

Lecture 3 recap

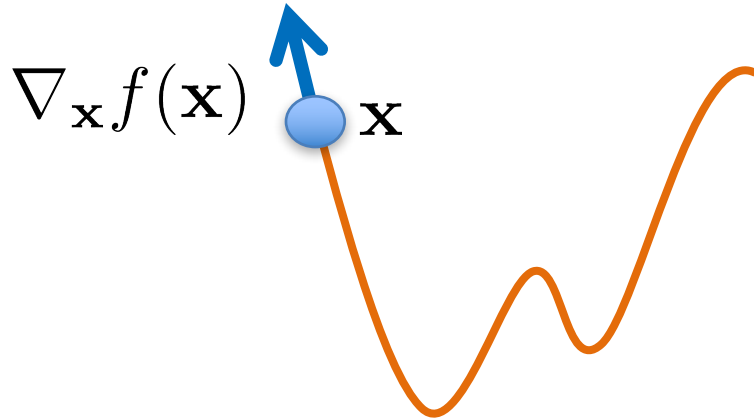
Gradient Descent

- From derivative to gradient

$$\frac{df(x)}{dx} \longrightarrow \nabla_{\mathbf{x}} f(\mathbf{x})$$

Direction of
greatest
increase of
the function

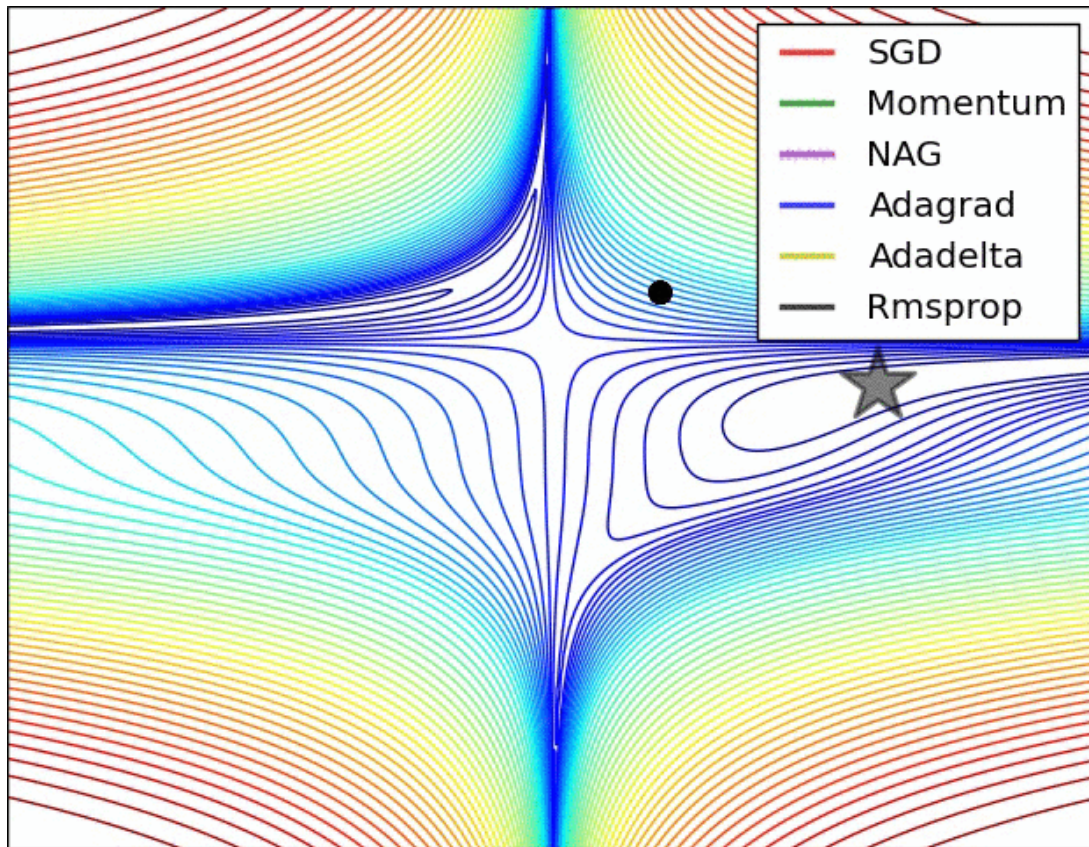
- Gradient steps in direction of negative gradient



$$\mathbf{x}' = \mathbf{x} - \epsilon \nabla_{\mathbf{x}} f(\mathbf{x})$$

Learning rate

Convergence



Numerical vs Analytical Gradient

Given the function $f(x) = \frac{1}{x}$
-> compute $f'(1.0)$

Numerical vs Analytical Gradient

Given the function $f(x) = \frac{1}{x}$
-> compute $f'(1.0)$

Numerical gradient:

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

for $h = 0.001$

$$f'(1.0) = \frac{f(1.001) - f(1.0)}{0.001} = -0.999$$

Numerical vs Analytical Gradient

Given the function $f(x) = \frac{1}{x}$

-> compute $f'(1.0)$

Numerical gradient:

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

for $h = 0.001$

$$f'(1.0) = \frac{f(1.001) - f(1.0)}{0.001} = -0.999$$

Analytical gradient:

$$\frac{d}{dx}f(x) = -\frac{1}{x^2}$$

$$f'(1.0) = -\frac{1}{1.0^2} = -1$$

Numerical vs Analytical Gradient

Given the function $f(x) = \frac{1}{x}$

-> compute $f'(1.0)$

Numerical gradient:

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

for $h = 0.001$

$$f'(1.0) = \frac{f(1.001) - f(1.0)}{0.001} = -0.999$$

Slow ☹️, approximate ☹️, easy-to-write 😊

Analytical gradient:

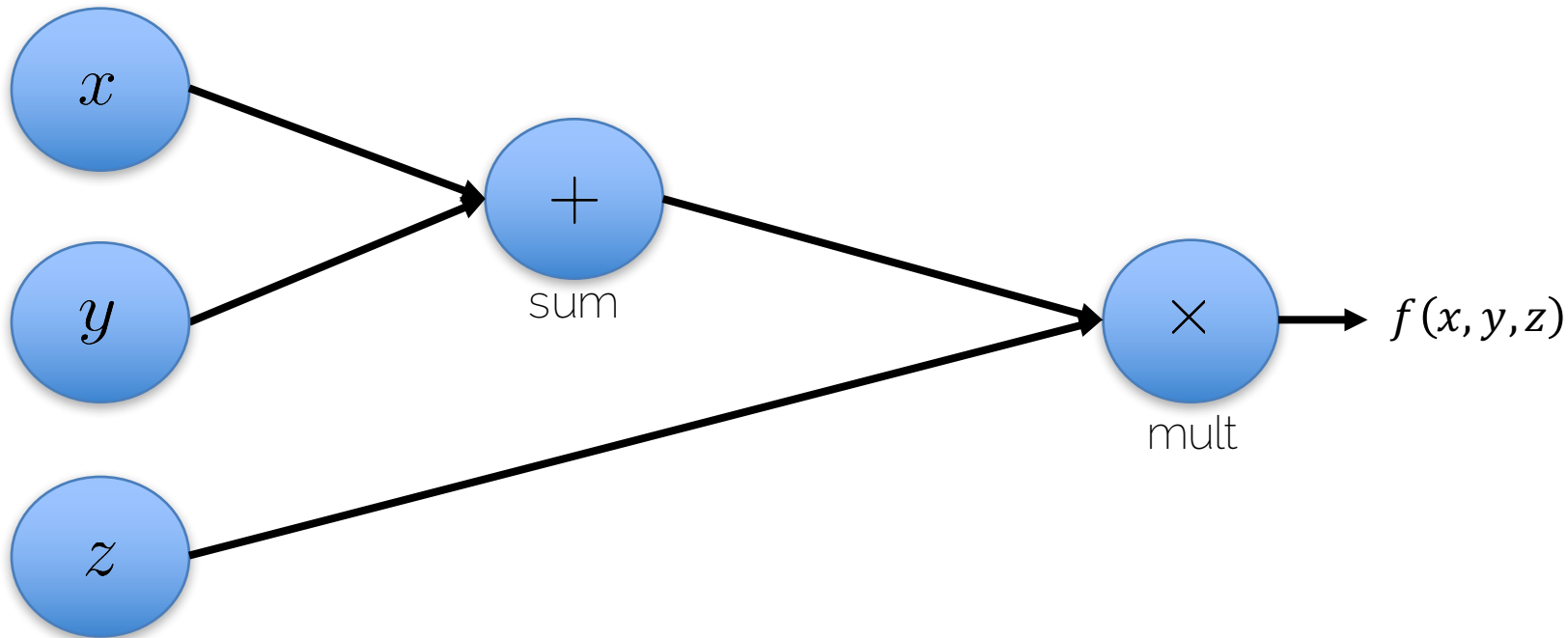
$$\frac{d}{dx}f(x) = -\frac{1}{x^2}$$

$$f'(1.0) = -\frac{1}{1.0^2} = -1$$

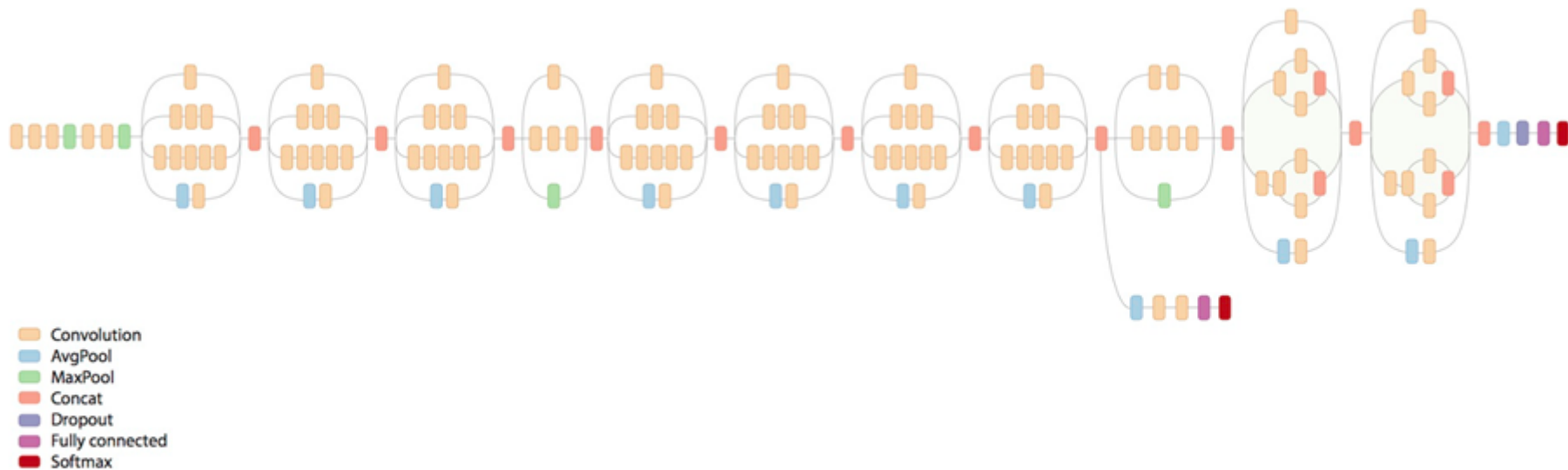
Fast 😊, exact* 😊, error-prone ☹️

Computational Graphs

- $f(x, y, z) = (x + y) \cdot z$



Computational Graphs



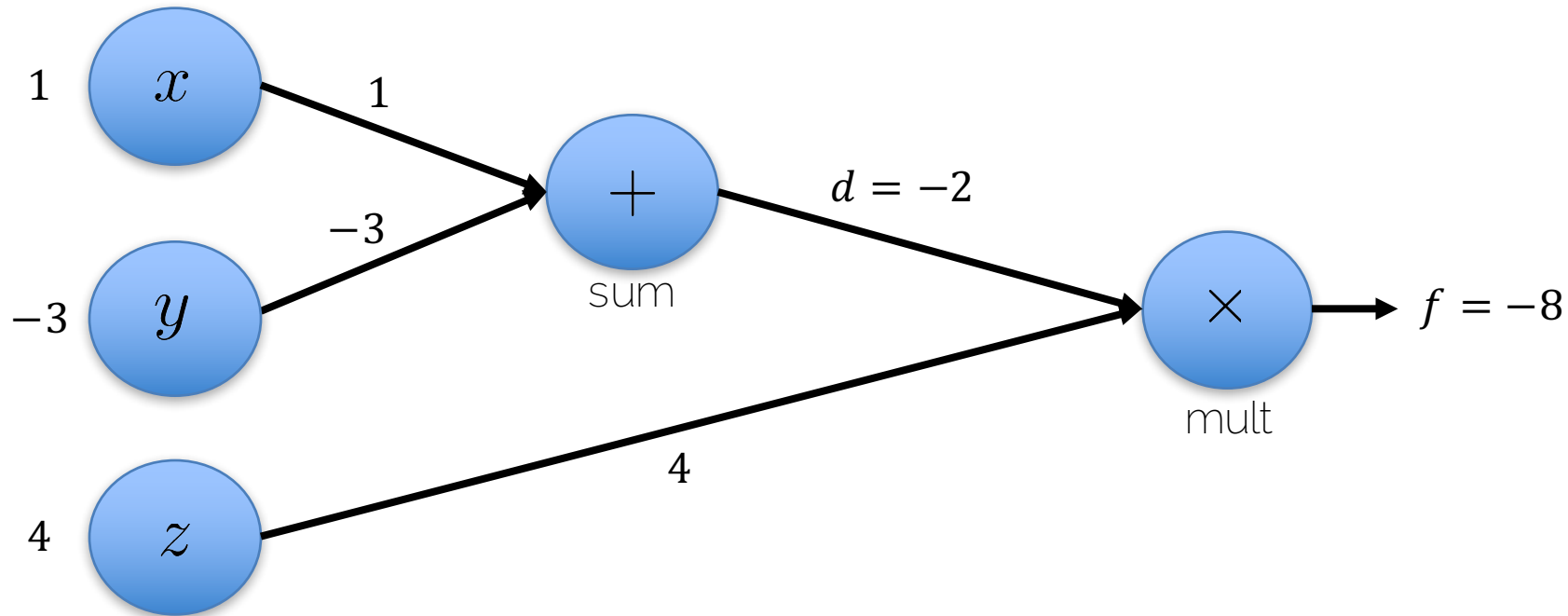
Another view of GoogLeNet's architecture.

Backprop: cont'd

Backprop: Forward Pass

- $f(x, y, z) = (x + y) \cdot z$

Initialization $x = 1, y = -3, z = 4$



Backprop: Backward Pass

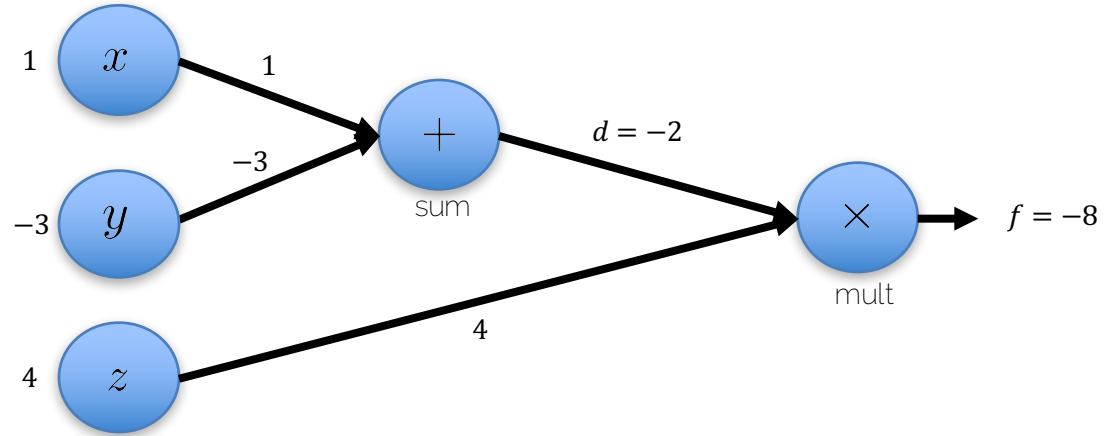
$$f(x, y, z) = (x + y) \cdot z$$

with $x = 1, y = -3, z = 4$

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

$$f = d \cdot z \quad \frac{\partial f}{\partial d} = z, \frac{\partial f}{\partial z} = d$$

What is $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$?



Backprop: Backward Pass

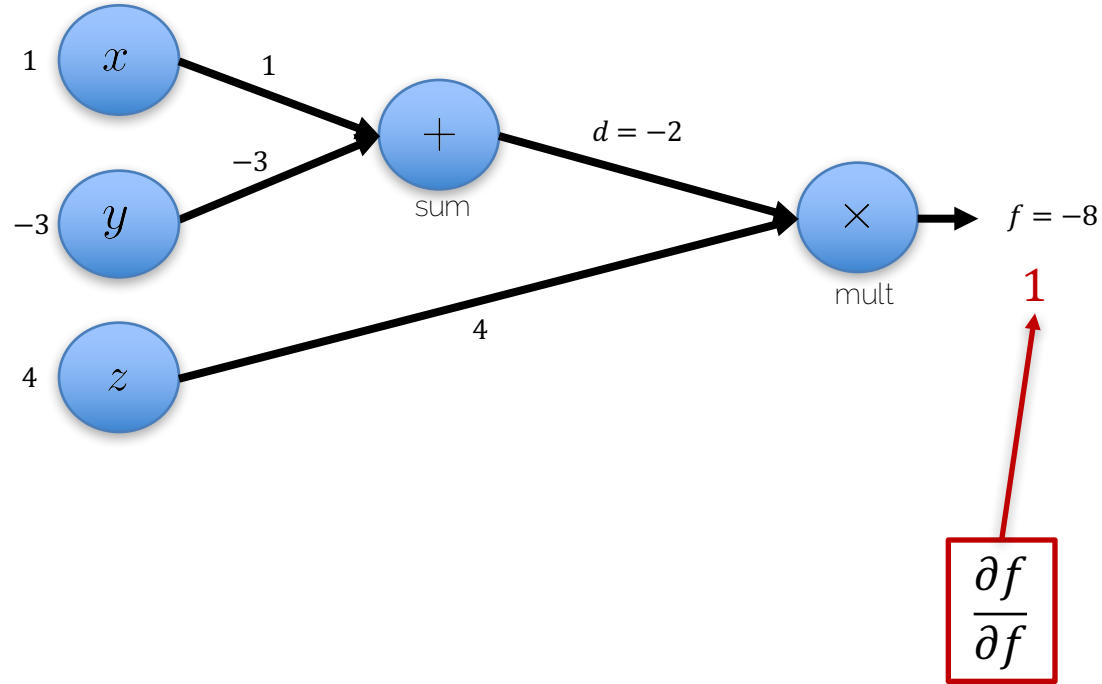
$$f(x, y, z) = (x + y) \cdot z$$

with $x = 1, y = -3, z = 4$

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

$$f = d \cdot z \quad \frac{\partial f}{\partial d} = z, \frac{\partial f}{\partial z} = d$$

What is $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$?



