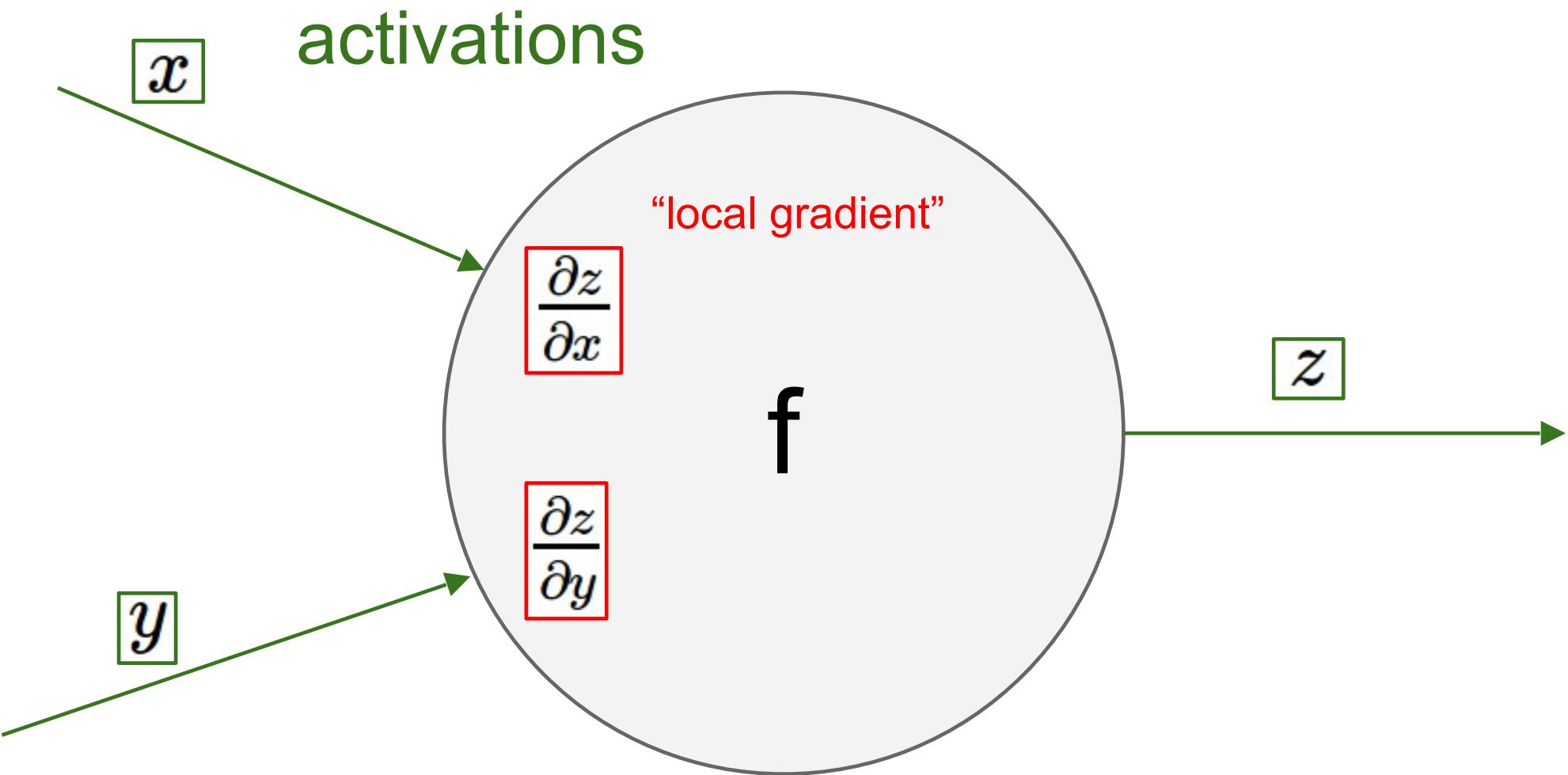


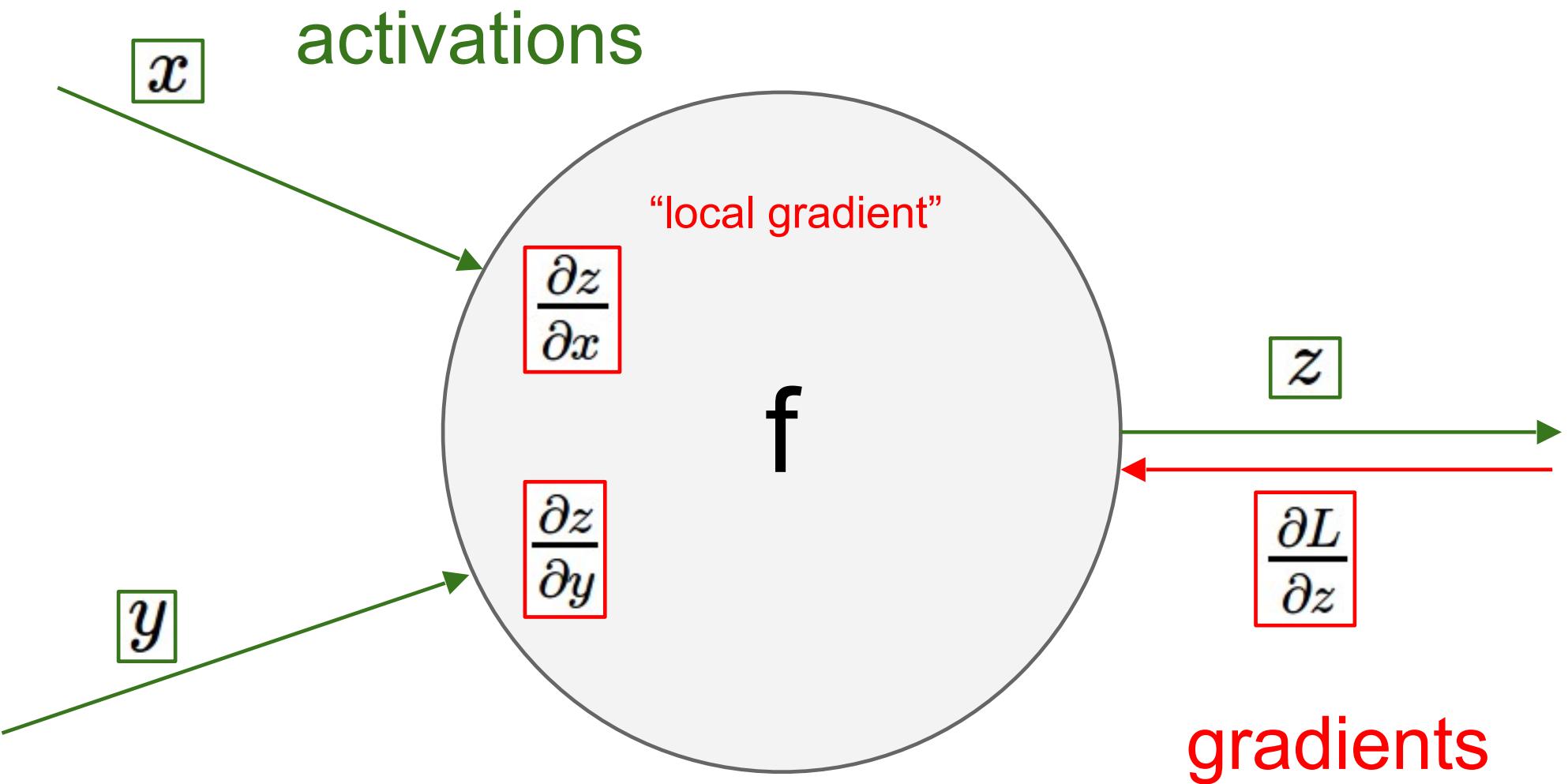
Chain rule:

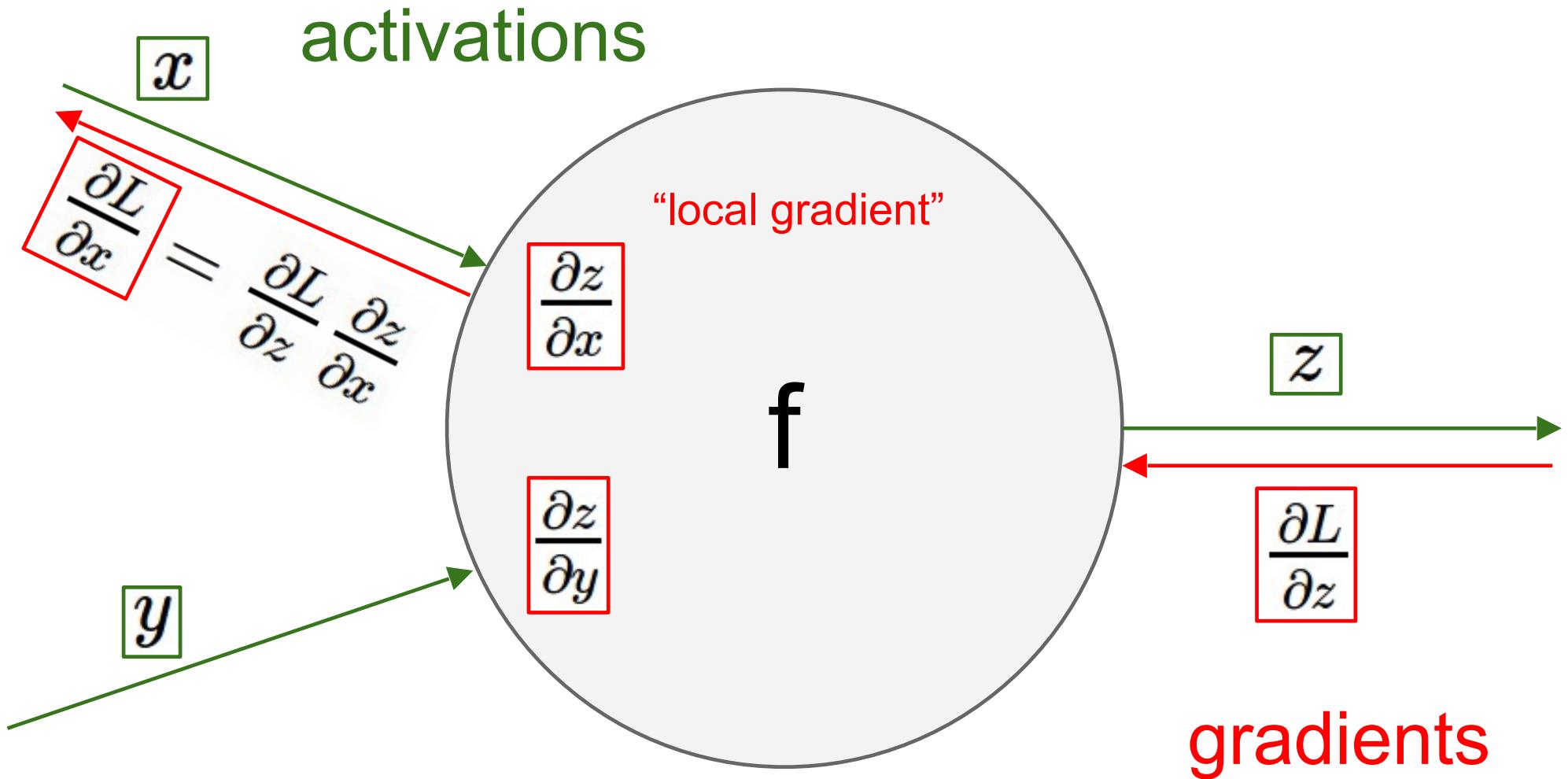
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

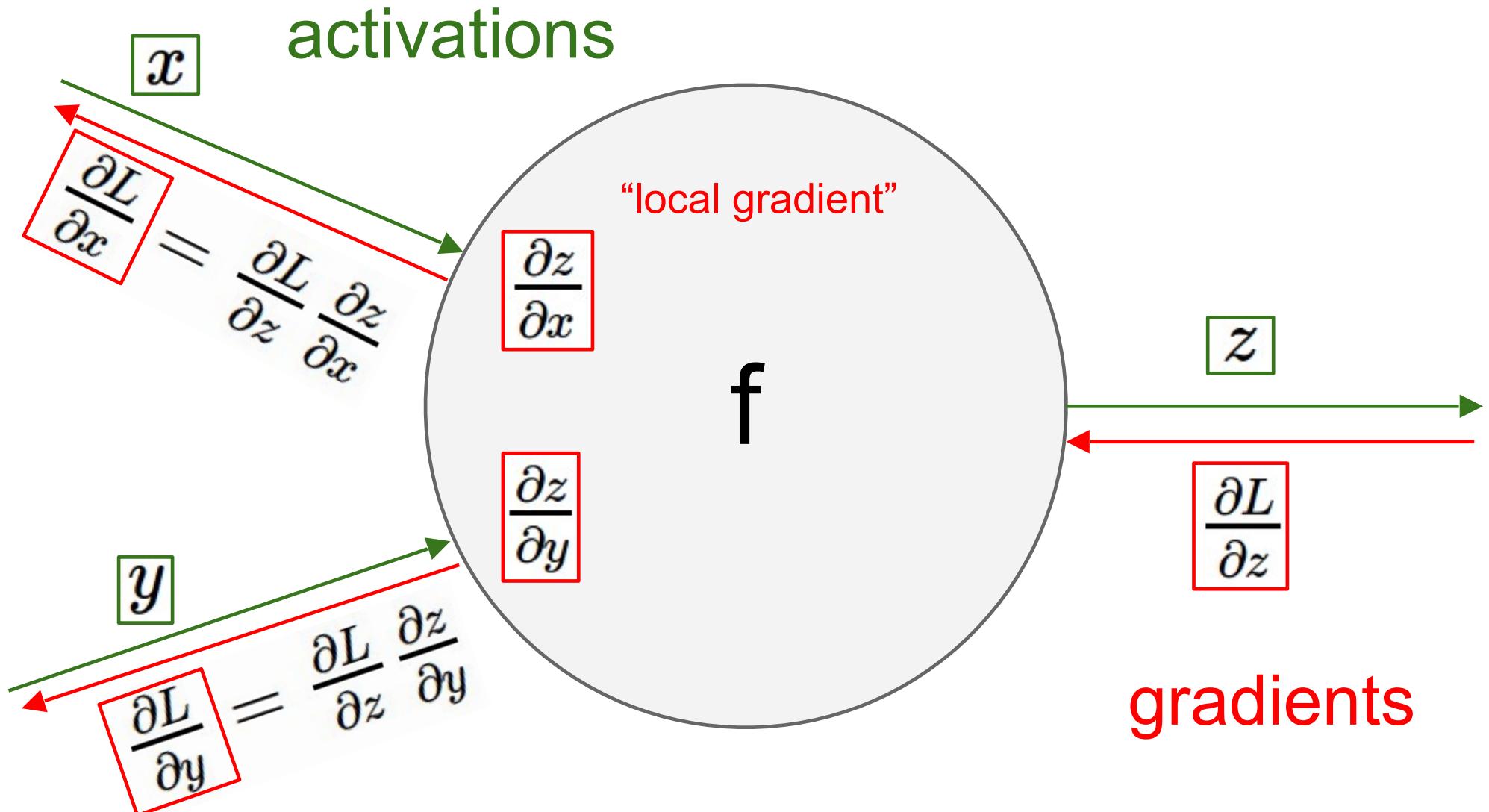
$$\frac{\partial f}{\partial x}$$











activations

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

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

“local gradient”

$f$

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

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

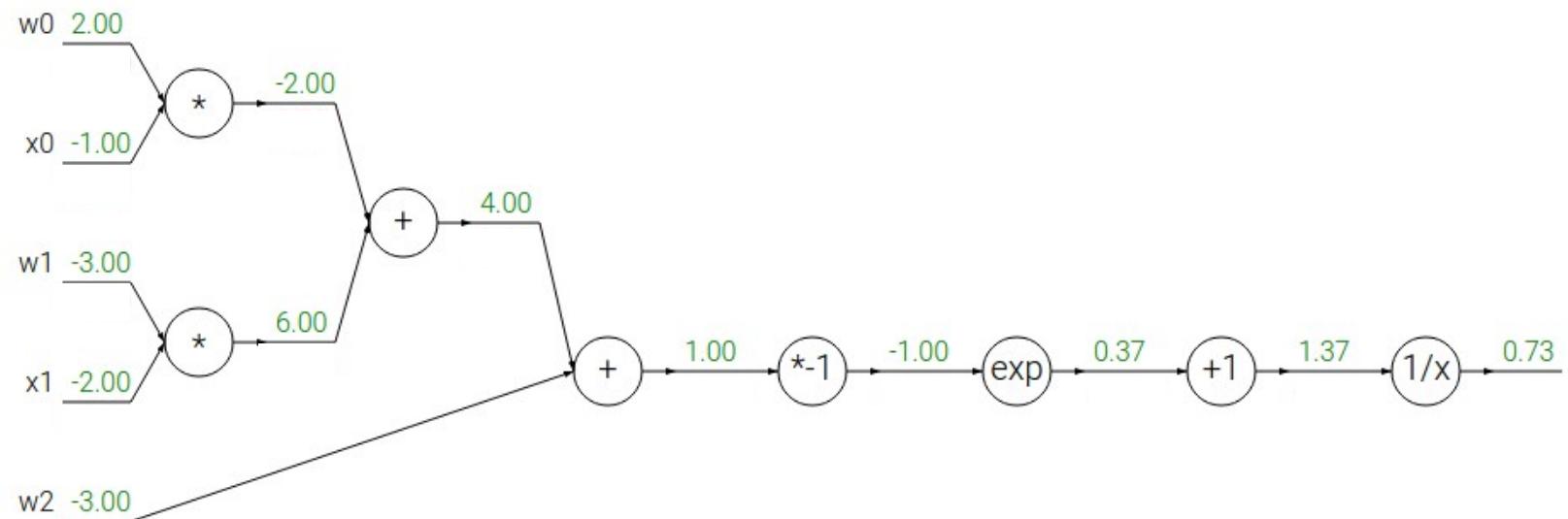
$$z$$

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

gradients

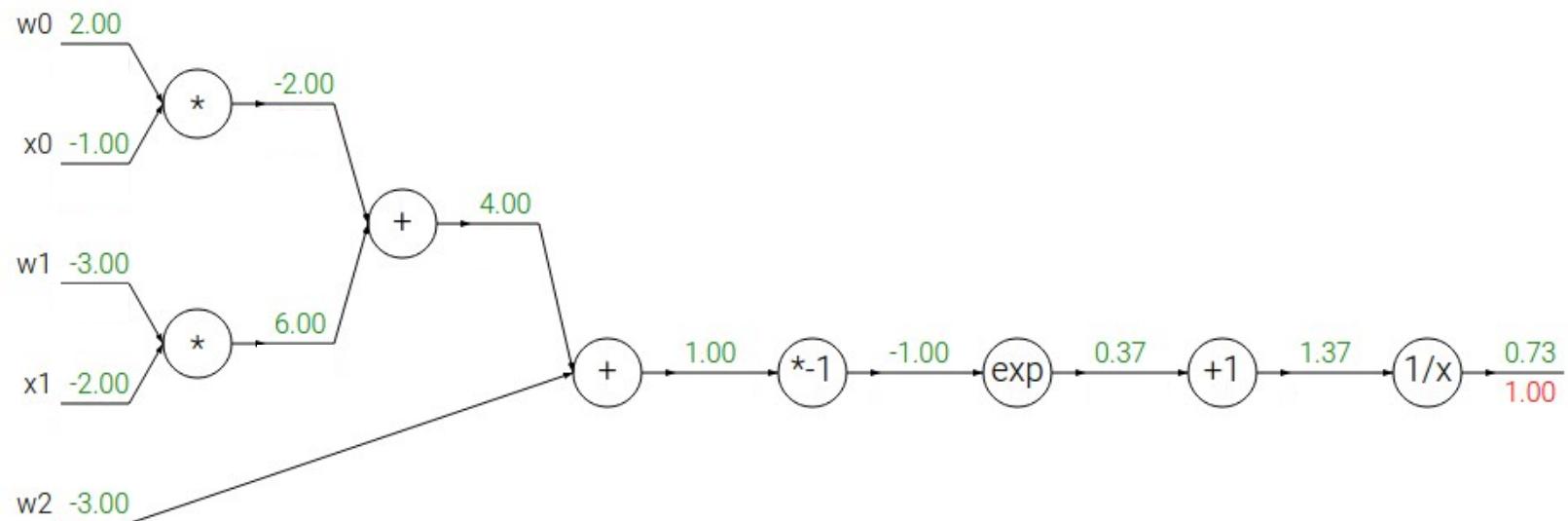
Another example:

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



Another example:

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



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

$\rightarrow$

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

$$f_a(x) = ax$$

$\rightarrow$

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

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

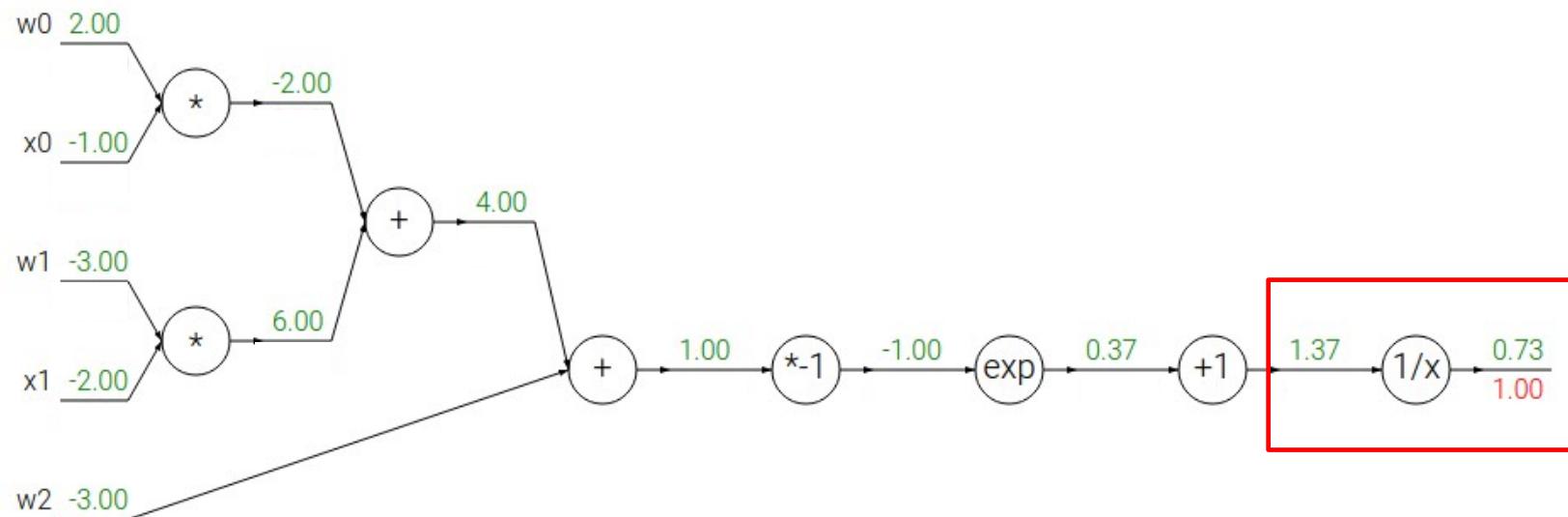
$$f_c(x) = c + x$$

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

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

Another example:

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



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

$\rightarrow$

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

$$f_a(x) = ax$$

$\rightarrow$

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

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

$\rightarrow$

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

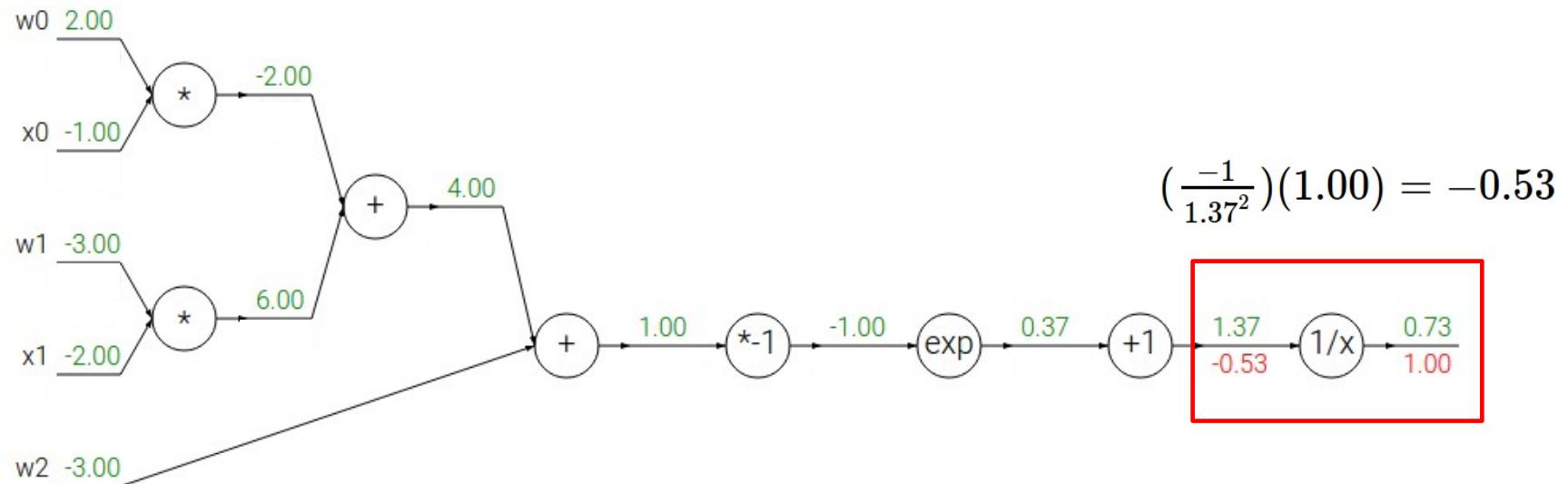
$$f_c(x) = c + x$$

$\rightarrow$

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

Another example:

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



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

$\rightarrow$

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

$$f_a(x) = ax$$

$\rightarrow$

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

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

$\rightarrow$

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

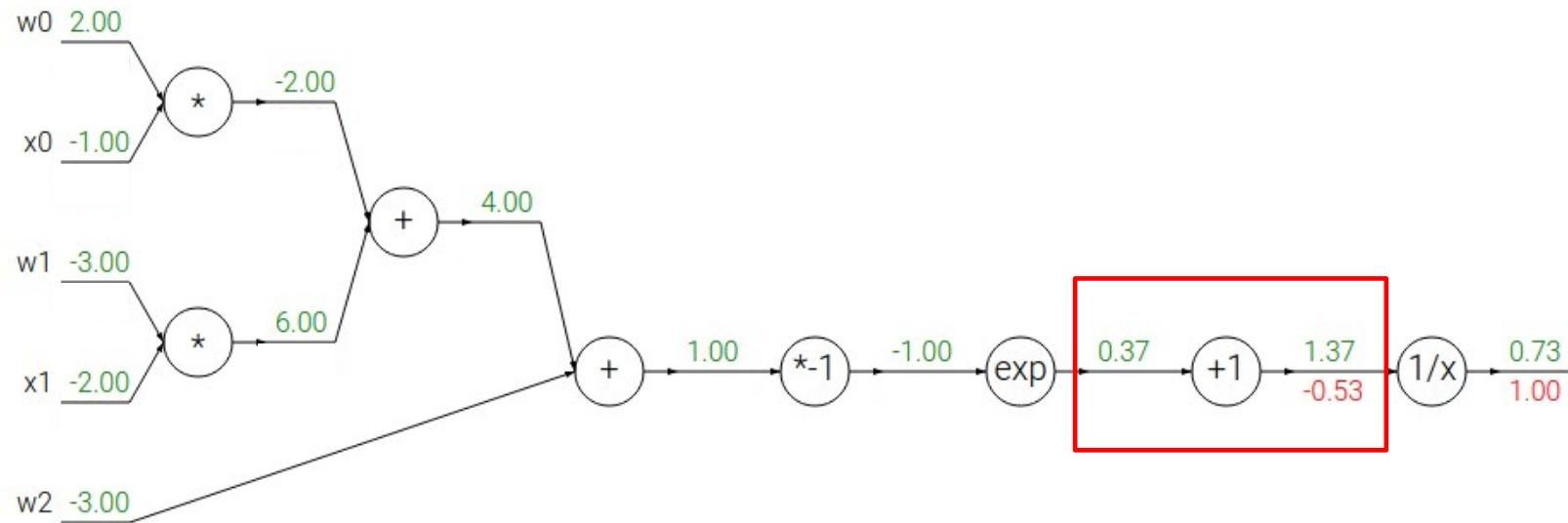
$$f_c(x) = c + x$$

$\rightarrow$

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

Another example:

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



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

$\rightarrow$

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

$$f_a(x) = ax$$

$\rightarrow$

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

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

$\rightarrow$

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

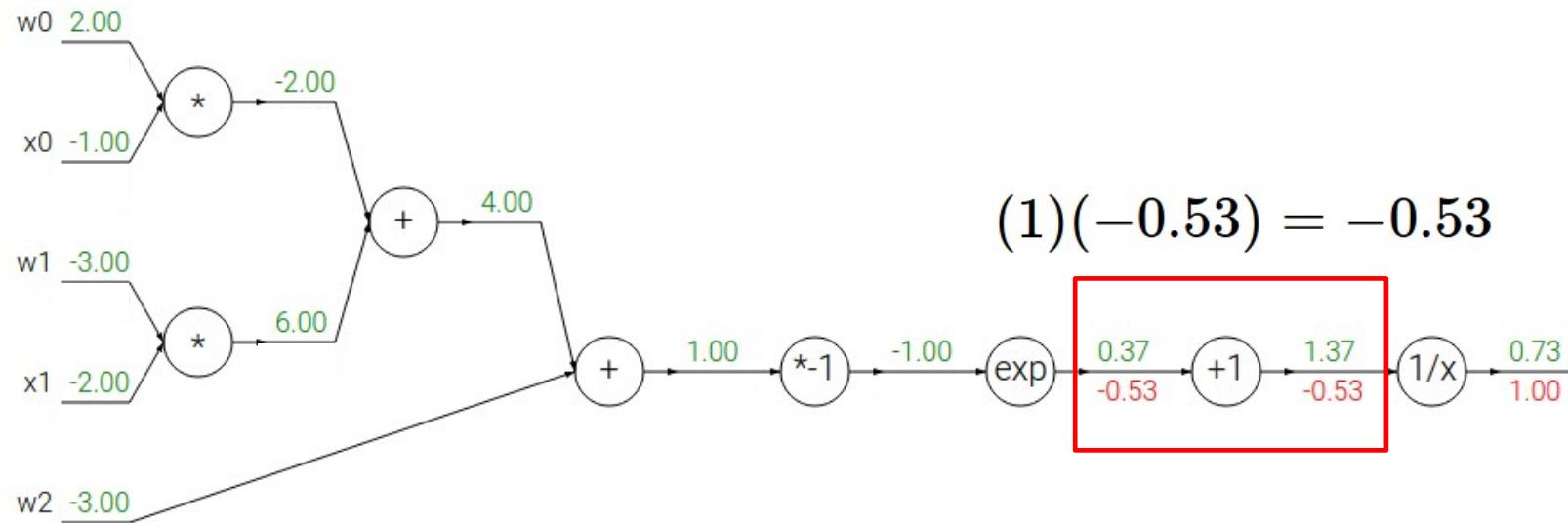
$$f_c(x) = c + x$$

$\rightarrow$

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

Another example:

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



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

→

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

$$f_a(x) = ax$$

→

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

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

→

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

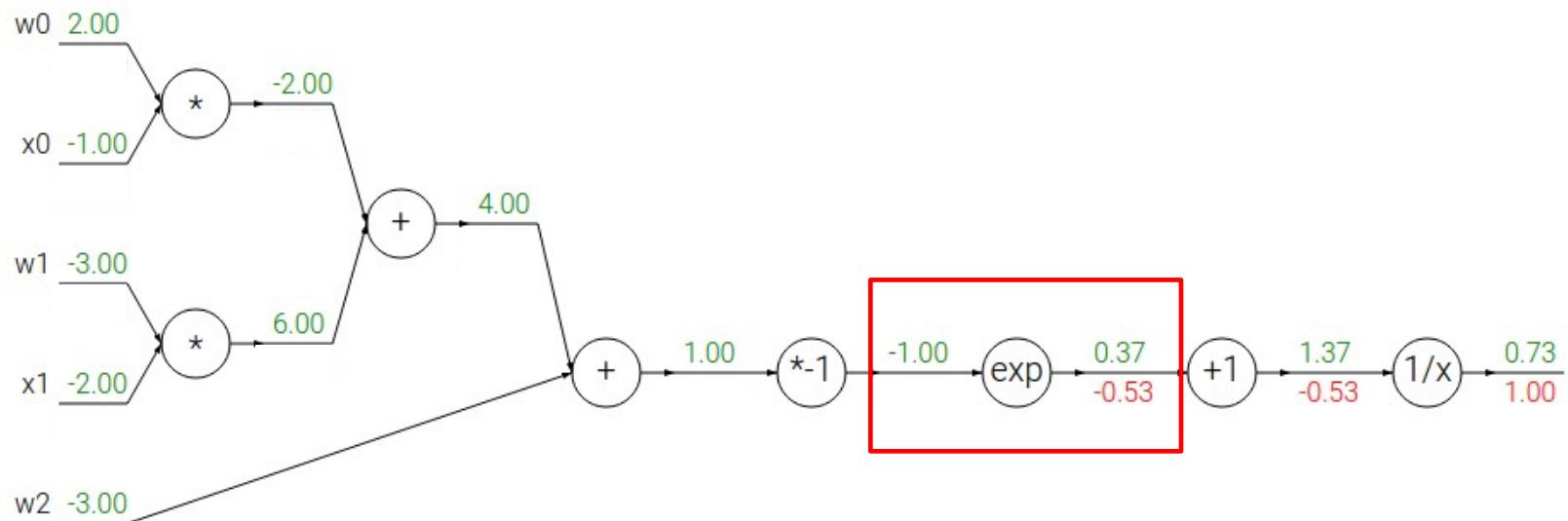
$$f_c(x) = c + x$$

→

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

Another example:

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



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

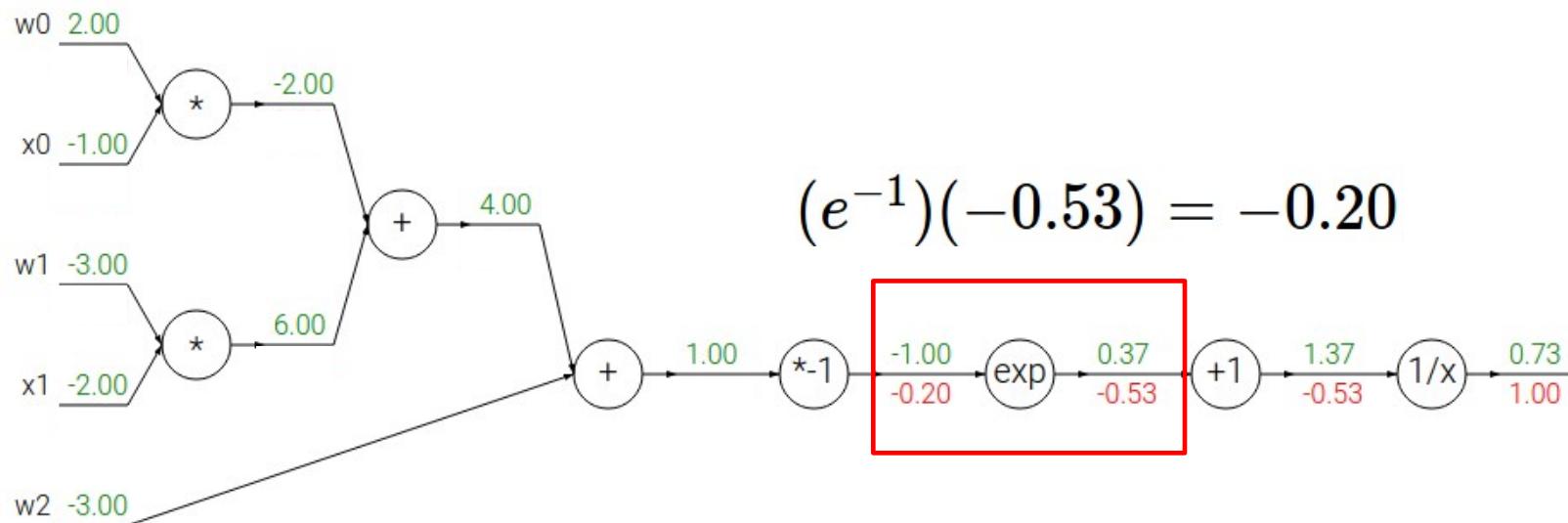
$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