Backprop: Backward Pass

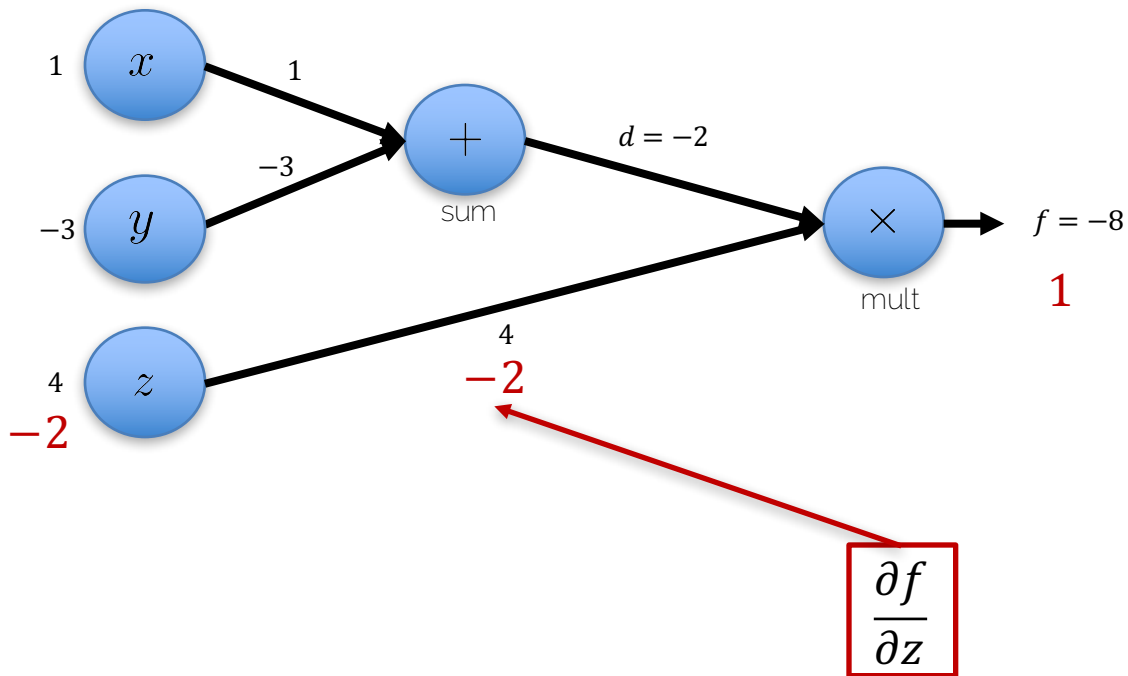
$$f(x, y, z) = (x + y) \cdot z$$

with $x = 1, y = -3, z = 4$

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

$$f = d \cdot z \quad \frac{\partial f}{\partial d} = z, \frac{\partial f}{\partial z} = d$$

What is $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$?



Backprop: Backward Pass

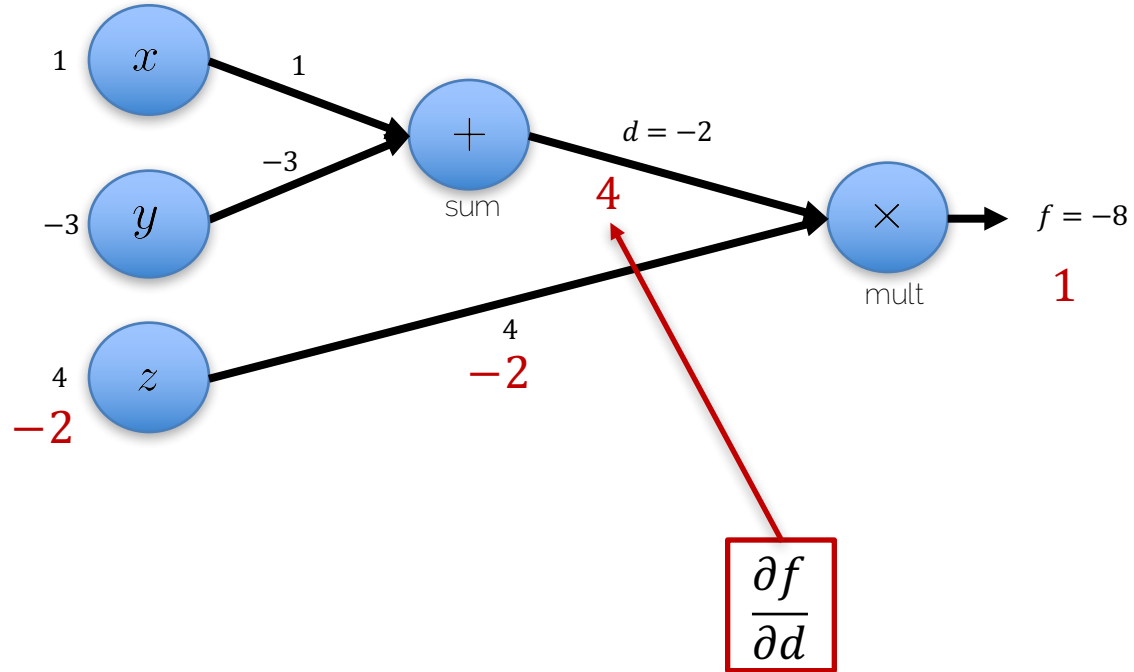
$$f(x, y, z) = (x + y) \cdot z$$

with $x = 1, y = -3, z = 4$

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

$$f = d \cdot z \quad \frac{\partial f}{\partial d} = z, \frac{\partial f}{\partial z} = d$$

What is $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$?



Backprop: Backward Pass

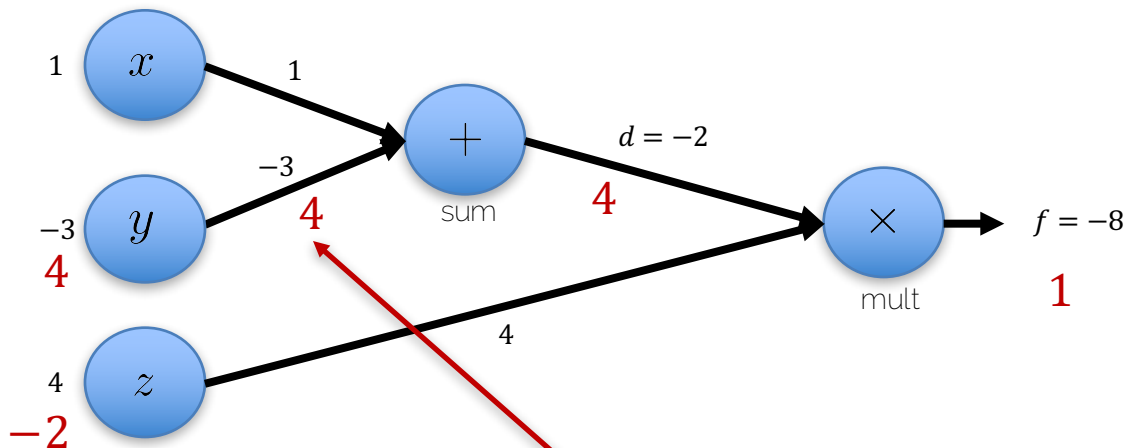
$$f(x, y, z) = (x + y) \cdot z$$

with $x = 1, y = -3, z = 4$

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

$$f = d \cdot z \quad \frac{\partial f}{\partial d} = z, \quad \frac{\partial f}{\partial z} = d$$

What is $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$?



Chain Rule:

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial d} \cdot \frac{\partial d}{\partial y}$$

$$\rightarrow \frac{\partial f}{\partial y} = 4 \cdot 1 = 4$$

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

Backprop: Backward Pass

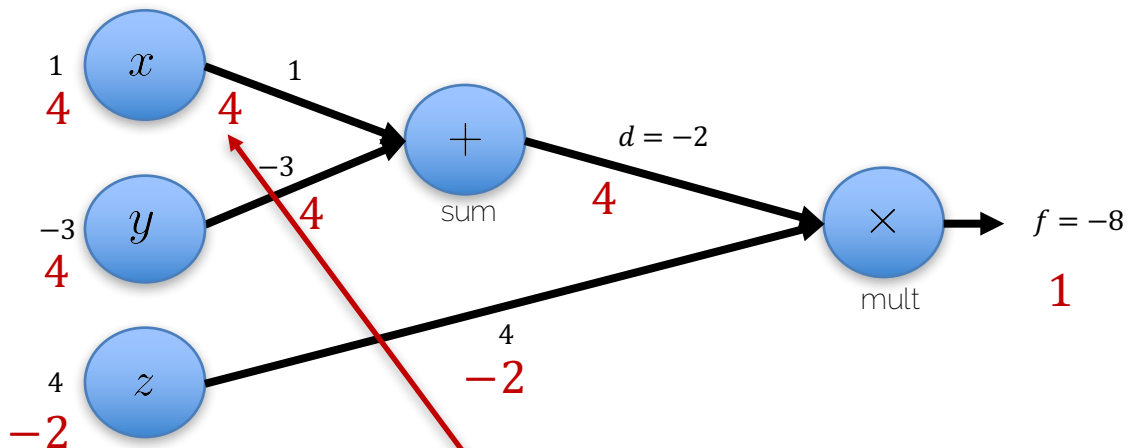
$$f(x, y, z) = (x + y) \cdot z$$

with $x = 1, y = -3, z = 4$

$$d = x + y \quad \boxed{\frac{\partial d}{\partial x} = 1} \quad \frac{\partial d}{\partial y} = 1$$

$$f = d \cdot z \quad \frac{\partial f}{\partial d} = z, \quad \frac{\partial f}{\partial z} = d$$

What is $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$?



Chain Rule:

$$\boxed{\frac{\partial f}{\partial x} = \frac{\partial f}{\partial d} \cdot \frac{\partial d}{\partial x}}$$

$$\rightarrow \frac{\partial f}{\partial x} = 4 \cdot 1 = -4$$

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

Backprop: Backward Pass

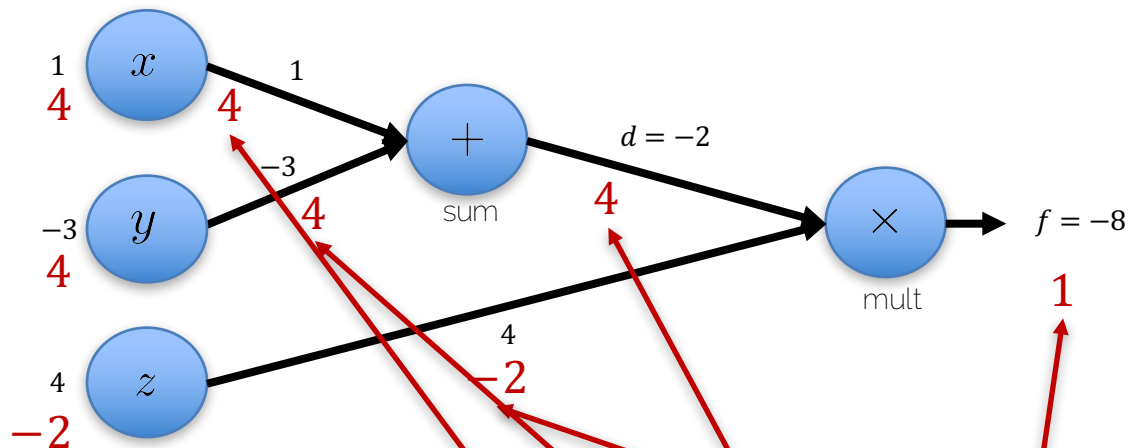
$$f(x, y, z) = (x + y) \cdot z$$

with $x = 1, y = -3, z = 4$

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

$$f = d \cdot z \quad \frac{\partial f}{\partial d} = z, \frac{\partial f}{\partial z} = d$$

What is $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$?



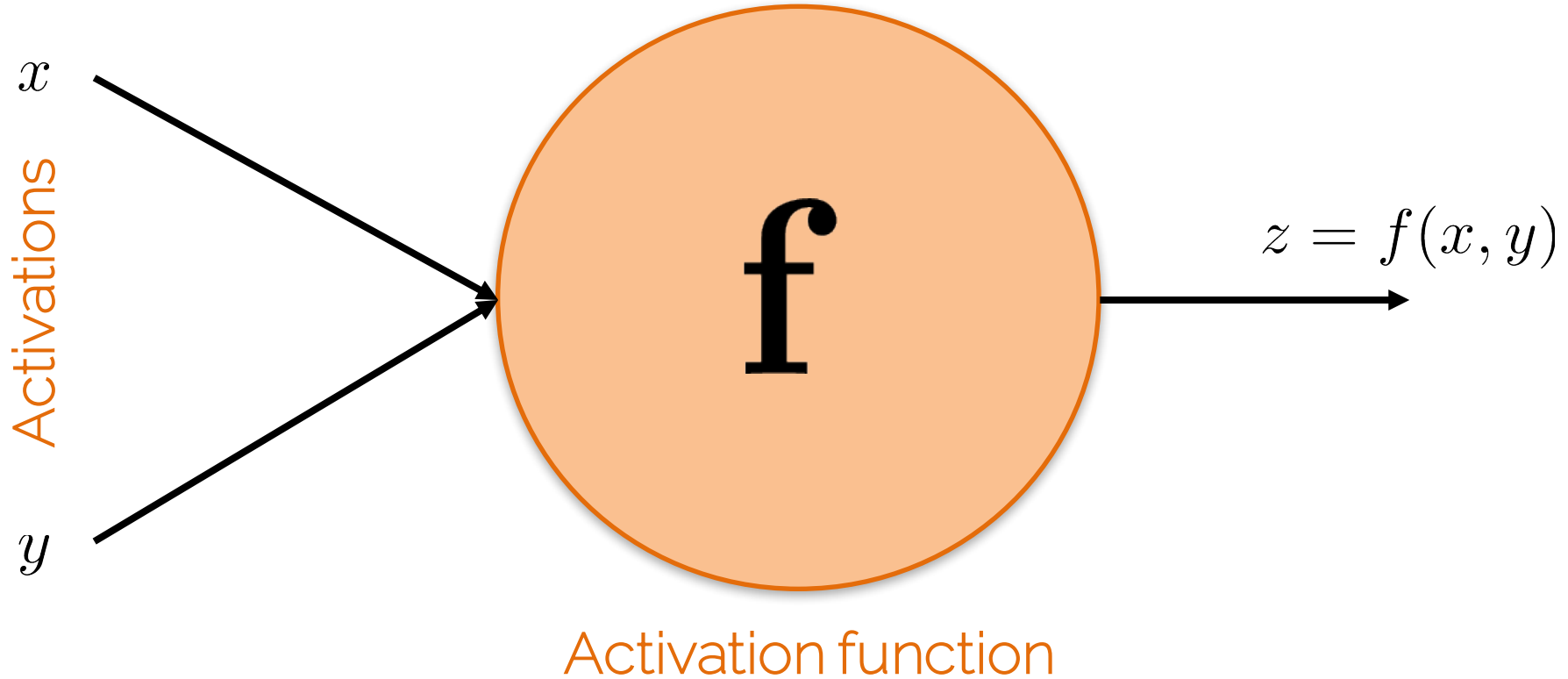
Chain Rule:

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial d} \cdot \frac{\partial d}{\partial y}$$