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

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

Another example:

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



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

$\rightarrow$

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

$$f_a(x) = ax$$

$\rightarrow$

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

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

$\rightarrow$

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

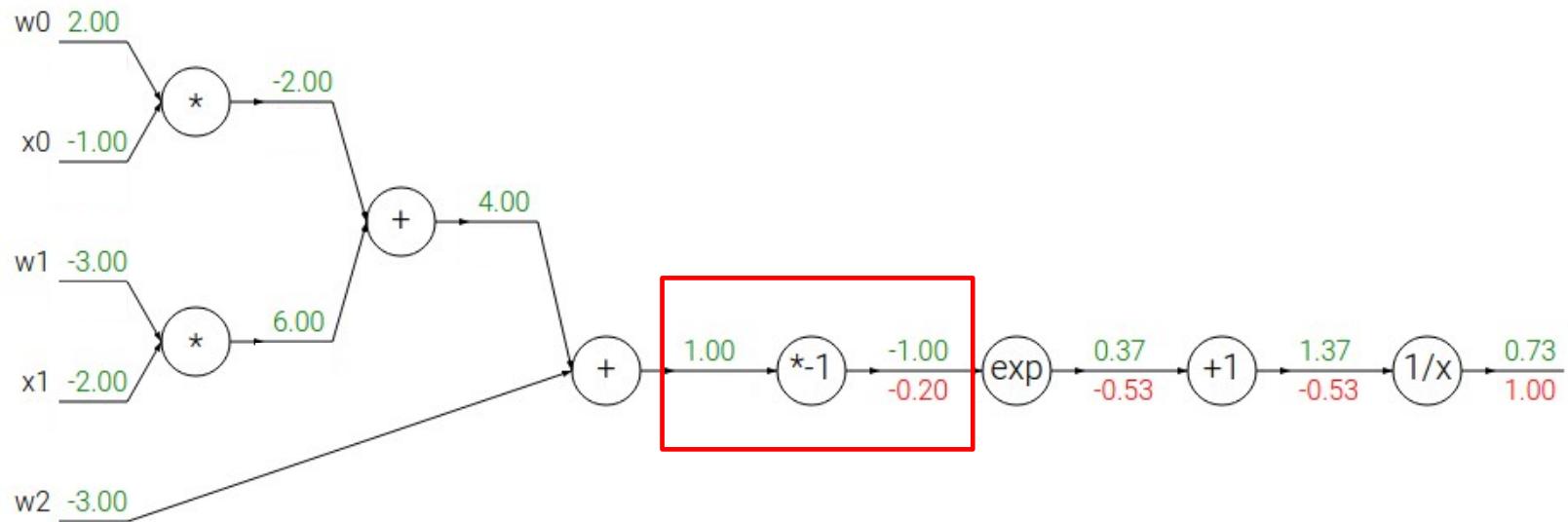
$$f_c(x) = c + x$$

$\rightarrow$

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

Another example:

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



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

$\rightarrow$

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

$$f_a(x) = ax$$

$\rightarrow$

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

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

$\rightarrow$

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

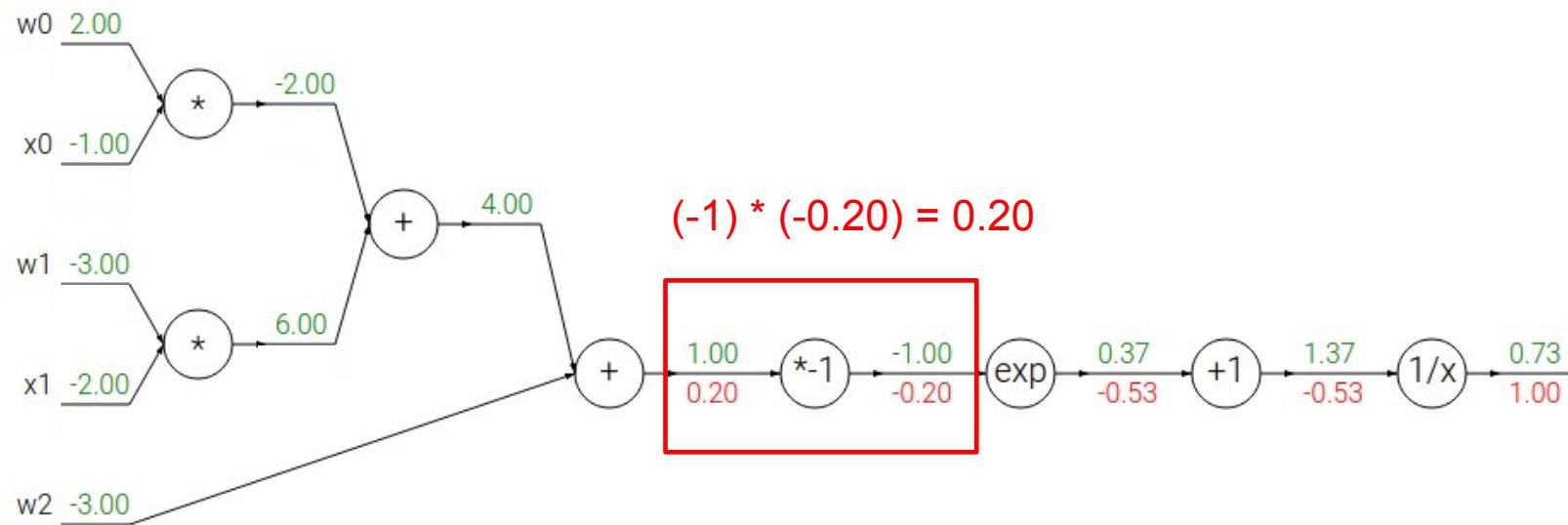
$$f_c(x) = c + x$$

$\rightarrow$

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

Another example:

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



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

$\rightarrow$

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

$$f_a(x) = ax$$

$\rightarrow$

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

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

$\rightarrow$

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

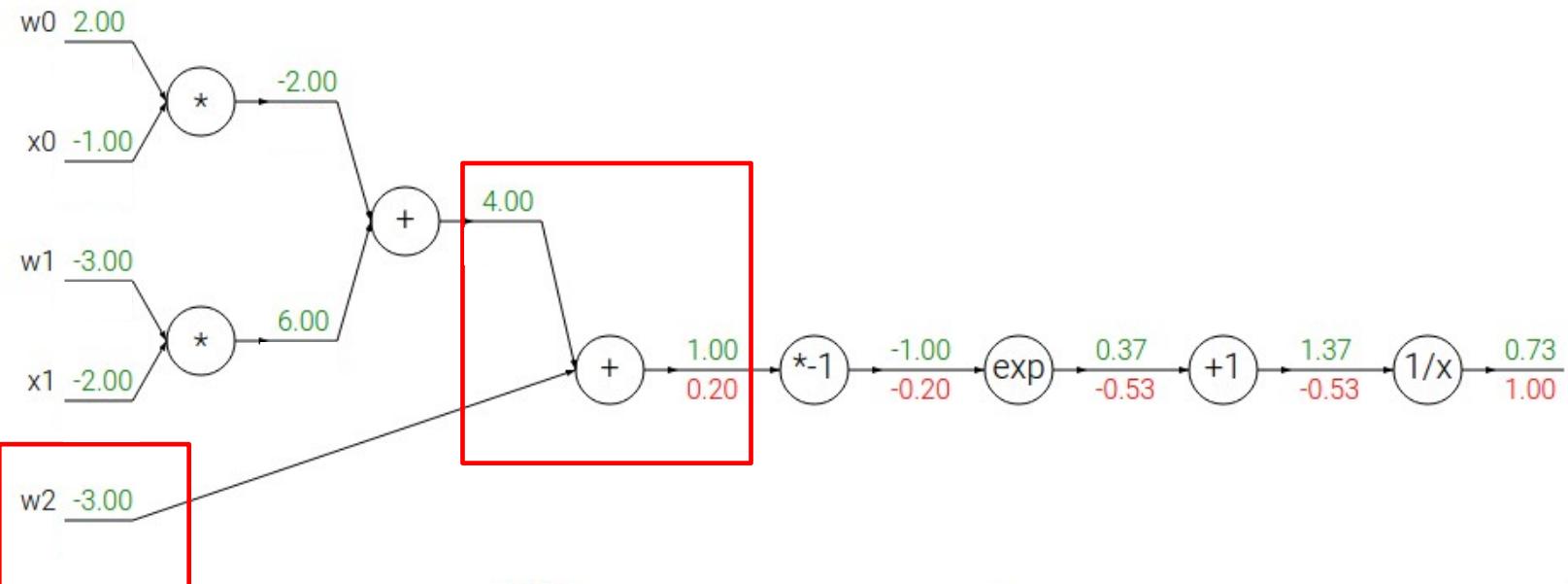
$$f_c(x) = c + x$$

$\rightarrow$

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

Another example:

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



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

$\rightarrow$

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

$$f_a(x) = ax$$

$\rightarrow$

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

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

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

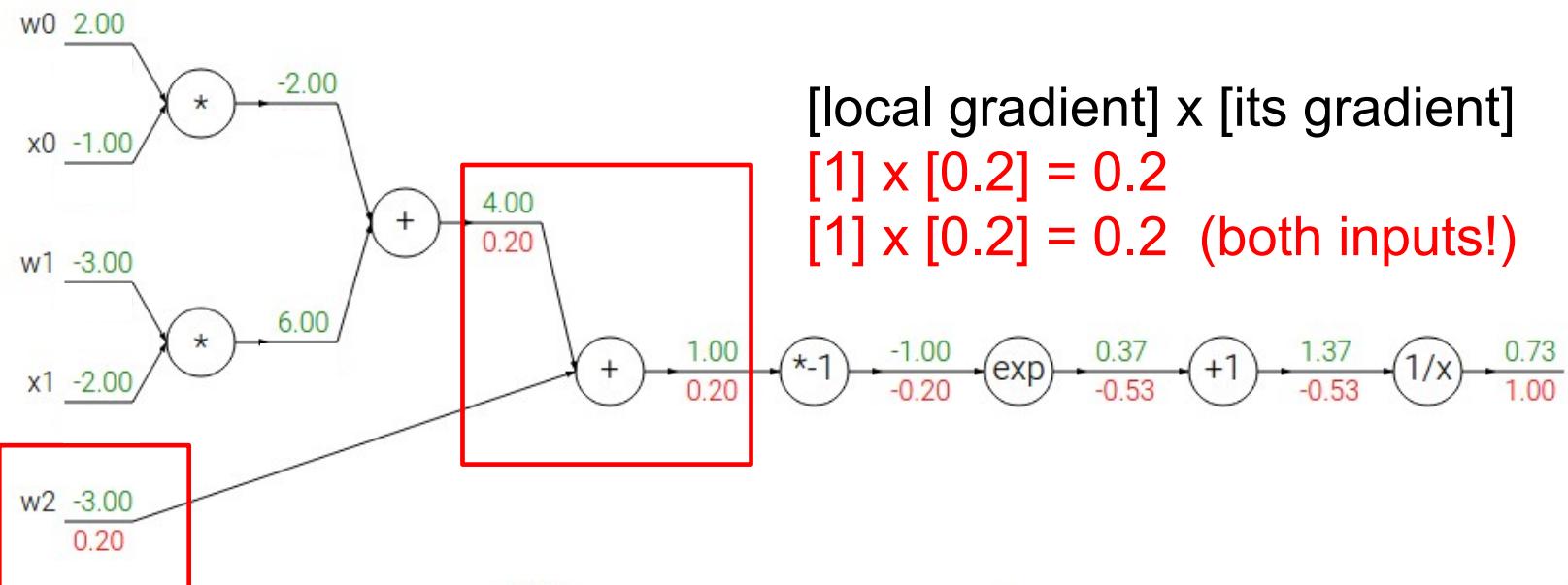
$\rightarrow$

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

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

Another example:

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



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

$\rightarrow$

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

$$f_a(x) = ax$$

$\rightarrow$

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

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

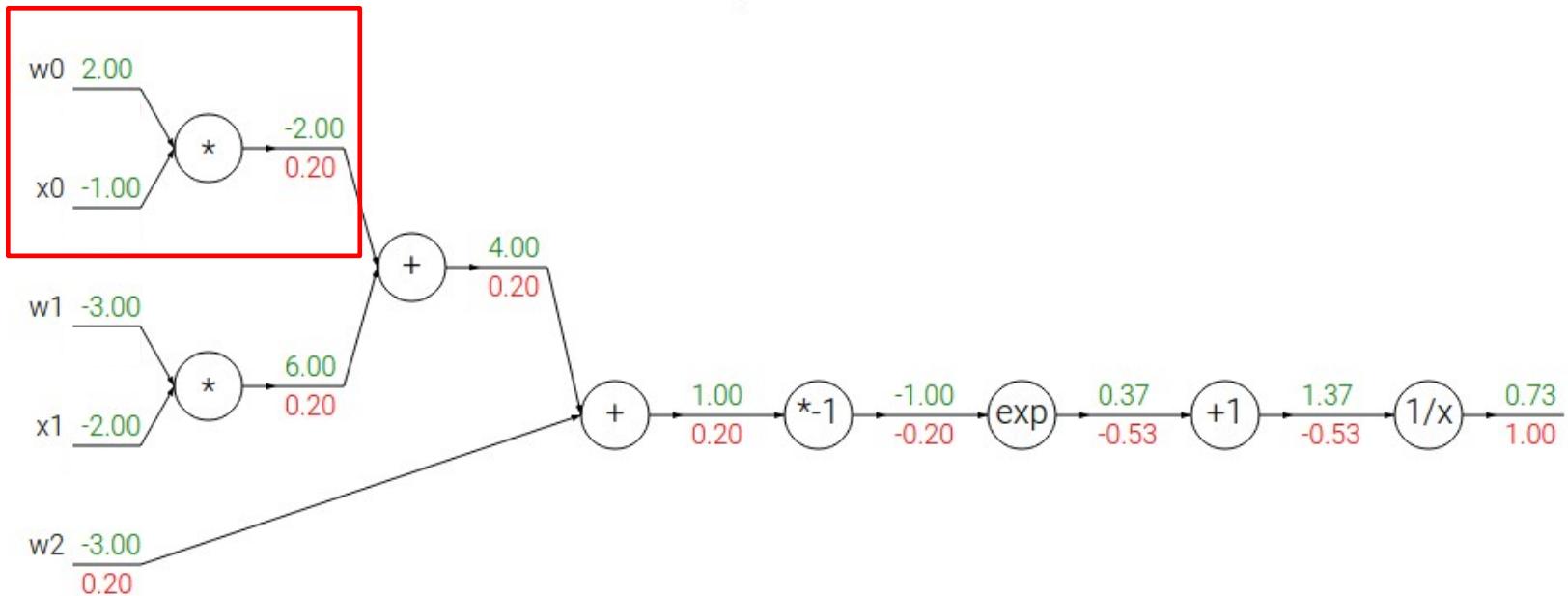
$$f_c(x) = c + x$$

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

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

Another example:

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



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

$\rightarrow$

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

$$f_a(x) = ax$$

$\rightarrow$

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

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

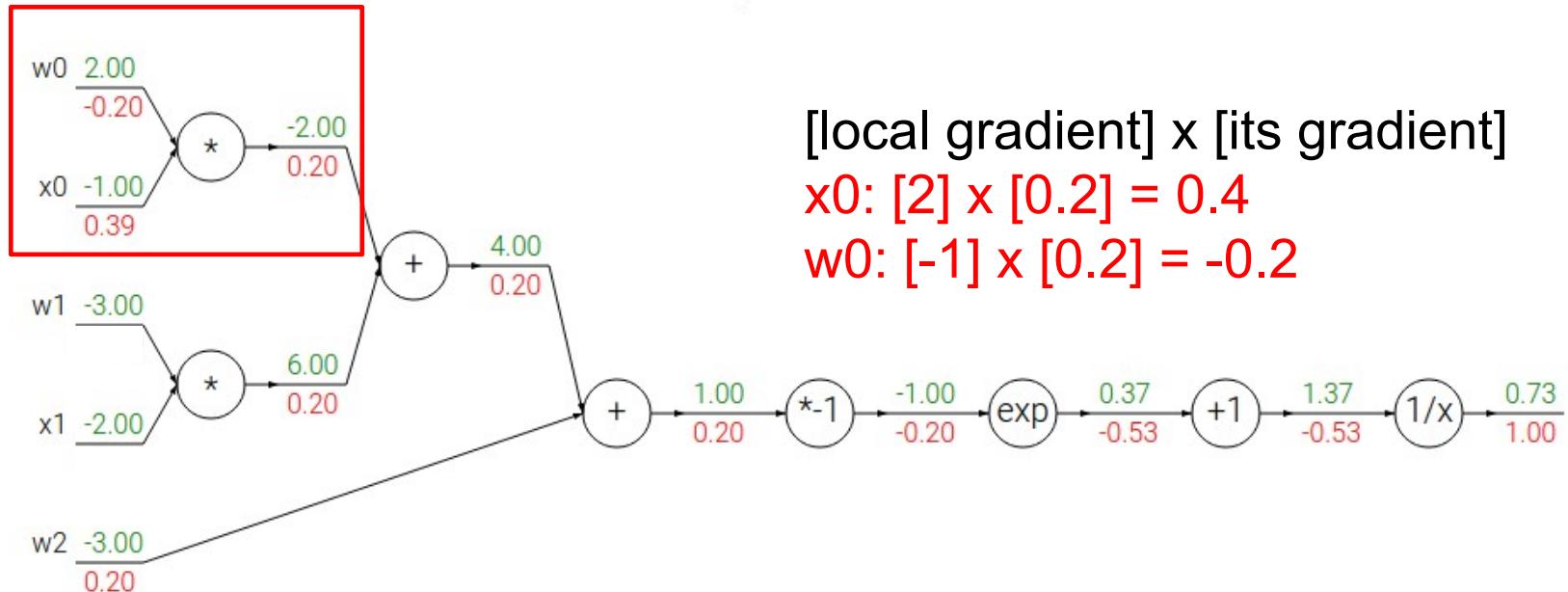
$$f_c(x) = c + x$$

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

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

Another example:

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



[local gradient] x [its gradient]

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

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

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

→

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

$$f_a(x) = ax$$

→

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

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

→

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

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

→

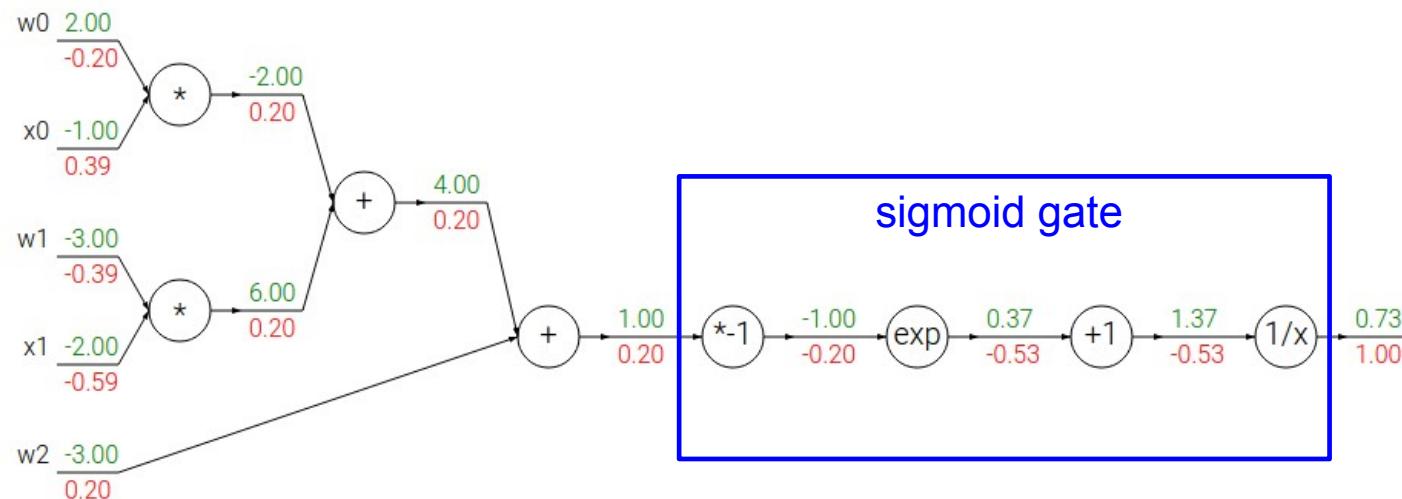
$$\frac{df}{dx} = 1$$

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

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

sigmoid function

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

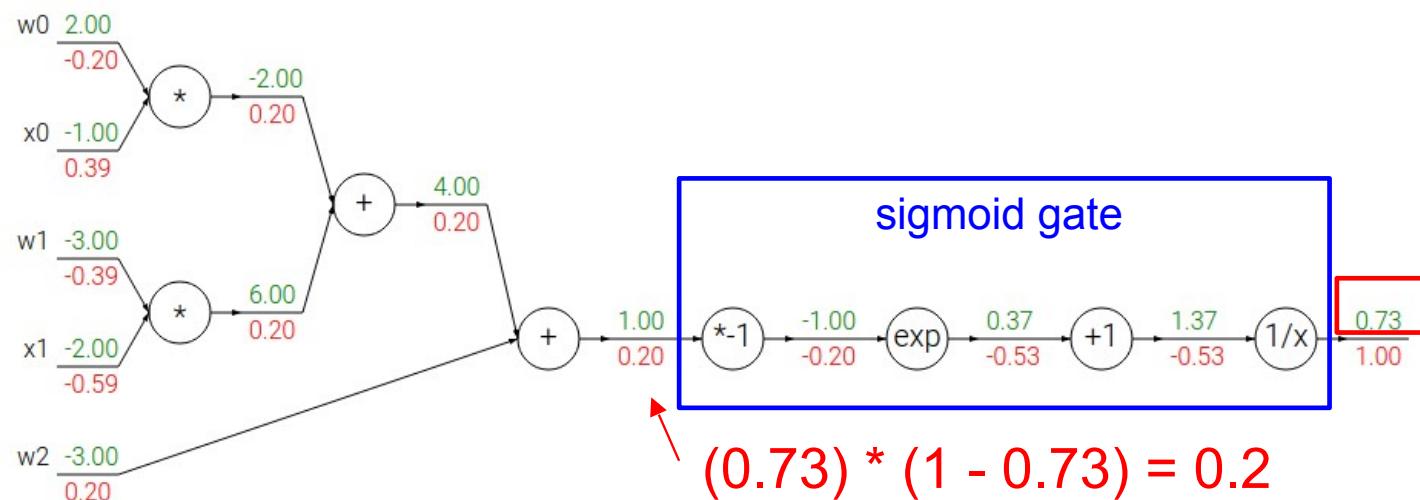


$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

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

sigmoid function

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

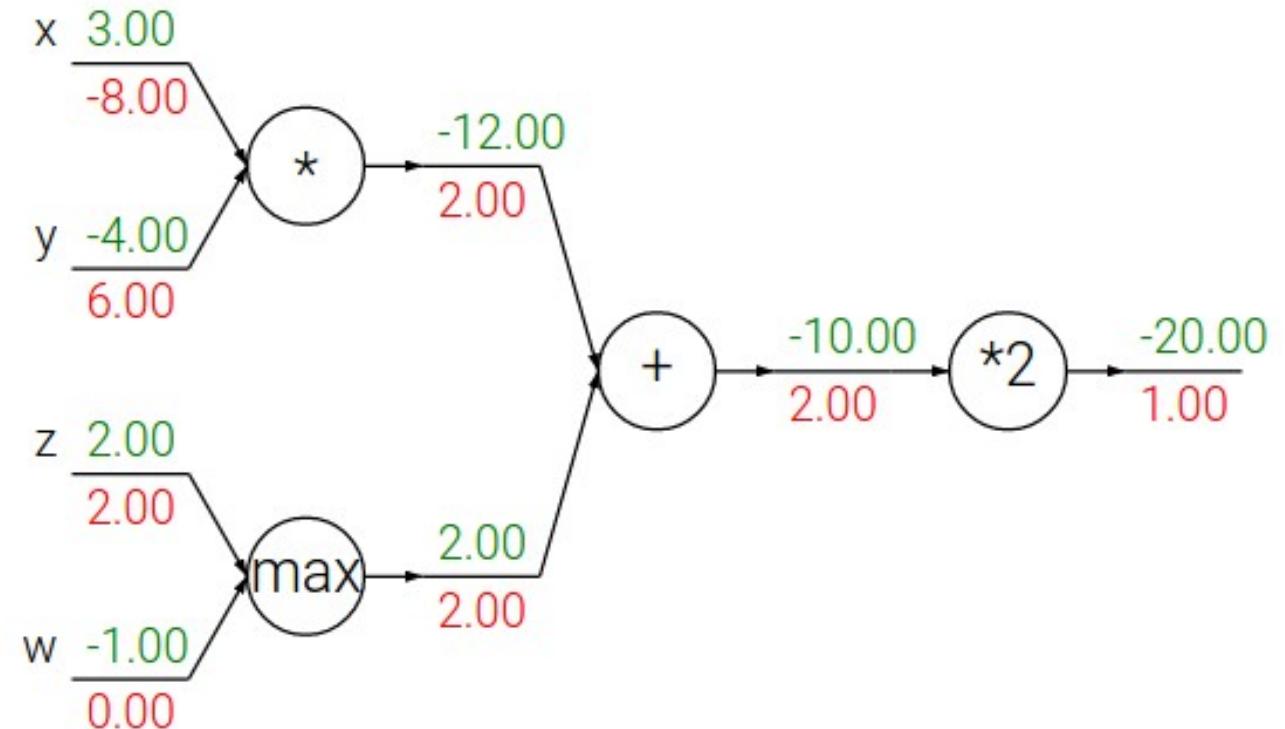


## Patterns in backward flow

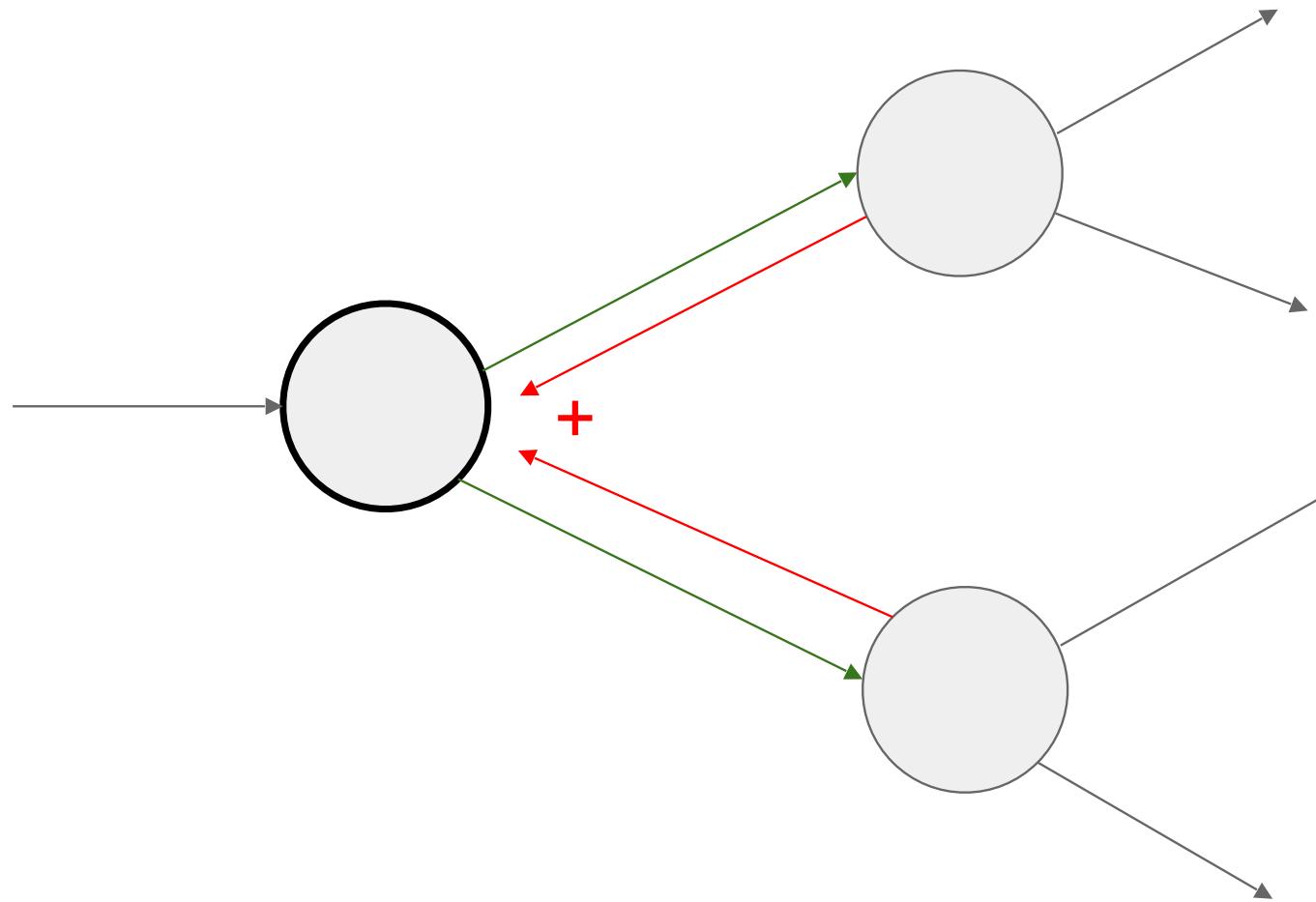
**add** gate: gradient distributor

**max** gate: gradient router

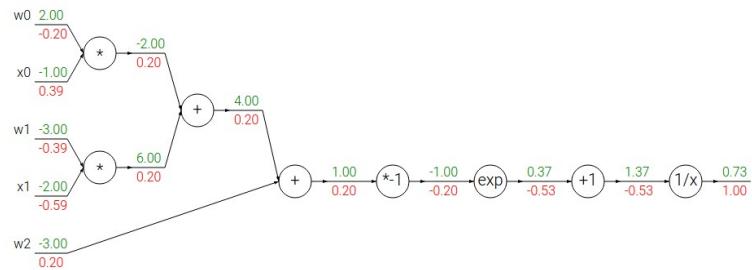
**mul** gate: gradient... “switcher”?



# Gradients add at branches



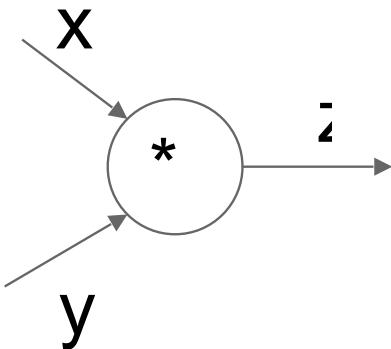
# Implementation: forward/backward API



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



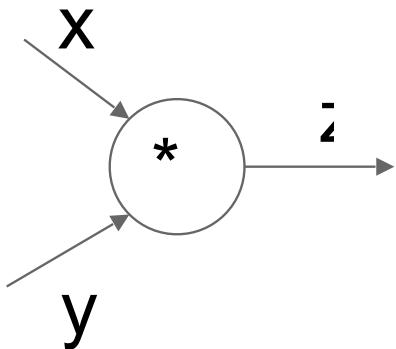
(x,y,z are scalars)

```
class MultiplyGate(object):  
    def forward(x,y):  
        z = x*y  
        return z  
    def backward(dz):  
        # dx = ... #todo  
        # dy = ... #todo  
        return [dx, dy]
```