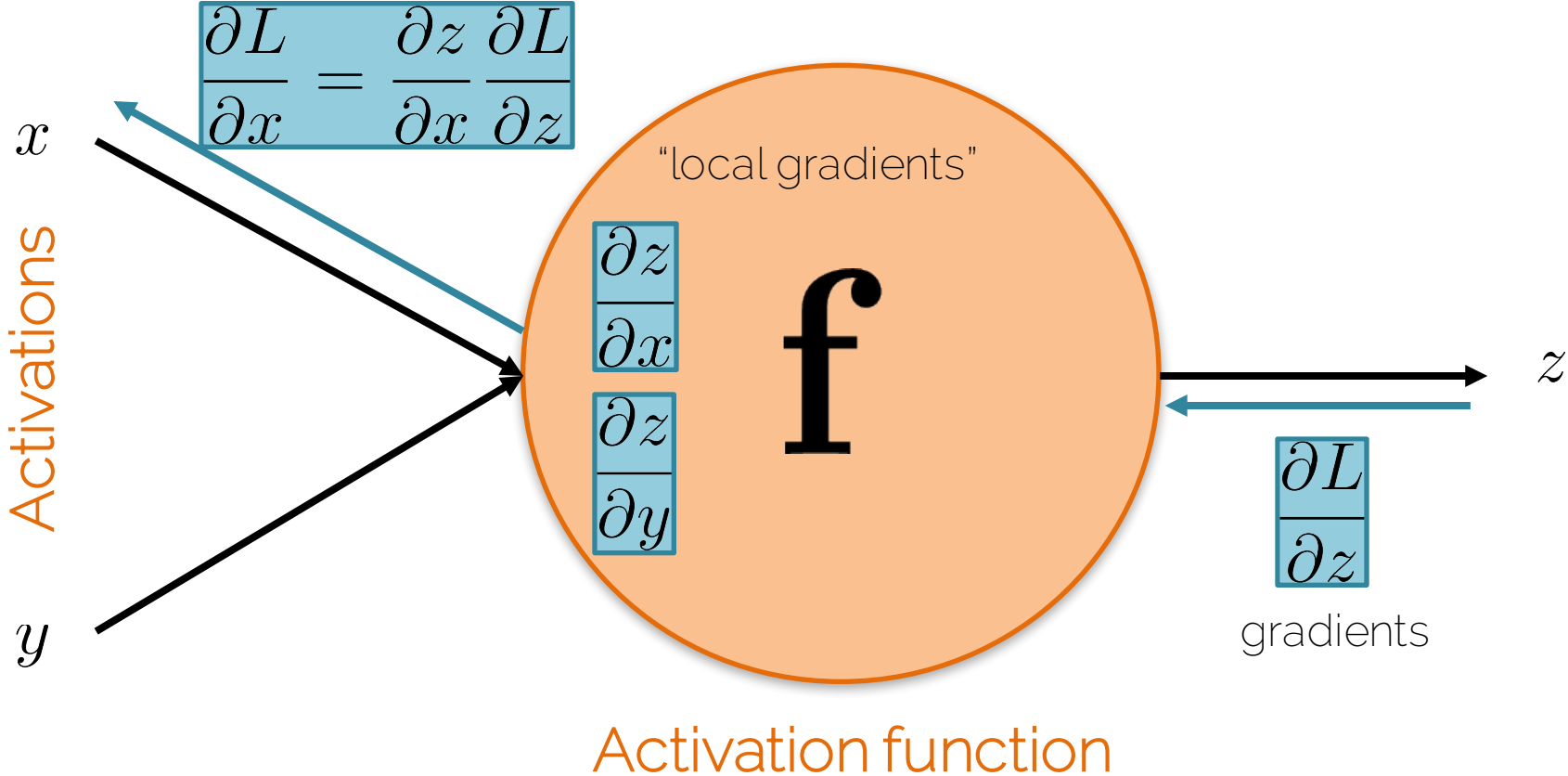
$$\rightarrow \frac{\partial f}{\partial y} = 4 \cdot 1 = 4$$

$$\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial d} \quad \frac{\partial f}{\partial z} \quad \frac{\partial f}{\partial f}$$

The Flow of Gradients

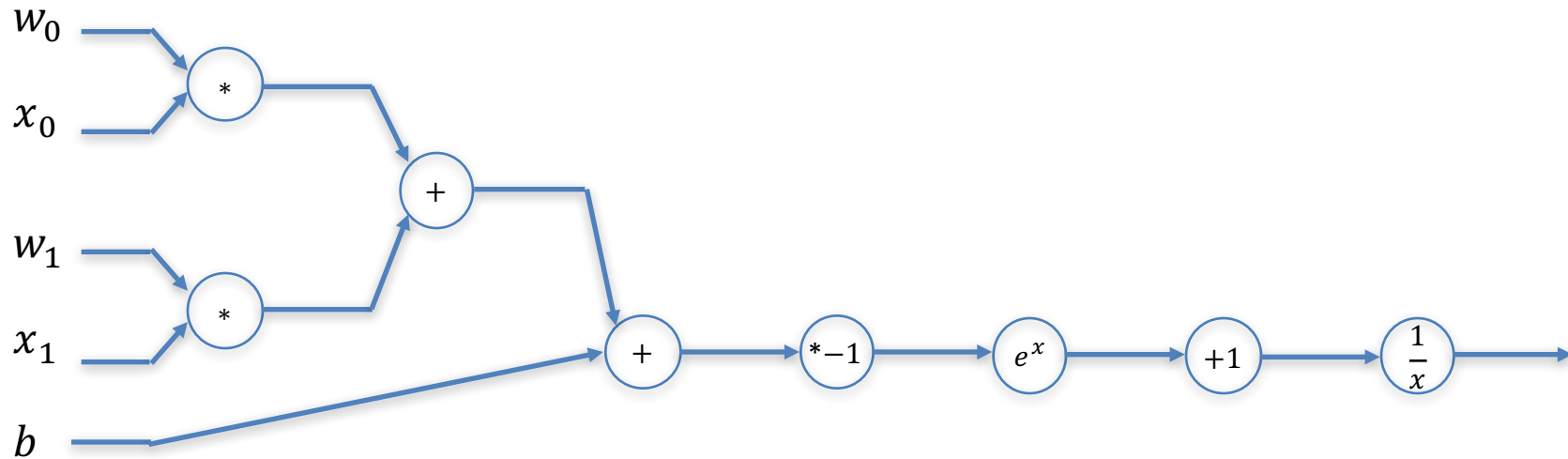


The Flow of Gradients



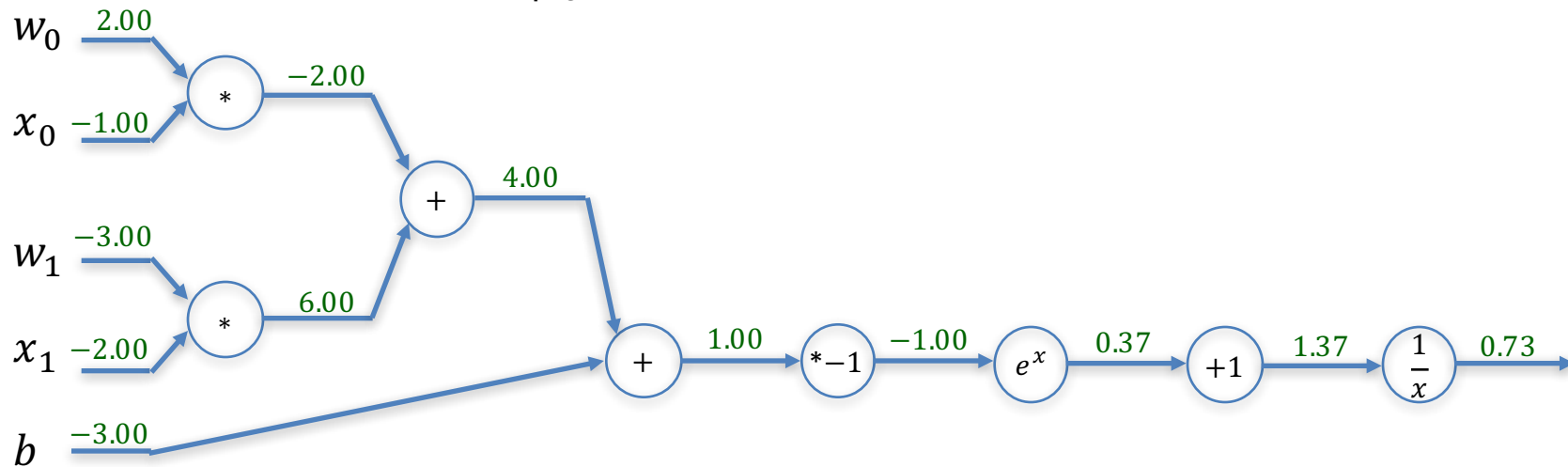
Backprop

$$f(w_0, x_0, w_1, x_1, b) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + b)}}$$



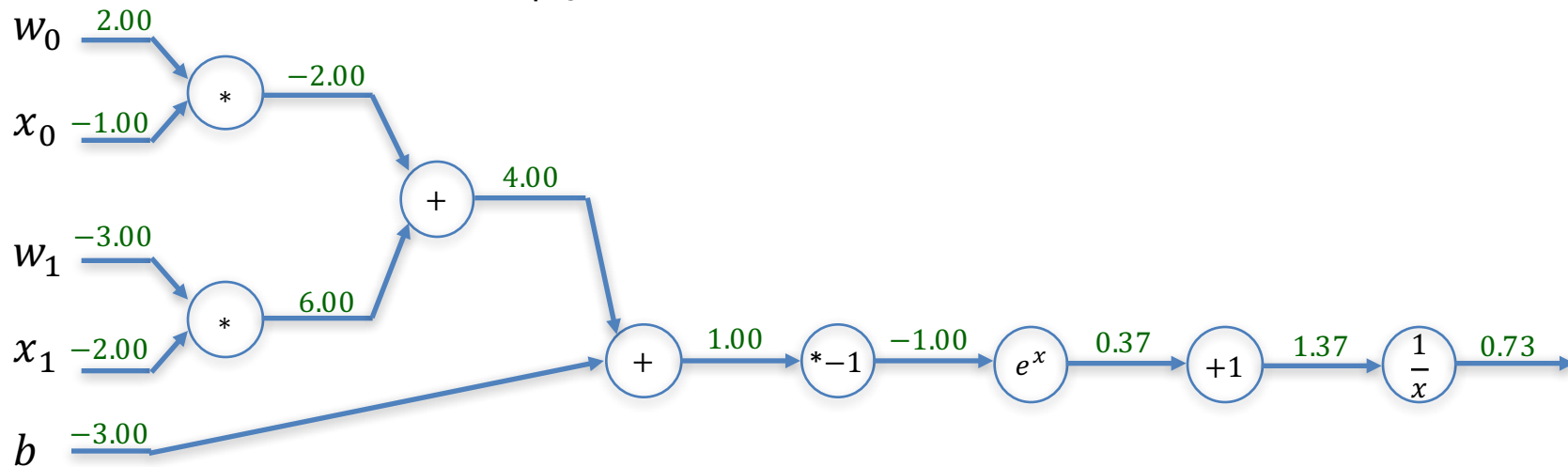
Backprop

$$f(w_0, x_0, w_1, x_1, b) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + b)}}$$



Backprop

$$f(w_0, x_0, w_1, x_1, b) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + b)}}$$



$$f(x) = e^x \quad \rightarrow \quad \frac{\partial f}{\partial x} = e^x$$

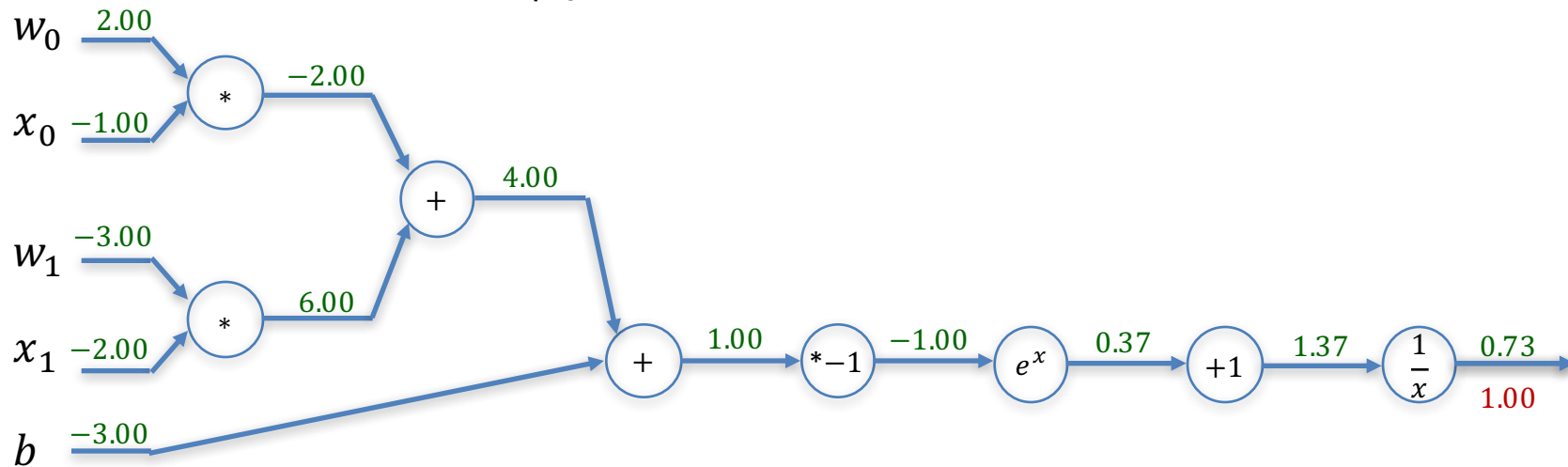
$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_a(x) = ax \quad \rightarrow \quad \frac{\partial f_a}{\partial x} = a$$

$$f_c(x) = c + x \quad \rightarrow \quad \frac{\partial f_c}{\partial x} = 1$$

Backprop

$$f(w_0, x_0, w_1, x_1, b) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + b)}}$$



$$f(x) = e^x \quad \rightarrow \quad \frac{\partial f}{\partial x} = e^x$$

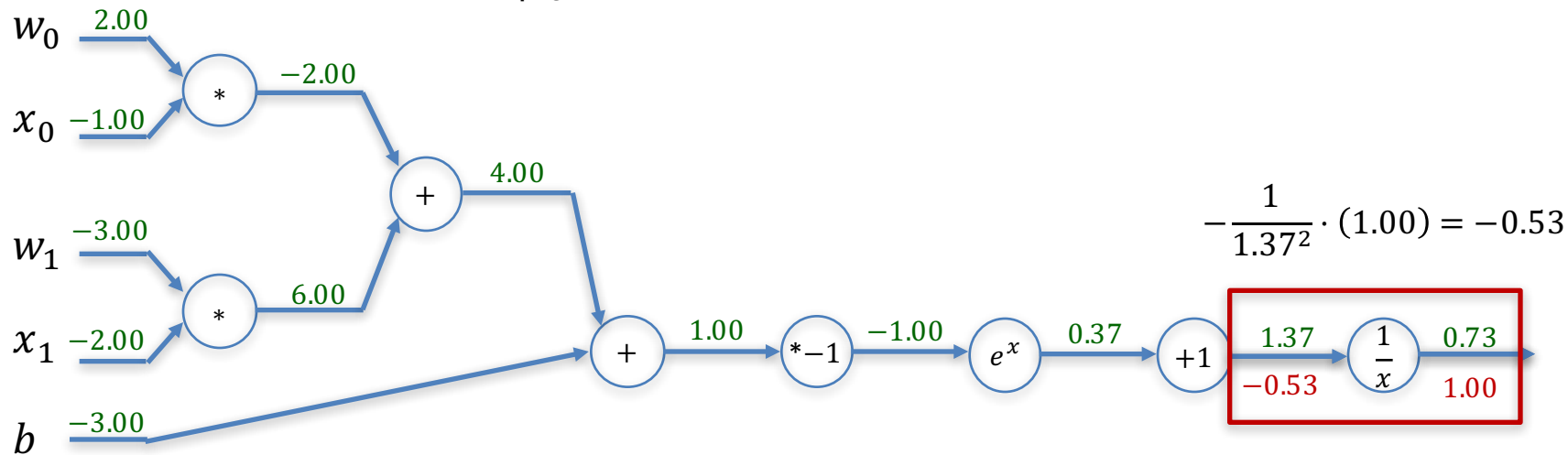
$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_a(x) = ax \quad \rightarrow \quad \frac{\partial f_a}{\partial x} = a$$

$$f_c(x) = c + x \quad \rightarrow \quad \frac{\partial f_c}{\partial x} = 1$$

Backprop

$$f(w_0, x_0, w_1, x_1, b) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + b)}}$$



$$f(x) = e^x \quad \rightarrow \quad \frac{\partial f}{\partial x} = e^x$$

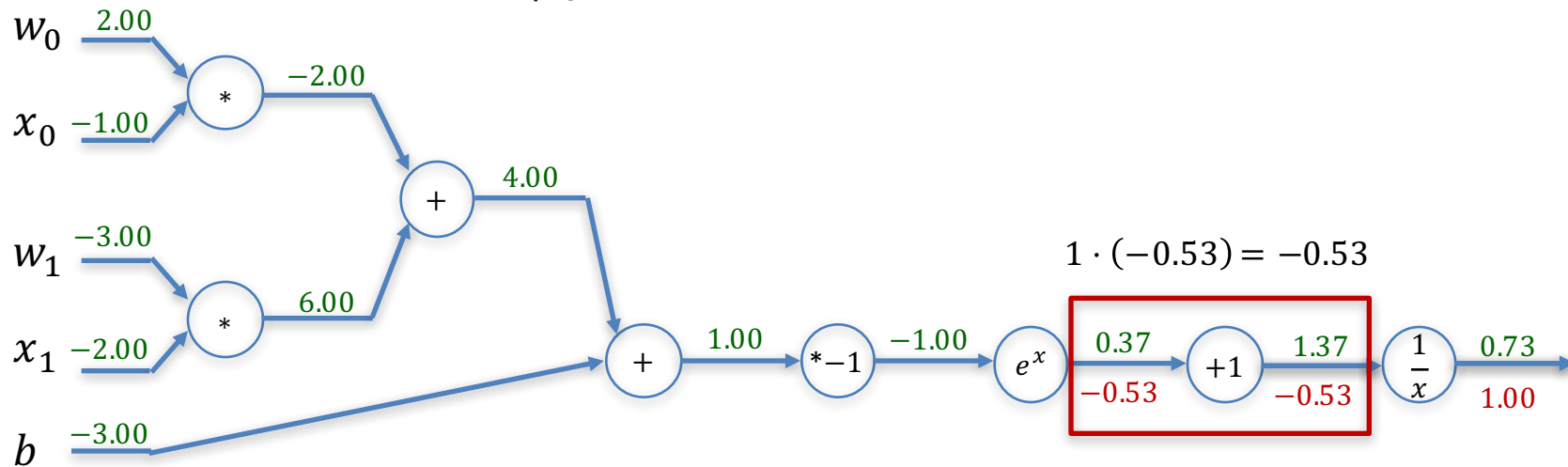
$$f_a(x) = ax \quad \rightarrow \quad \frac{\partial f_a}{\partial x} = a$$

$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_c(x) = c + x \quad \rightarrow \quad \frac{\partial f_c}{\partial x} = 1$$