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

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

# 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)

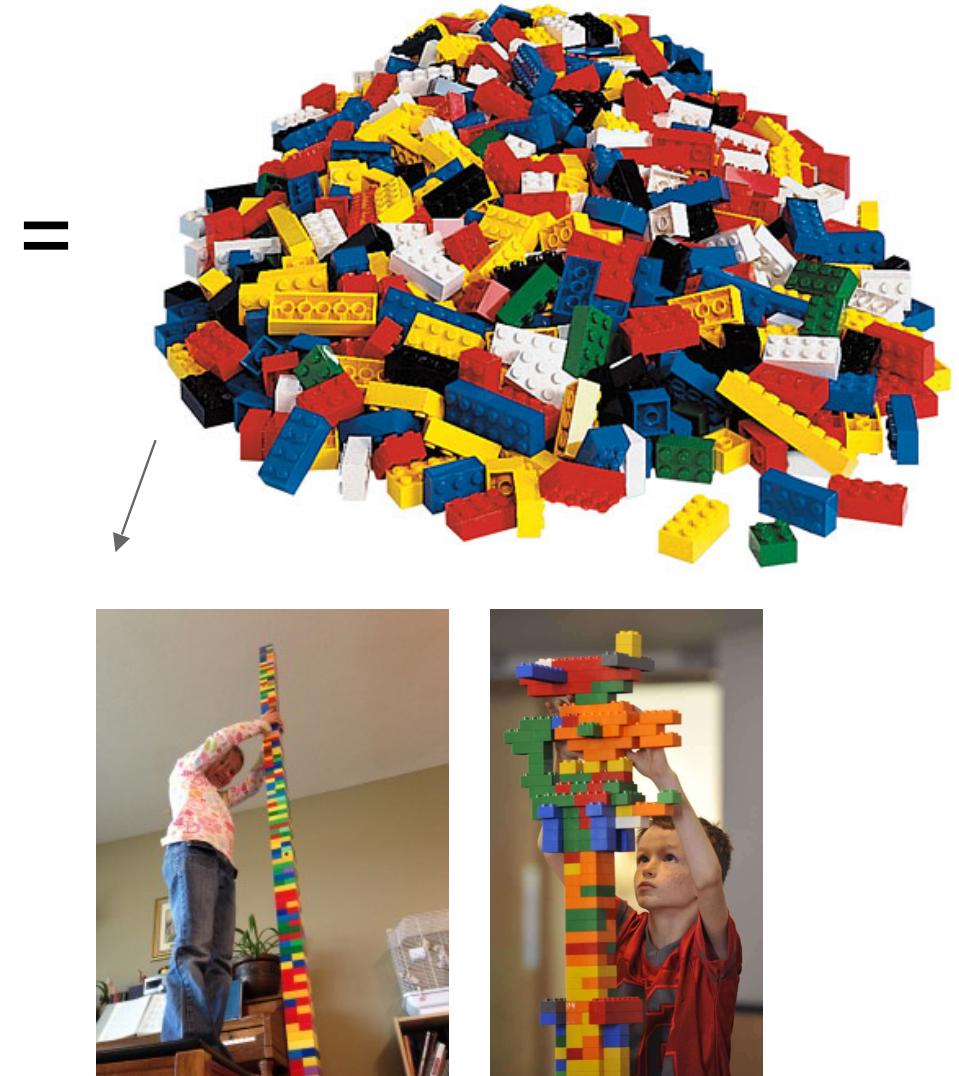


# Example: Torch Layers

Torch / nn			
Code		Issues 27	Star 278
Branch: master		New pull request	New file Find file HTTPS https://github.com/search Download ZIP
No description or website provided.			
LogSigmoid.lua	Add THNN conversion of [ELU, LeakyReLU, LogSigmoid, LogSoftMax, Look... 7 days ago		
LogSoftMax.lua	Add THNN conversion of [ELU, LeakyReLU, LogSigmoid, LogSoftMax, Look... 7 days ago		
LookupTable.lua	Harmonize LookupTable signature with cunn impl 5 days ago		
MM.lua	Rename unpack to table.unpack for Lua 5.2 8 months ago		
MSECriterion.lua	Add ScaLverage to criterion in the constructor 2 months ago		
MarginCriterion.lua	modernized MarginCriterion a year ago		
MarginRankingCriterion.lua	Fix batch mode in MarginRankingCriterion 4 days ago		
Max.lua	Merge pull request #444 from gigahester 2 months ago		
Mean.lua	Add support for negative dimension and both batch and non batch input... 2 months ago		
Min.lua	Merge pull request #444 from gigahester 2 months ago		
MinTable.lua	cancel unused table and useless expression 29 days ago		
Module.lua	Revert "Don't re-factor parameters if they are already flattened" 15 hours ago		
Mul.lua	removing the requirement for providing size in nn.Mul a year ago		
MultiConstant.lua	Ignore updateGradInput if self.gradInput is nil 3 months ago		
MultiCriterion.lua	assets in MultiCriterion and ParallelCriterion add 2 months ago		
MultiLabelMarginCriterion.lua	initial levancy of torch tree 4 years ago		
MultiMarginCriterion.lua	multimargin supports pr2 11 months ago		
Narrow.lua	type# in Narrow not done in place. 6 months ago		
NarrowTable.lua	NarrowTable 6 months ago		
Normalise.lua	Remove bren and badbren from Normalise, because they allocate memory.... 20 days ago		
PRelu.lua	Buffers for PRelu CUDA implementation 8 months ago		
Padding.lua	fixed pull nn.Padding() if self returned in backward 5 months ago		
ParwiseDistance.lua	Merge pull request #532 from gigahester 29 days ago		
Parallel.lua	fix a bug in conditional expression a month ago		
ParallelCriterion.lua	assets in MultiCriterion and ParallelCriterion add 2 months ago		
ParallelTable.lua	Parallel optimization. ParallelTable inherits Container, unit tests a year ago		
Power.lua	Use UNIX line endings 7 months ago		
README.md	docs readability 5 months ago		
RReLU.lua	Add randomized leaky rectified linear unit (RReLU) 3 months ago		
ReLU.lua	adds in-place ReLU and fixes a potential divide-by-zero in nn.Sqz 9 months ago		
Replicate.lua	Replicate table modules 8 months ago		
Reshape.lua	Added more informative pretty-printing. a year ago		
Select.lua	initial levancy of torch tree 4 years ago		
SelectTable.lua	nn.Module preserve type sharing semantics (#187); add nn.Module apply 4 months ago		
Sequential.lua	fixing Sequential.remove corner case 6 months ago		
Sigmoid.lua	initial levancy of torch tree 4 years ago		
SmoothL1Criterion.lua	Add ScaLverage to criterion in the constructor 2 months ago		
SoftMax.lua	Fix various unused variables in nn a year ago		
SoftMin.lua	Fix various unused variables in nn a year ago		
SoftPlus.lua	fixed a numerical issue in the SoftPlus module (it breaks for input g... 2 years ago		
SoftShrink.lua	initial levancy of torch tree 4 years ago		
SoftSign.lua	initial levancy of torch tree 4 years ago		
SparseJacobian.lua	Fix various unused variables in nn a year ago		
SpatialLinear.lua	Using sparse implementation of zeroGradParameters for SparseLinear a month ago		
SpatialMaxPooling.lua	Added SpatialAdaptiveMaxPooling 4 years ago		
SpatialAdaptiveMaxPooling.lua	SpatialAveragePooling supports padding, cell mode and exclude_pad div... 29 days ago		
SpatialBatchNormalization.lua	Add C implementation of SpatialBatchNormalization 7 days ago		
SpatialConvolutionNormalize.lua	Make type# truly recursive. 9 months ago		
SpatialConvolution.lua	Fix type# in SpatialConvolution 3 months ago		
SpatialConvolutionMM.lua	Fix type# in SpatialConvolution 3 months ago		
SpatialConvolutionMap.lua	Remove unused and expensive initialization logic from nn.SpatialConv... 8 months ago		
SpatialCrossMagGRN.lua	cuda consistency 18 days ago		
SpatialConvolutionNormalize...	SpatialConvolutionDivisiveSubtractiveNormalization work with bat...8 months ago		
SpatialDropout.lua	small fix on error message 6 months ago		
SpatialFractionalMaxPooling...	Adding Fractional Max Pooling 3 months ago		
SpatialFractionalConvolution...	Add adjustment term to SpatialFractionalConvolution to control the size of... 5 days ago		
SpatialFractionalConvolutionMap.lua	New NN classes 3 years ago		
SpatialFPHPooling.lua	SpatialMaxPooling divides by KNNH 10 months ago		
SpatialMaxPooling.lua	SpatialMaxPooling supports padding and cell mode 6 months ago		
SpatialMaxUnpooling.lua	Add SpatialMaxUnpooling 28 days ago		
SpatialSoftMax.lua	Update SoftMax to work in spatial mode 4 months ago		
SpatialSoftSigmoid.lua	Merge branch 'trt_fix_revert' 3 years ago		
SpatialSubActiveNormalize...	SpatialConvolutionDivisiveSubtractiveNormalization work with bat...8 months ago		
SpatialSamplingNearest...	Use UNIX line endings 7 months ago		
SpatialZeroPadding.lua	Added more informative pretty-printing. a year ago		
SplitTable.lua	Add support for negative indices in nn.SplitTable 7 months ago		

# Example: Torch Layers

torch / nn		Watch	s2	star	278	fork	307
Code	Issues 27	Pull requests 11	Wiki	Pulse	Graphs		
No description or website provided.							
1,039 commits	7 branches	releases	88 contributors				
Branch: master · <a href="#">New pull request</a>	<a href="#">New file</a>	<a href="#">Find file</a>	<a href="#">HTTPS</a> · <a href="https://github.com/search">https://github.com/search</a>	<a href="#">Download ZIP</a>			
<b>sohumit Merge pull request #603 from torchrevert/563-master</b>	LATEST COMMIT 23d61d 15 hours ago						
diff	Fix batch mode in MarginRankingCriterion	4 days ago					
genetic	Improve error message in SpatialConvolutionMM	a day ago					
lb	THNN: add missing OpenMP include	2 days ago					
rocks	Add 'tunf' dependency	14 days ago					
gignore	git ls ignore build output	4 months ago					
luachack	[Torch] Move test to the top level	a year ago					
travis.yml	small fix to test paths	2 months ago					
abs.lua	Add THNN conversion of {ELU, LeakyReLU, LogSigmoid, LogSoftMax, Look...}	7 days ago					
AbsCriterion.lua	Add THNN conversion of {ELU, LeakyReLU, LogSigmoid, LogSoftMax, Look...}	7 days ago					
Add.lua	fix Add with multi-dim bias	10 months ago					
AddConstant.lua	Adding in-place AddConstant and MultiConstant	9 months ago					
BCECriterion.lua	Remove unnecessary softmax from BCECriterion	3 months ago					
BatchNormalization.lua	fix batchnorm reset	3 months ago					
CAddTable.lua	Fixing table modules to return correct number of gradients	6 months ago					
CDNTable.lua	Fixing table modules to return correct number of gradients	6 months ago					
CMACLoss.tst	Add C implementation of SpatialBatchNormalization	7 days ago					
CMU.lua	nn.Module preserve type sharing semantics (#187); add nn.Module.apply	4 months ago					
CMUTable.lua	fixing table modules to return correct number of gradients	6 months ago					
CONTRIBUTING.md	added developing tips	3 months ago					
COPYRIGHT.rst	add copyright file	2 years ago					
CScaleTable.lua	fixing table modules to return correct number of gradients	6 months ago					
Clip.lua	Use custom range in HardTanh and mask it as Clip	3 months ago					
CrossNLLCriterion.lua	Add functional conversion of CrossNLLCriterion	13 days ago					
Concat.lua	fix a bug in conditional expression	a month ago					
ConcatTable.lua	bug fix in ConcatTable variable length	4 months ago					
Container.lua	Adding applyToModule() to nn.Container, which is like apply() but...	3 months ago					
Copy.lua	nn.Module preserve type sharing semantics (#187); add nn.Module.apply	4 months ago					
Create.lua	Fx typed in Create	a month ago					
CosineDistance.lua	Do not change state variables in CosineDistance/CosineEmbeddingCriterion	2 months ago					
CosineEmbeddingCriterion.lua	Do not change state variables in CosineDistance/CosineEmbeddingCriterion	2 months ago					
Criterion.lua	nn.Module preserve type sharing semantics (#187); add nn.Module.apply	4 months ago					
CriterionTable.lua	Rename unsafe table to unsafe for Luas 5.2	8 months ago					
CrossEntropyCriterion.lua	Check for nn.Module and nn.Criterion in recursive type.	8 months ago					
DepthConcurrent.lua	adding direct backward to Concurrent.DepthConcurrent, Sequential	9 months ago					
DepthConvolution.lua	Use tensor for THNN functions for certain element outputs	10 days ago					
DProdProduct.lua	Add batch mode in DProdProduct + unit test	2 months ago					
Dropout.lua	In-place dropout	4 months ago					
ELU.lua	Add THNN conversion of {ELU, LeakyReLU, LogSigmoid, LogSoftMax, Look...}	7 days ago					
ErrorMessages.lua	Give better error messages when trying to use the wrong kind of Tensor	a year ago					
Euclidean.lua	nn.Module preserve type sharing semantics (#187); add nn.Module.apply	4 months ago					
Exp.lua	Exp module fix only	9 months ago					
FlattenTable.lua	nn.Module preserve type sharing semantics (#187); add nn.Module.apply	4 months ago					
GradientReversal.lua	Add GradientReversal layer	4 months ago					
HardShrink.lua	Add functional conversion of HardShrink	10 days ago					
HardTanh.lua	Add functional conversion of HardTanh	10 days ago					
HingeEmbeddingCriterion.lua	remove HingeEmbeddingCriterion to support batch mode	6 months ago					
Identity.lua	Revert to previous identity loss implementation	2 months ago					
Index.lua	Simplifying and more efficient nn.Index	2 months ago					
Jacobian.lua	Add util tests for hessian.lua, to bugs detected by the tests.	6 months ago					
JoinTable.lua	nn.Module preserve type sharing semantics (#187); add nn.Module.apply	4 months ago					
LCatTable.lua	Use tensor for THNN functions even for single element outputs	10 days ago					
LinearEmbeddingCriterion.lua	Make type() truly recursive	9 months ago					
L1Penalty.lua	fix L1Penalty constructor arguments	a year ago					
LeakyReLU.lua	Add THNN conversion of {ELU, LeakyReLU, LogSigmoid, LogSoftMax, Look...}	7 days ago					
Linear.lua	Remove sparse matrics from inLinear	4 months ago					
LogSigmoid.lua	Add THNN conversion of {ELU, LeakyReLU, LogSigmoid, LogSoftMax, Look...}	7 days ago					
LogSoftMax.lua	Add THNN conversion of {ELU, LeakyReLU, LogSigmoid, LogSoftMax, Look...}	7 days ago					
LookupTable.lua	Harmonize LookupTable signature with cum input	8 months ago					
Mul.lua	Remove unsafe to batch for unsafe for Luas 5.2	8 months ago					
MSECriterion.lua	Add SoftAverage to criterion in the constructor	2 months ago					
MarginCriterion.lua	modernized MarginCriterion	a year ago					
MarginRankingCriterion.lua	Fix batch mode in MarginRankingCriterion	4 days ago					
Max.lua	Merge pull request #604 from vghemster	2 months ago					
Mean.lua	Add support for negative dimension and both batch and non batch input...	2 months ago					
Min.lua	Merge pull request #604 from vghemster	2 months ago					
MultiTable.lua	cancel unused variable and useless expression	29 days ago					
Module.lua	Revert "Don't re-factor parameters if they are already shared"	15 hours ago					
Mul.lua	removing the requirement for providing size in m.Mul	a year ago					
MulConstant.lua	Ignore update/gradient if selfградient is nil	3 months ago					
MultiCriterion.lua	assets in MultiCriterion and ParallelCriterion add	2 months ago					
MultLabMarginCriterion.lua	initial reavng of torch tree	4 years ago					
MultMarginCriterion.lua	multMargin supports p=2	11 months ago					
NewRow.lua	typical in Name not done in place	6 months ago					
NewRowTable.lua	NameTable	6 months ago					
Normalise.lua	Remove bmm and baddmm from Normalise, because they allocate memory....	20 days ago					
PRelu.lua	Buffers for PRelu CUDA implementation	8 months ago					
Padding.lua	fixed broken nn.Padding: was returned in backprop	5 months ago					
ParwiseDistance.lua	Merge pull request #602 from vghemster	29 days ago					
Parallel.lua	fix a bug in conditional expression	a month ago					
ParallelCriterion.lua	assets in MultiCriterion and ParallelCriterion add	2 months ago					
ParallelTable.lua	Parallel optimization. ParallelTable inherits Container, unit tests	a year ago					
Power.lua	Use UNIX line endings	7 months ago					
READEME.md	doc readme	5 months ago					
RReLU.lua	Add randomized leaky rectified linear unit (RReLU)	3 months ago					
ReLU.lua	adds in-place ReLU and uses a partial divide-by-zero in nn.Soft...	9 months ago					
Replace.lua	Replicate batchMode	8 months ago					
Reshape.lua	Add more informative pretty-printing	a year ago					
Select.lua	intel reavng of torch tree	4 years ago					
SelectTable.lua	nn.Module preserve type sharing semantics (#187); add nn.Module.apply	4 months ago					
Sequential.lua	fix Sequential remove corner case	6 months ago					
SignGrad.lua	initial reavng of torch tree	4 years ago					
Smooth1Criterion.lua	Add Smooth1 to criterion in the constructor	2 months ago					
SoftMax.lua	Fix various unused variables in nn	a year ago					
SoftMin.lua	Fix various unused variables in nn	a year ago					
SoftPlus.lua	feed a numerical result in the SoftPlus module (it breaks for input g...	2 years ago					
SoftShrink.lua	intel reavng of torch tree	4 years ago					
SoftSign.lua	intel reavng of torch tree	4 years ago					
Sparsesocabian.lua	Fix various unused variables in nn	a year ago					
SpatialLinear.lua	Using sparse implementation of nn.SparseParameters for SparseLinear	a month ago					
SpatialAdaptiveMaxPooling.lua	Add SpatialAdaptiveMaxPooling	a year ago					
SpatialAveragingPooling.lua	SpatialAveragingPooling supports padding, cell mode and exclude, pad, div...	29 days ago					
SpatialBatchNormalizable.lua	AD C implementation of SpatialBatchNormalizable	7 days ago					
SpatialContrastiveNormalizat...	Make type() truly recursive	9 months ago					
SpatialConvolution.lua	Fix type() in SpatialConvolution	3 months ago					
SpatialConvolutionMM.lua	Fix type() in SpatialConvolution	3 months ago					
SpatialConvolutionMap.lua	Remove unused and expensive initialization logic from nn.SpatialConv...	8 months ago					
SpatialCrossPoolRN.lua	code consistency	18 days ago					
SpatialCrossNormKd.lua	SpatialConvolution, Divide, SubtractiveNormalization work with bat...	8 months ago					
SpatialDropout.lua	small fix in error message	6 months ago					
SpatialFactorialPooling.lua	Adding Factorial Pooling	3 months ago					
SpatialFullConvolution.lua	Add adjustment term to SpatialFullConvolution to control the size of...	3 days ago					
SpatialConvolutionMap.lua	New NN classes	3 years ago					
SpatialPhasing.lua	SpatialAveragingPooling divides by N^M^H	10 months ago					
SpatialPooling.lua	SpatialAveragingPooling supports padding and cell mode	6 months ago					
SpatialReLU.lua	Add SpatialReLU	25 days ago					
SpatialSoftmax.lua	Update Softmax to work in spatial mode	4 months ago					
SpatialSubsampling.lua	Merge branch 'nn_fuse_neal'	3 years ago					
SpatialUnbatchNorm.lua	SpatialConvolution, Divide, SubtractiveNormalization work with bat...	8 months ago					
SpatialUnbatchNormNeal...	Use UNIX line endings	7 months ago					
SpatialZeroPadding.lua	Added more informative pretty-printing	a year ago					
SpatialTieTable.lua	Add support for negative indices in nn.TieTable	7 months ago					



# Fei-Fei Li & Andrej Karpathy & Justin Johnson

# Lecture 4 - 50

13 Jan 2016

# Example: Torch MulConstant

$$f(X) = aX$$

initialization

forward()

backward()

```
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
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
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
```