Backprop

$$f(w_0, x_0, w_1, x_1, b) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + b)}}$$



$$1 \cdot (-0.53) = -0.53$$

$$f(x) = e^x \quad \rightarrow \quad \frac{\partial f}{\partial x} = e^x$$

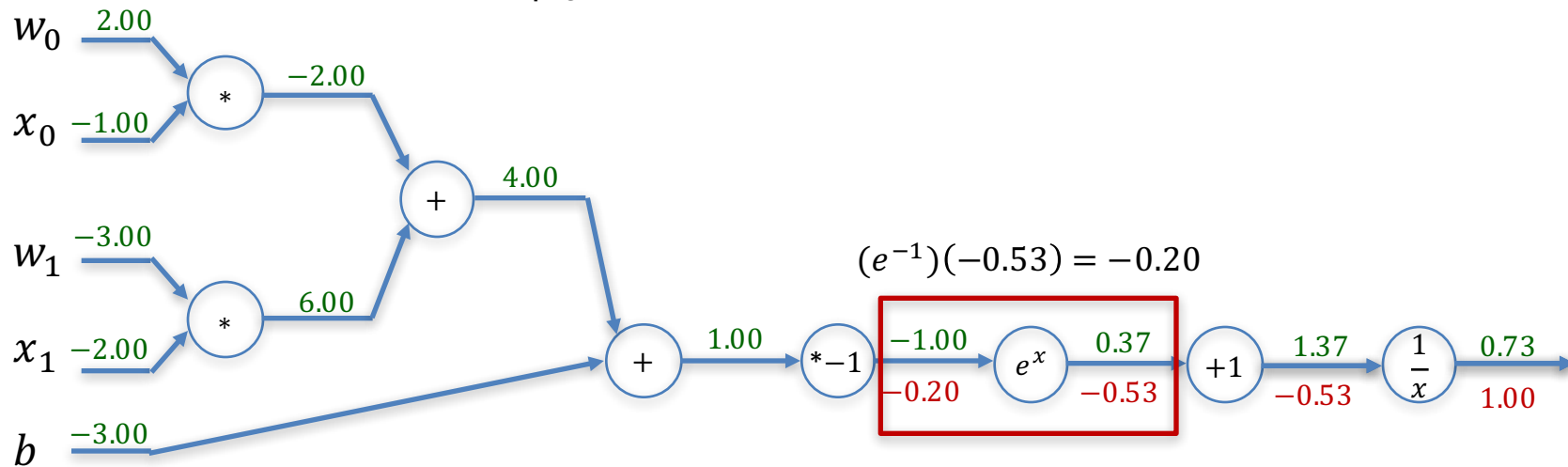
$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_a(x) = ax \quad \rightarrow \quad \frac{\partial f_a}{\partial x} = a$$

$$f_c(x) = c + x \quad \rightarrow \quad \frac{\partial f_c}{\partial x} = 1$$

Backprop

$$f(w_0, x_0, w_1, x_1, b) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + b)}}$$



$$f(x) = e^x \quad \rightarrow \quad \frac{\partial f}{\partial x} = e^x$$

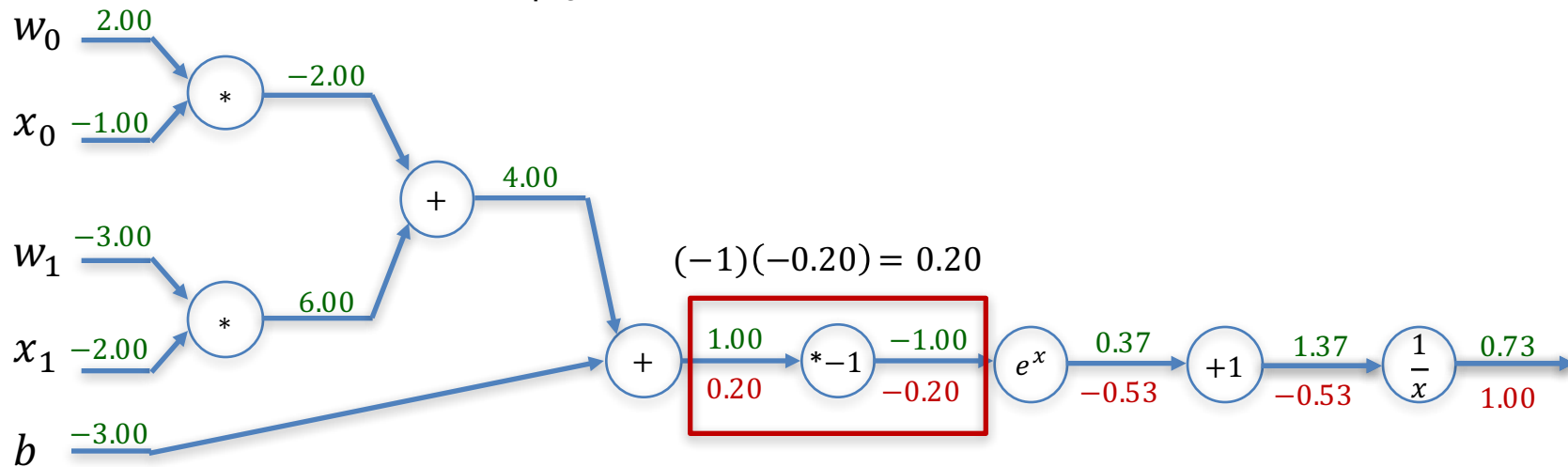
$$f_a(x) = ax \quad \rightarrow \quad \frac{\partial f_a}{\partial x} = a$$

$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_c(x) = c + x \quad \rightarrow \quad \frac{\partial f_c}{\partial x} = 1$$

Backprop

$$f(w_0, x_0, w_1, x_1, b) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + b)}}$$



$$f(x) = e^x \quad \rightarrow \quad \frac{\partial f}{\partial x} = e^x$$

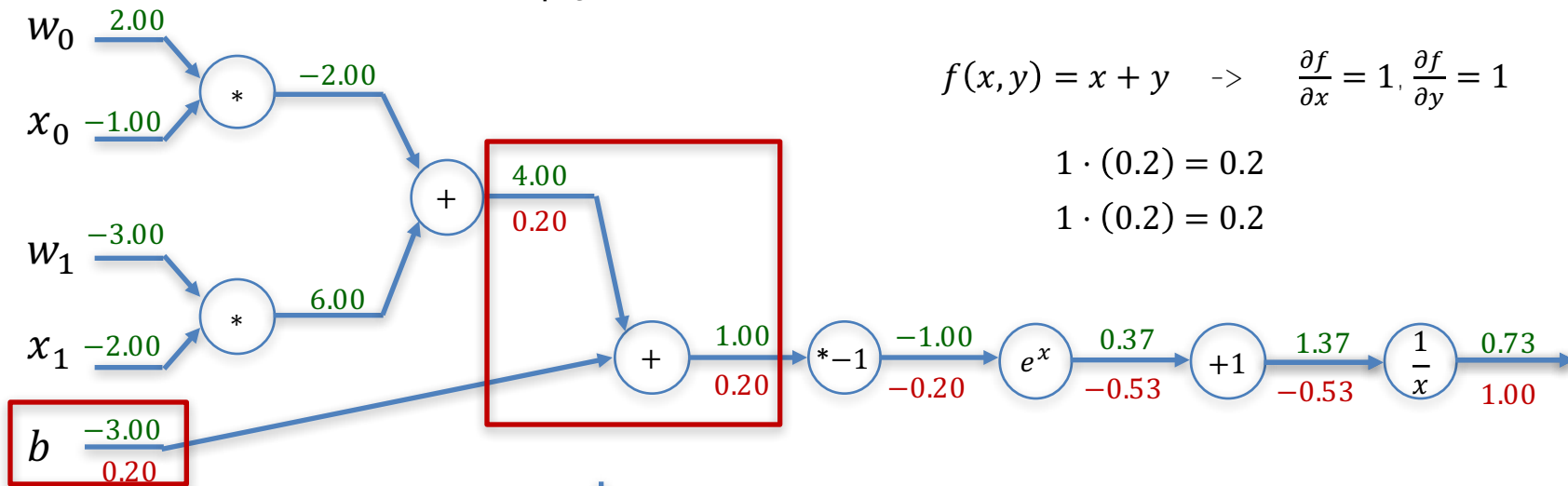
$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_a(x) = ax \quad \rightarrow \quad \frac{\partial f_a}{\partial x} = a$$

$$f_c(x) = c + x \quad \rightarrow \quad \frac{\partial f_c}{\partial x} = 1$$

Backprop

$$f(w_0, x_0, w_1, x_1, b) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + b)}}$$



$$f(x, y) = x + y \quad \rightarrow \quad \frac{\partial f}{\partial x} = 1, \quad \frac{\partial f}{\partial y} = 1$$