# Example: Caffe Layers

# Caffe Sigmoid Layer

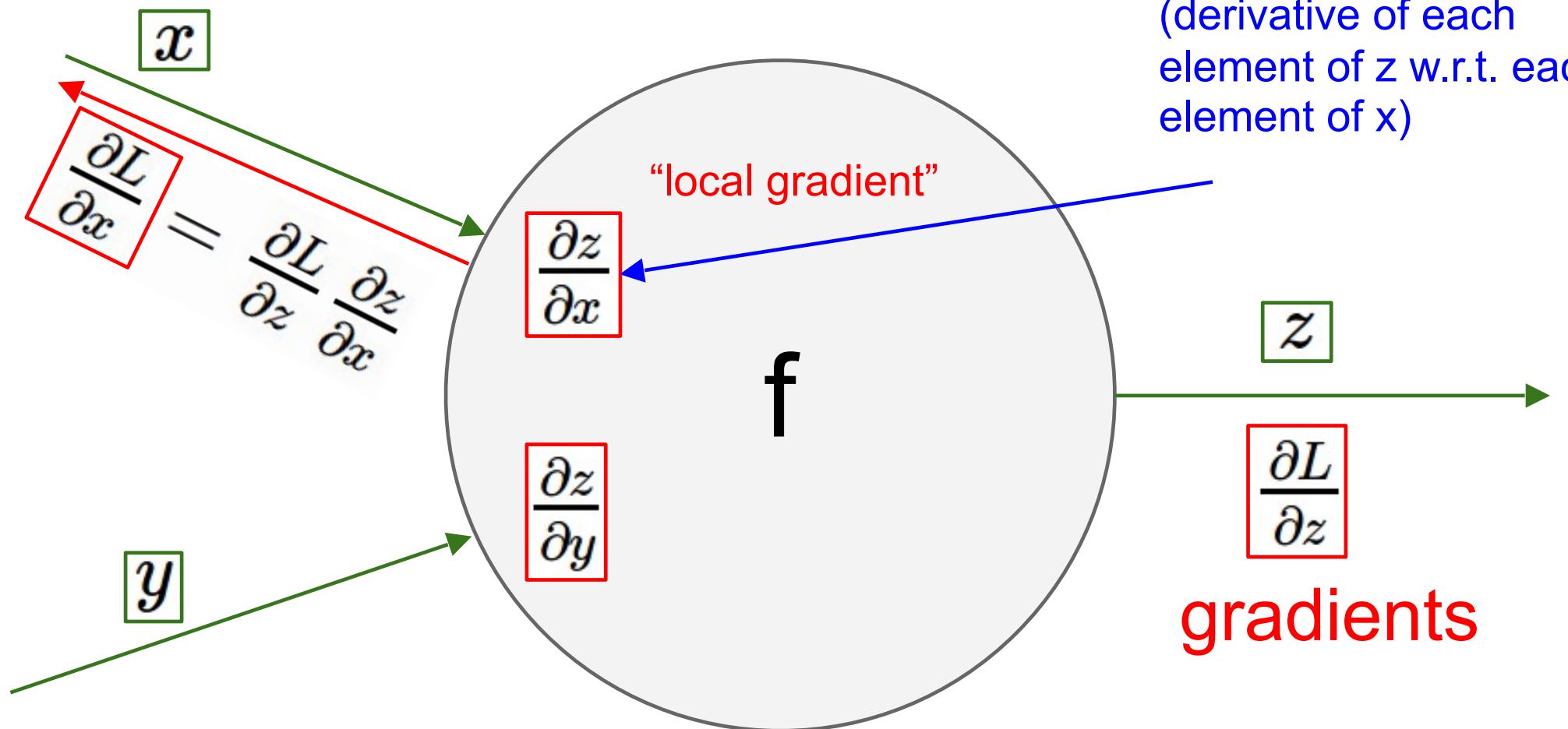
```
1 #include <cmath>
2 #include <vector>
3
4 #include "caffe/layers/sigmoid_layer.hpp"
5
6 namespace caffe {
7
8     template <typename Dtype>
9     inline Dtype sigmoid(Dtype x) {
10         return 1. / (1. + exp(-x));
11     }
12
13     template <typename Dtype>
14     void SigmoidLayer<Dtype>::Forward_cpu(const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

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

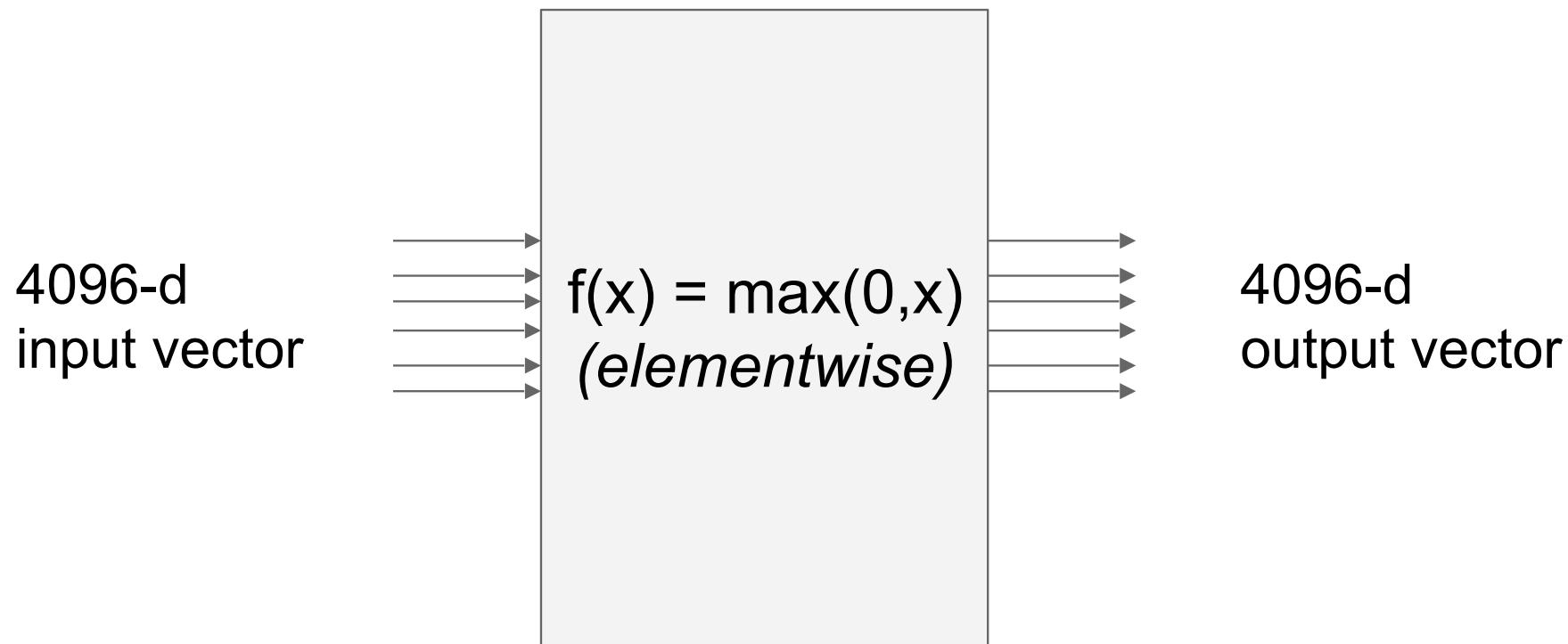
$$(1 - \sigma(x)) \sigma(x)$$

\*top\_diff (chain rule)

## Gradients for vectorized code (x,y,z are now vectors)



# Vectorized operations

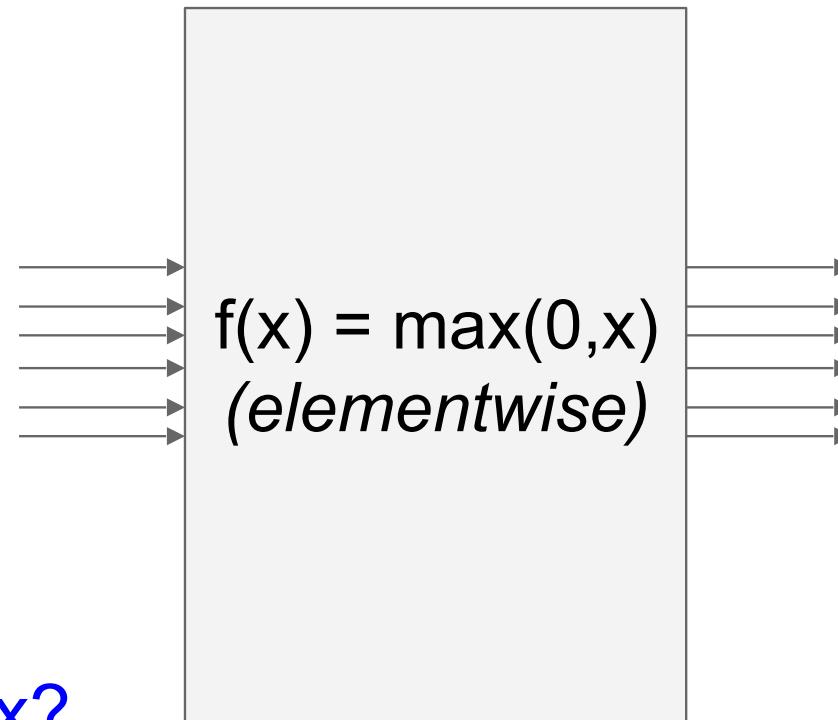


## Vectorized operations

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

Jacobian matrix

4096-d  
input vector



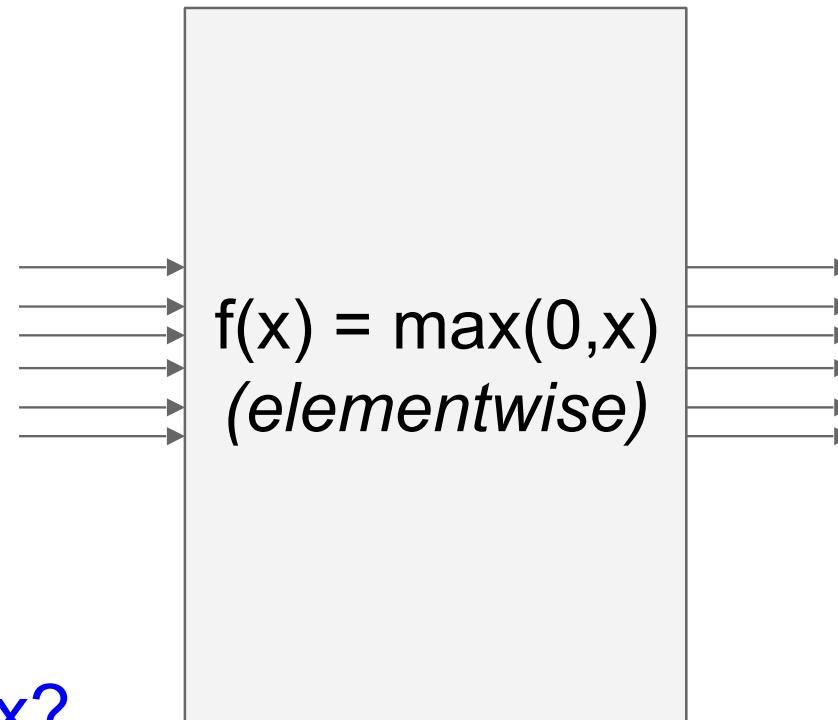
Q: what is the  
size of the  
Jacobian matrix?

## Vectorized operations

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

Jacobian matrix

4096-d  
input vector



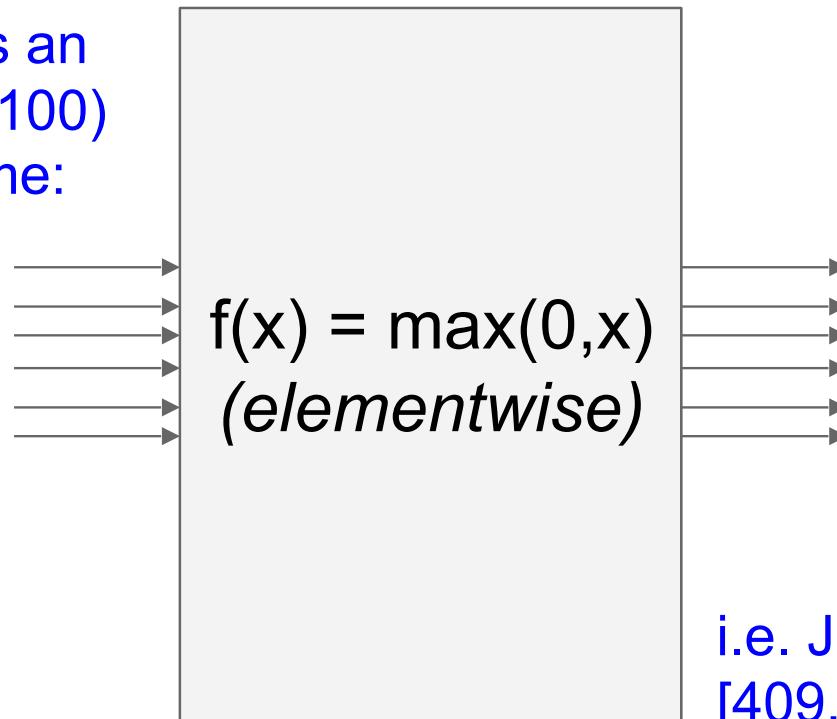
Q: what is the  
size of the  
Jacobian matrix?  
[4096 x 4096!]

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



100 4096-d  
output vectors

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

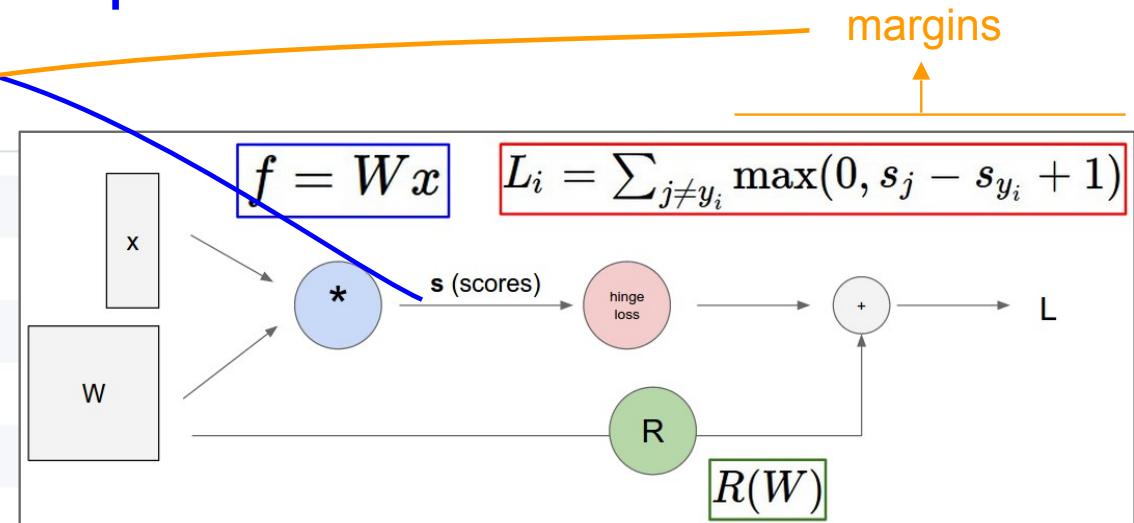
# Assignment: Writing SVM/Softmax

## 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()** 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: without the brain stuff

(Before) Linear score function:

$$f = Wx$$

# Neural Network: without the brain stuff

**(Before)** Linear score function:

$$f = Wx$$

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

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

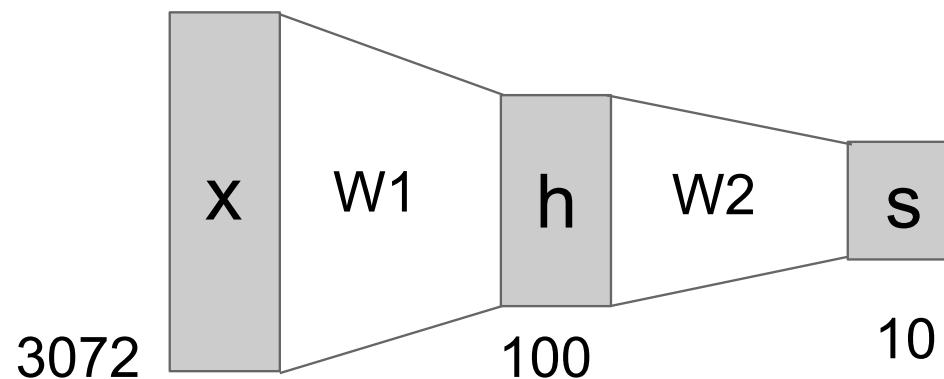
# Neural Network: without the brain stuff

**(Before)** Linear score function:

$$f = Wx$$

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

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



# Neural Network without the brain stuff

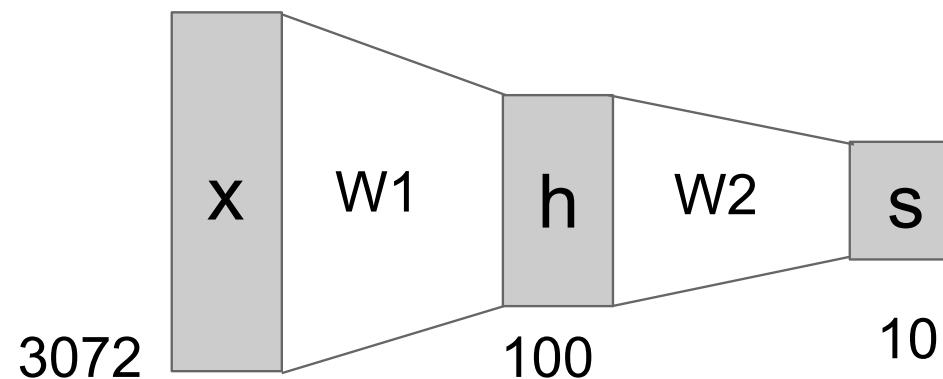


(Before) Linear score function:

$$f = Wx$$

(Now) 2-layer Neural Network

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



# Neural Network: without the brain stuff

**(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)
```

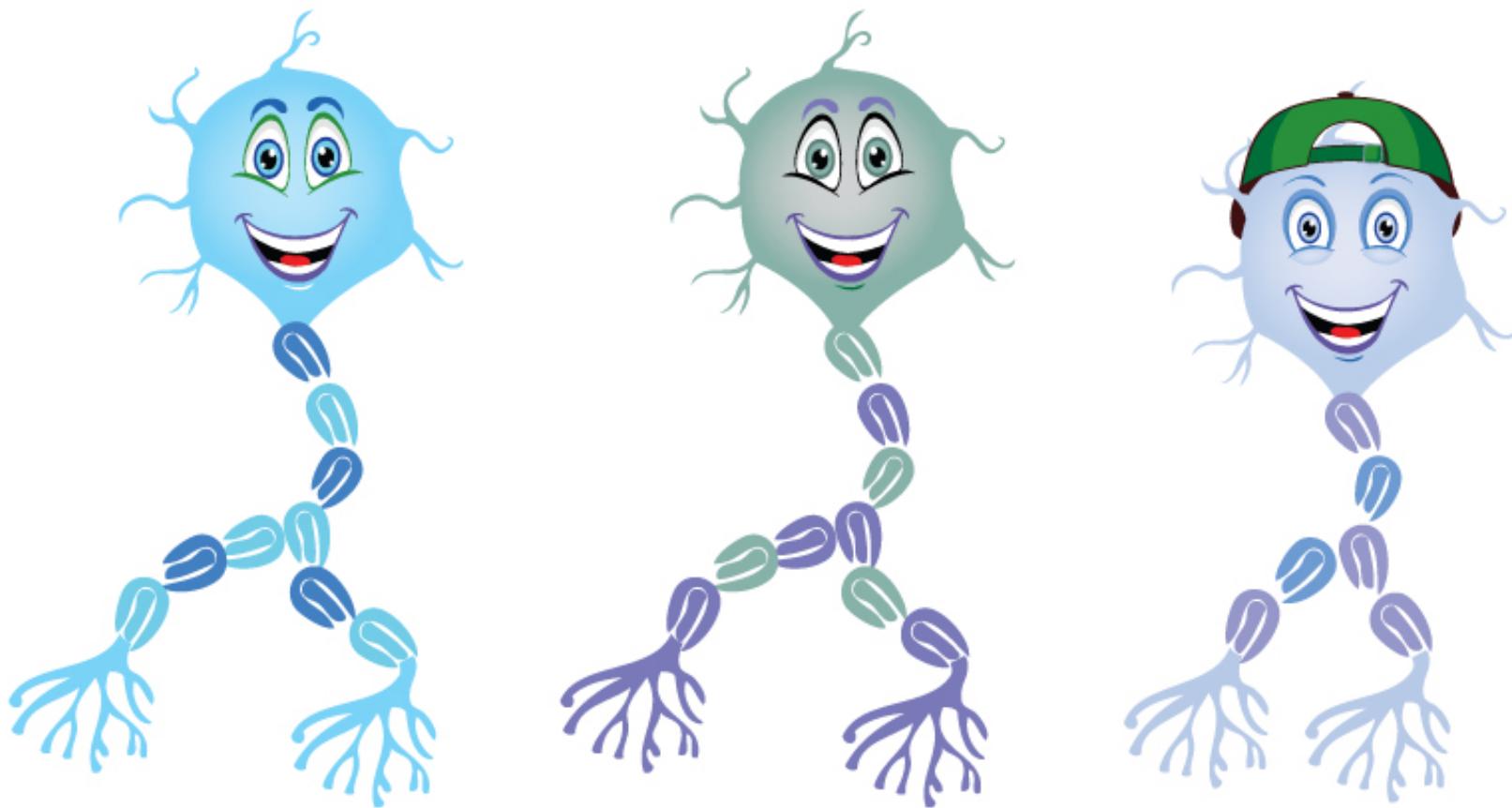
from @iamtrask, <http://iamtrask.github.io/2015/07/12/basic-python-network/>

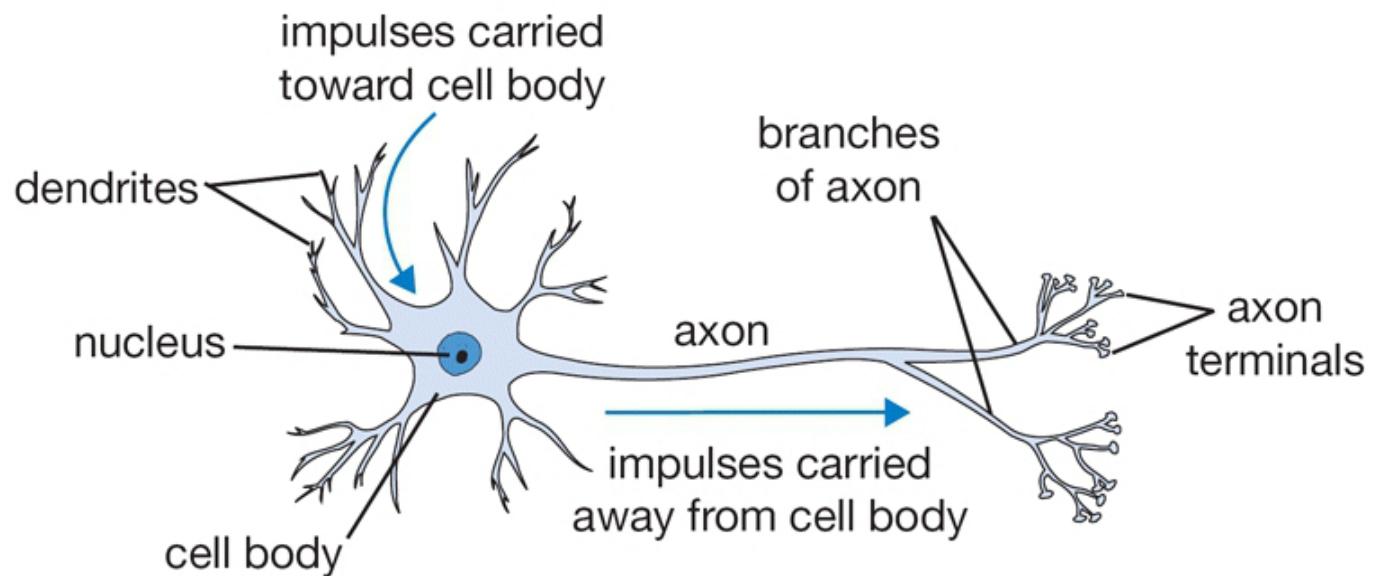
# Assignment: Writing 2layer Net

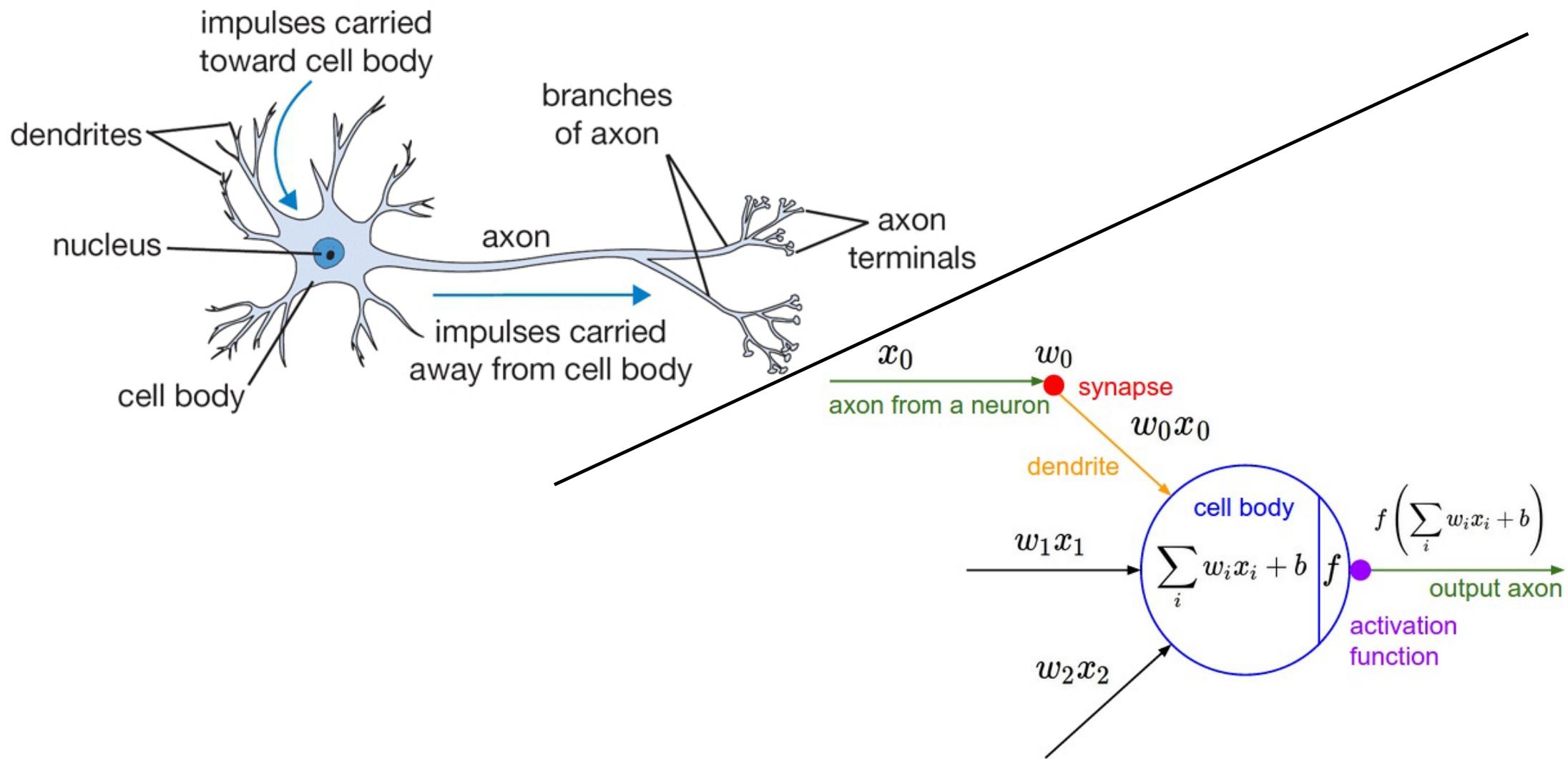
## Stage your forward/backward computation!

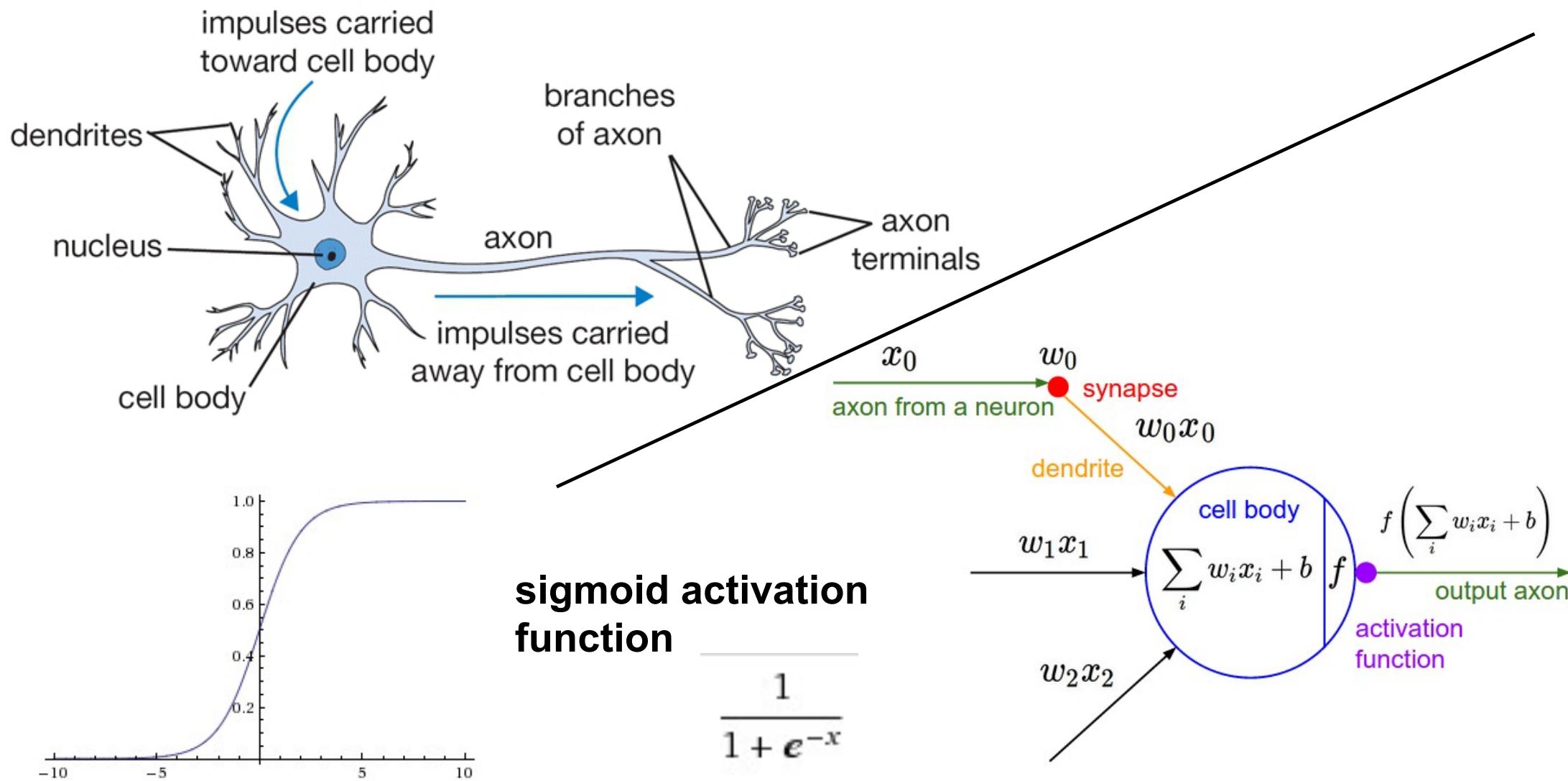
```
# receive W1,W2,b1,b2 (weights/biases), X (data)
# forward pass:
h1 = #... function of X,W1,b1
scores = #... function of h1,W2,b2
loss = #... (several lines of code to evaluate Softmax loss)
# backward pass:
dscores = #...
dh1,dW2,db2 = #...
dW1,db1 = #...
```

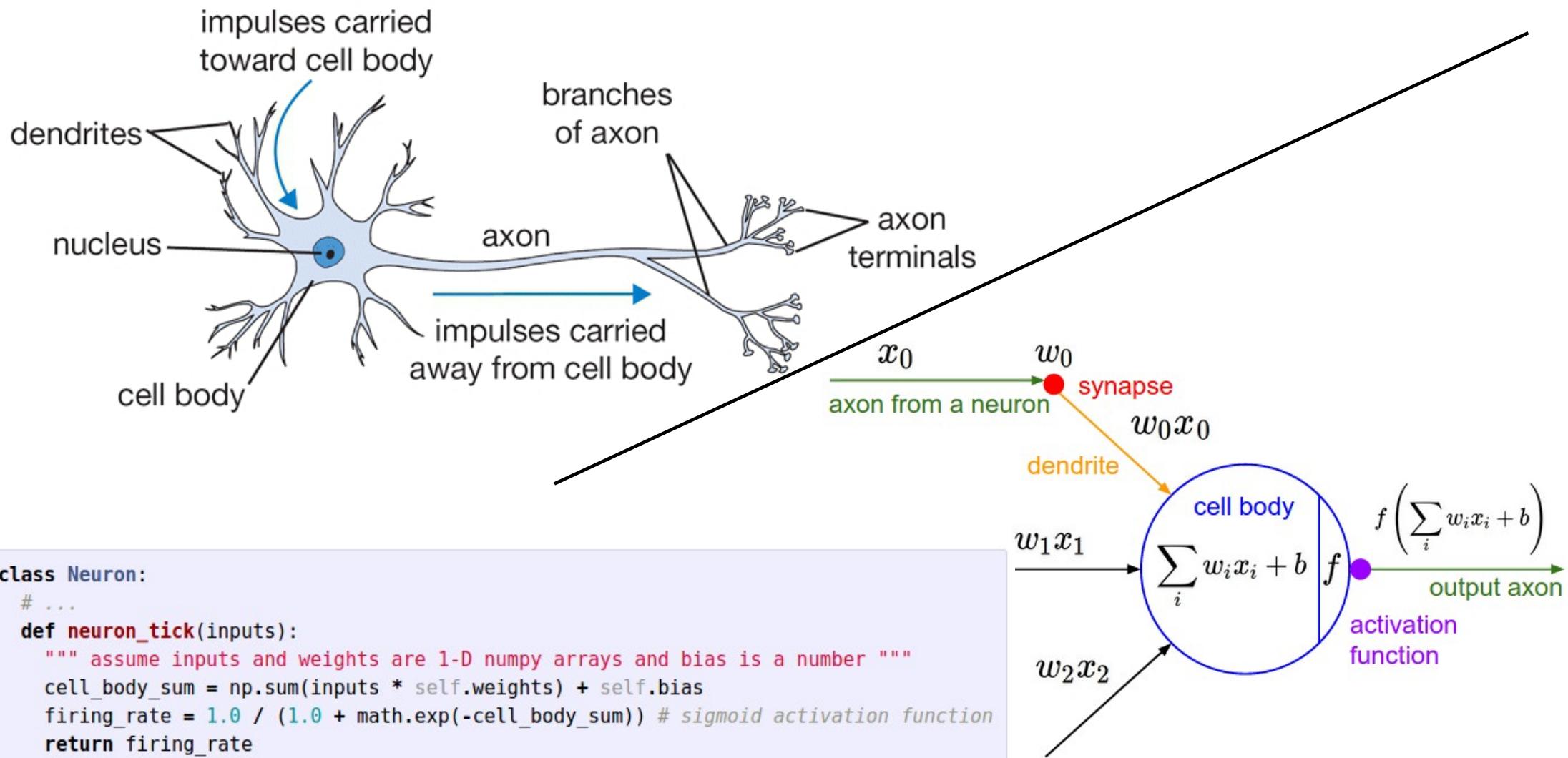








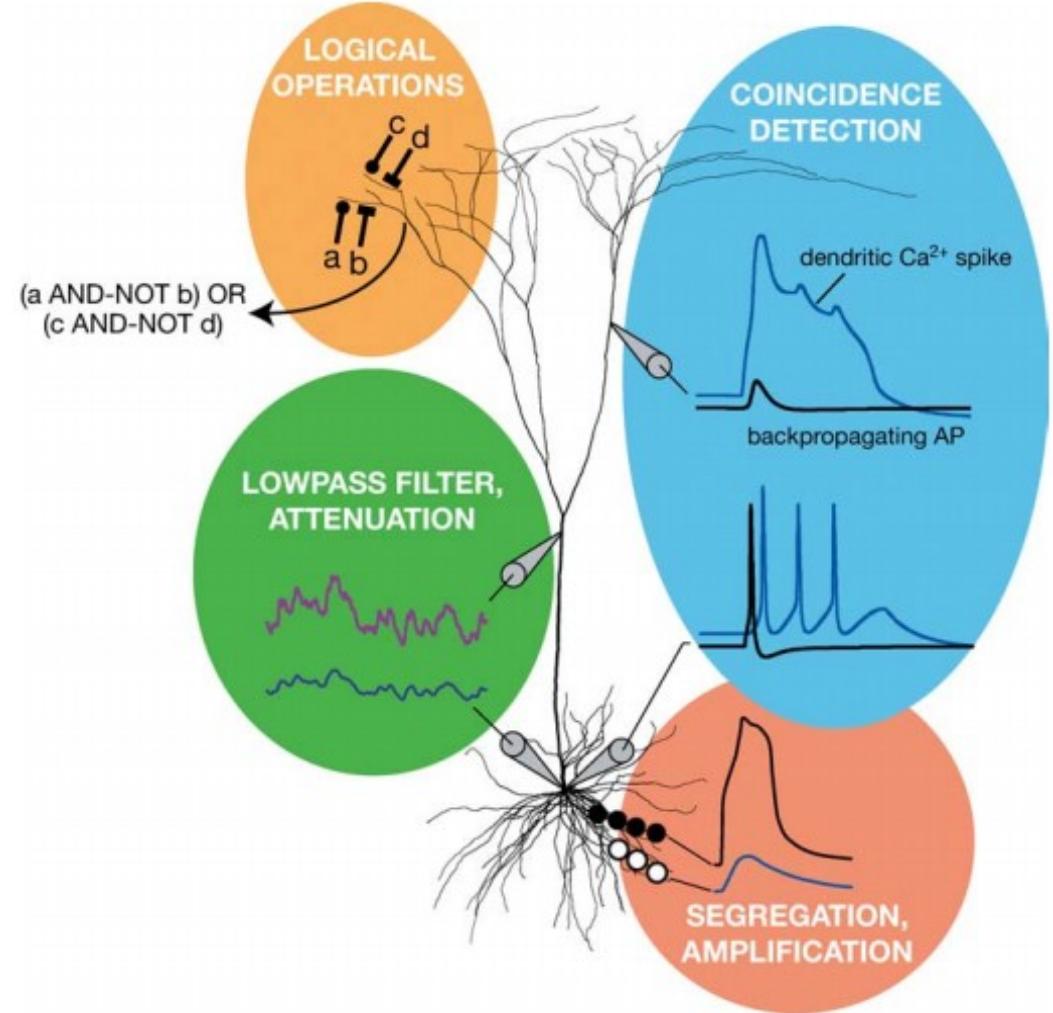




Be very careful with your Brain analogies:

### 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

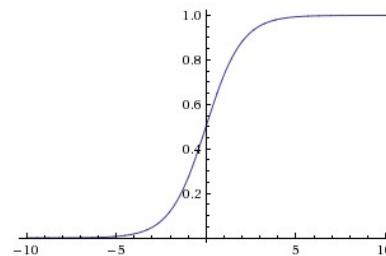


[*Dendritic Computation*. London and Häusser]

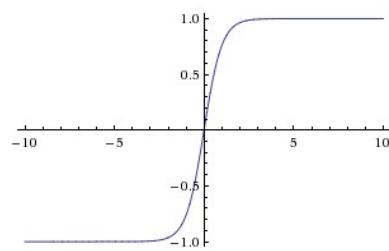
# Activation Functions

## Sigmoid

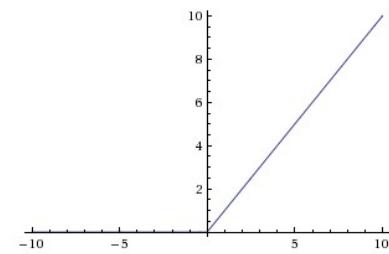
$$\sigma(x) = 1/(1 + e^{-x})$$



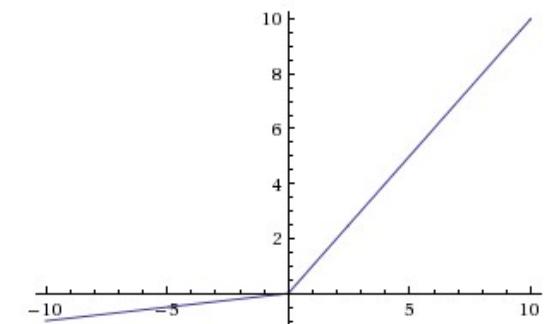
## tanh      tanh(x)



## ReLU      max(0,x)



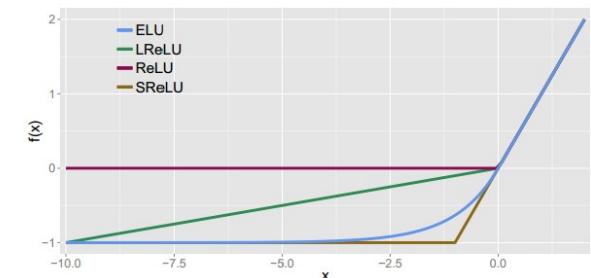
## Leaky ReLU $\max(0.1x, x)$



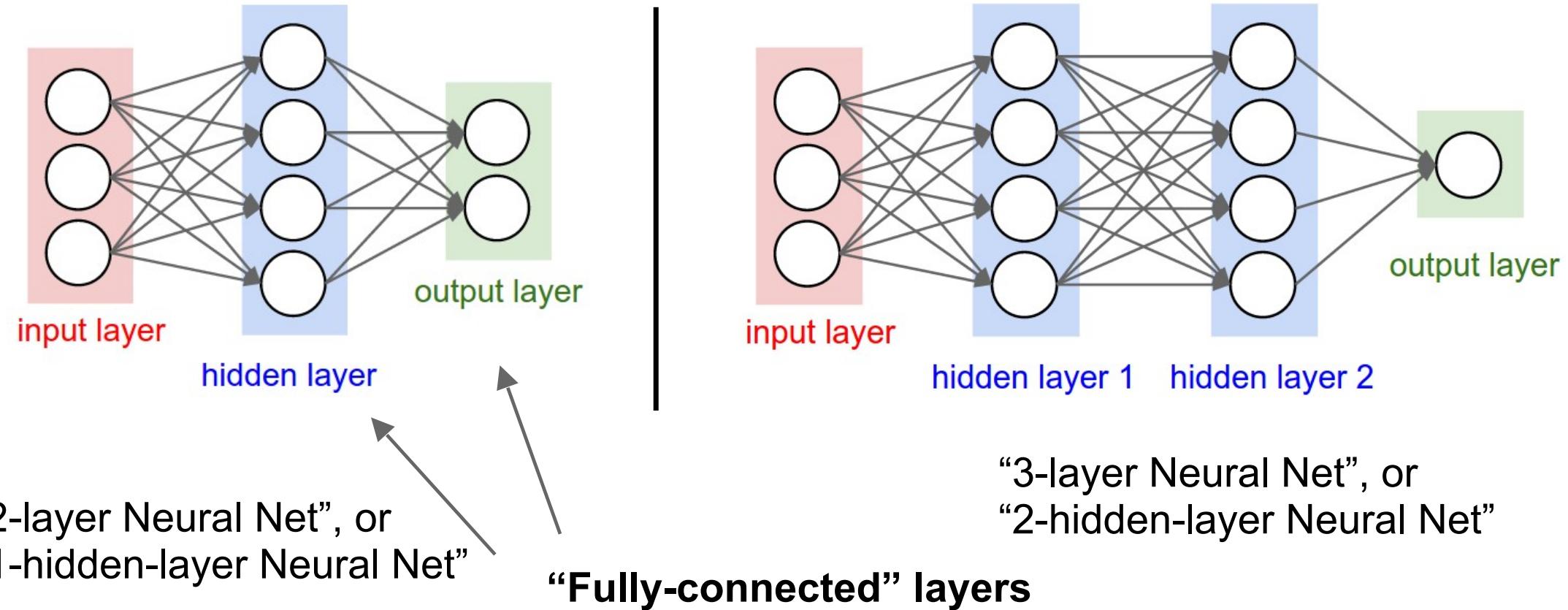
## Maxout      $\max(w_1^T x + b_1, w_2^T x + b_2)$

## ELU

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



# Neural Networks: Architectures

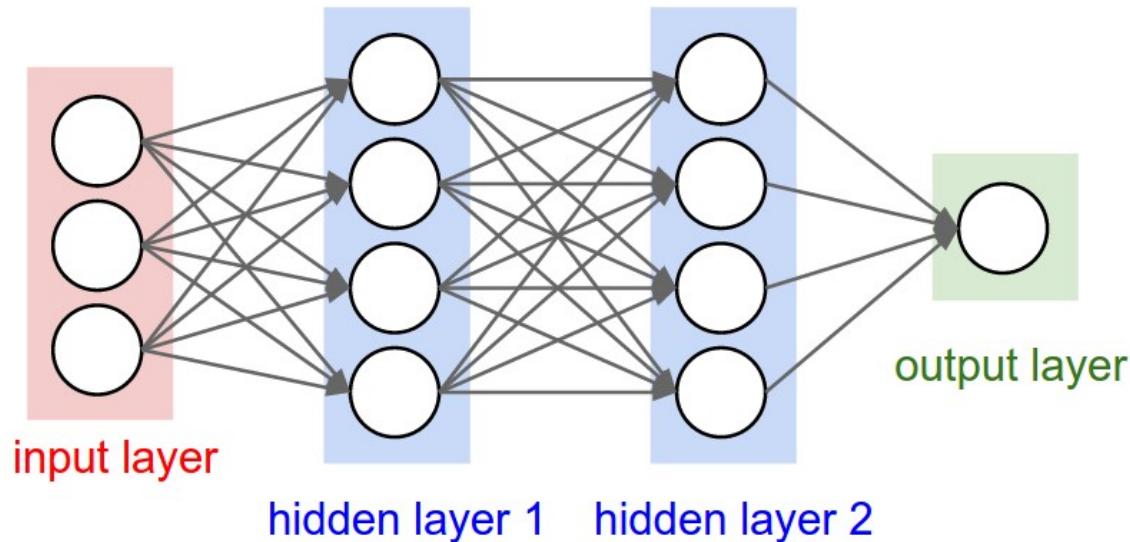


# Example Feed-forward computation of a Neural Network

```
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
```

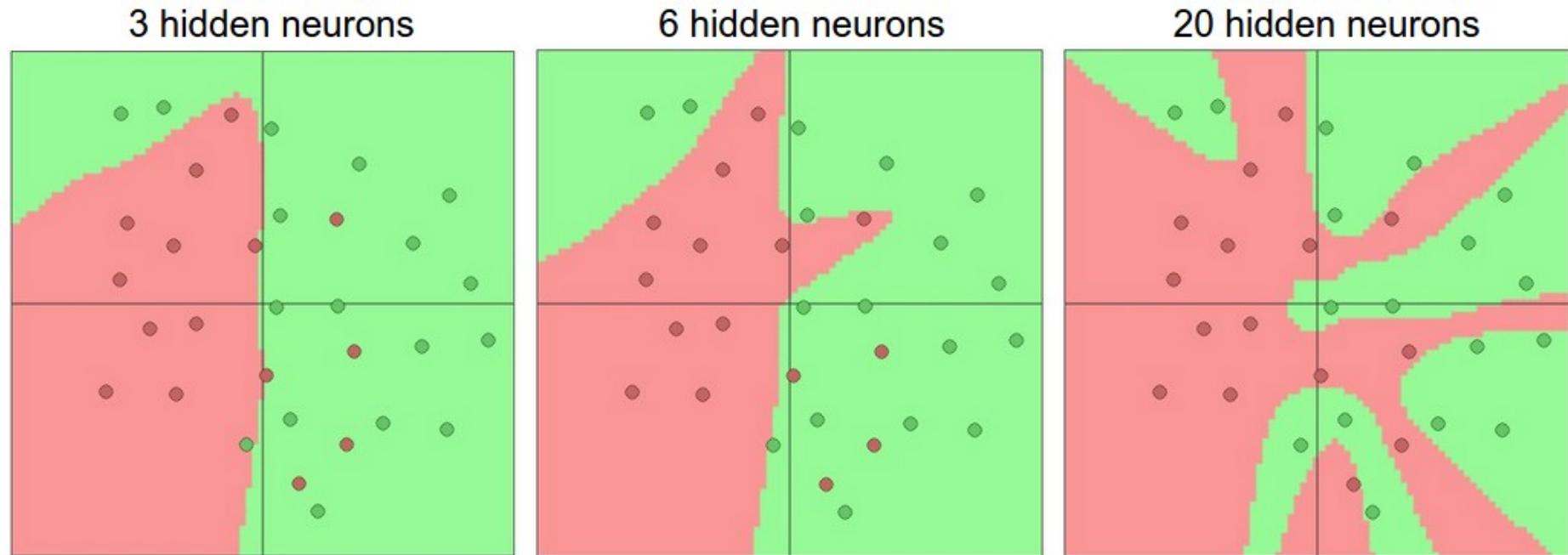
We can efficiently evaluate an entire layer of neurons.

# 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



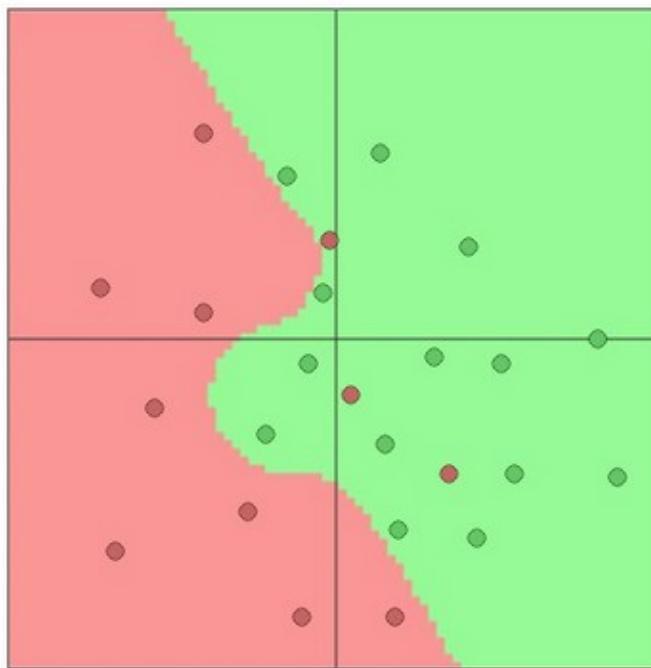
more neurons = more capacity

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

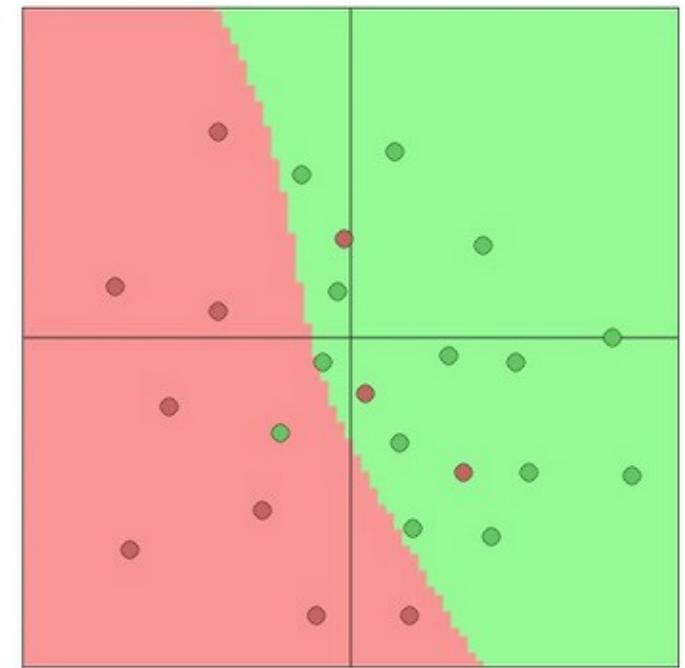
$\lambda = 0.001$



$\lambda = 0.01$



$\lambda = 0.1$



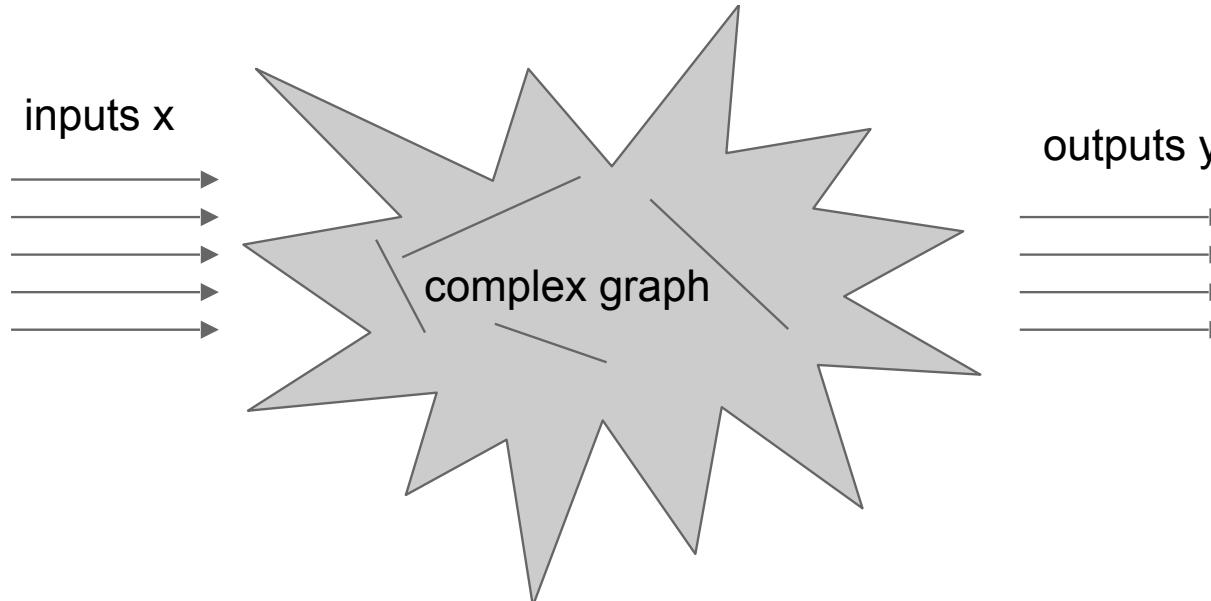
(you can play with this demo over at ConvNetJS: <http://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>)

# 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.



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

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

forward-mode differentiation (if you want effect of one thing on many things)

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