$$1 \cdot (0.2) = 0.2$$

$$1 \cdot (0.2) = 0.2$$

$$f(x) = e^x \quad \rightarrow \quad \frac{\partial f}{\partial x} = e^x$$

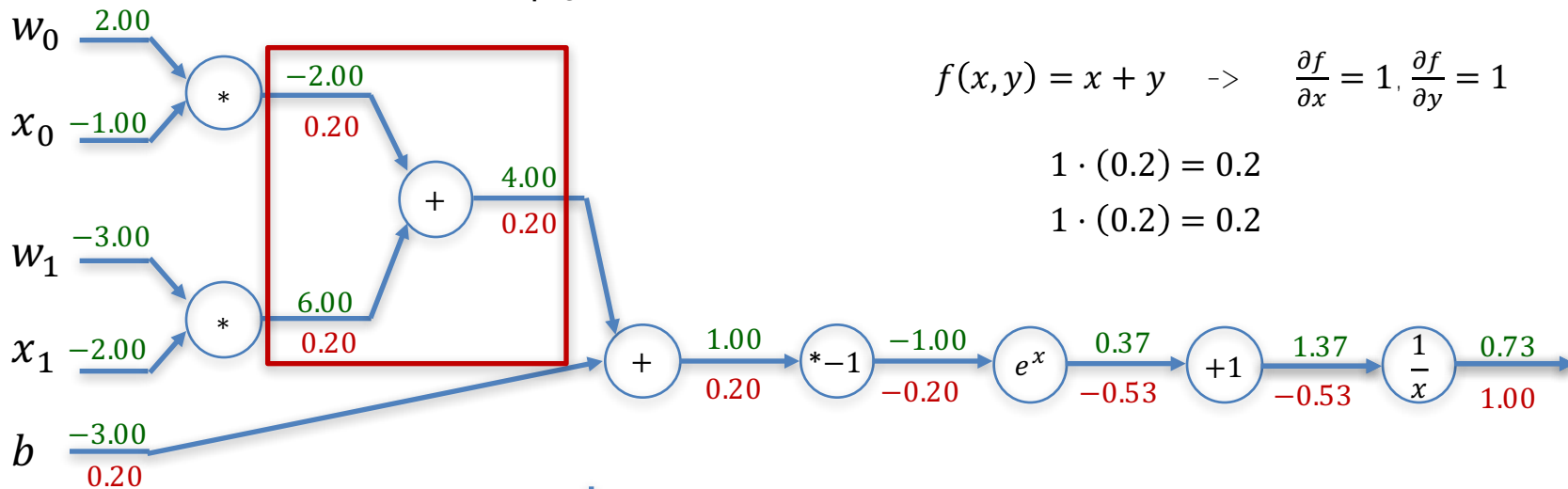
$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_a(x) = ax \quad \rightarrow \quad \frac{\partial f_a}{\partial a} = x$$

$$f_c(x) = c + x \quad \rightarrow \quad \frac{\partial f_c}{\partial c} = 1$$

Backprop

$$f(w_0, x_0, w_1, x_1, b) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + b)}}$$



$$f(x, y) = x + y \quad \rightarrow \quad \frac{\partial f}{\partial x} = 1, \quad \frac{\partial f}{\partial y} = 1$$

$$1 \cdot (0.2) = 0.2$$

$$1 \cdot (0.2) = 0.2$$

$$f(x) = e^x \quad \rightarrow \quad \frac{\partial f}{\partial x} = e^x$$

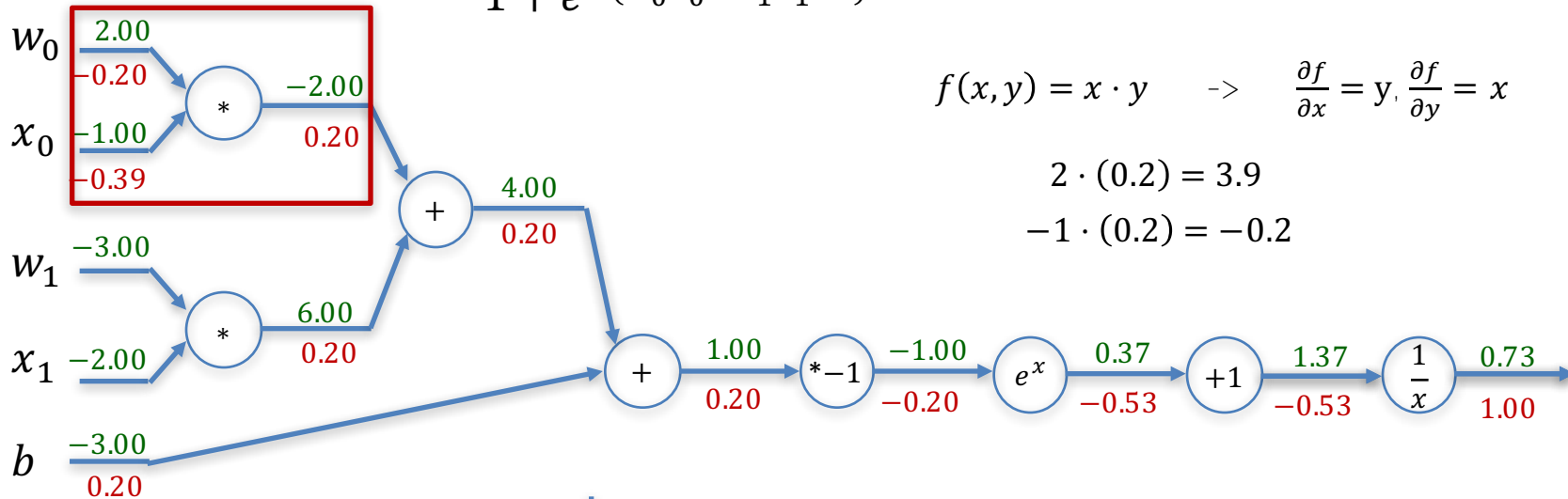
$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_a(x) = ax \quad \rightarrow \quad \frac{\partial f_a}{\partial x} = a$$

$$f_c(x) = c + x \quad \rightarrow \quad \frac{\partial f_c}{\partial x} = 1$$

Backprop

$$f(w_0, x_0, w_1, x_1, b) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + b)}}$$



$$f(x, y) = x \cdot y \quad \rightarrow \quad \frac{\partial f}{\partial x} = y, \quad \frac{\partial f}{\partial y} = x$$

$$2 \cdot (0.2) = 3.9$$

$$-1 \cdot (0.2) = -0.2$$

$$f(x) = e^x \quad \rightarrow \quad \frac{\partial f}{\partial x} = e^x$$

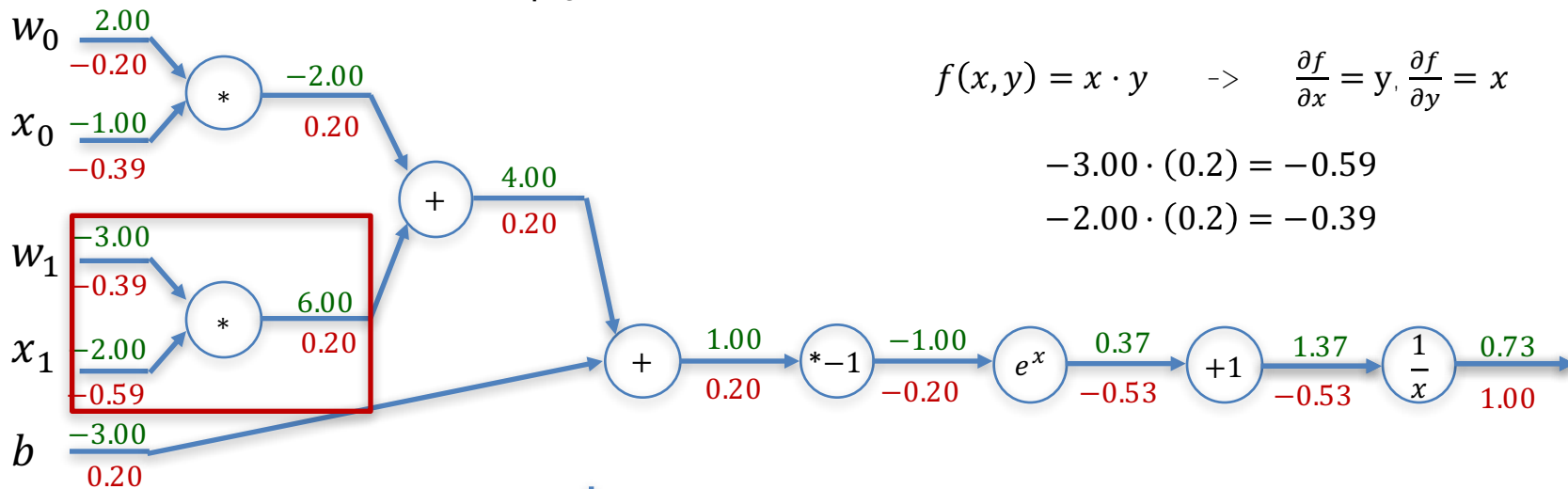
$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_a(x) = ax \quad \rightarrow \quad \frac{\partial f_a}{\partial x} = a$$

$$f_c(x) = c + x \quad \rightarrow \quad \frac{\partial f_c}{\partial x} = 1$$

Backprop

$$f(w_0, x_0, w_1, x_1, b) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + b)}}$$



$$f(x, y) = x \cdot y \quad \rightarrow \quad \frac{\partial f}{\partial x} = y, \quad \frac{\partial f}{\partial y} = x$$

$$-3.00 \cdot (0.2) = -0.59$$

$$-2.00 \cdot (0.2) = -0.39$$

$$f(x) = e^x \quad \rightarrow \quad \frac{\partial f}{\partial x} = e^x$$

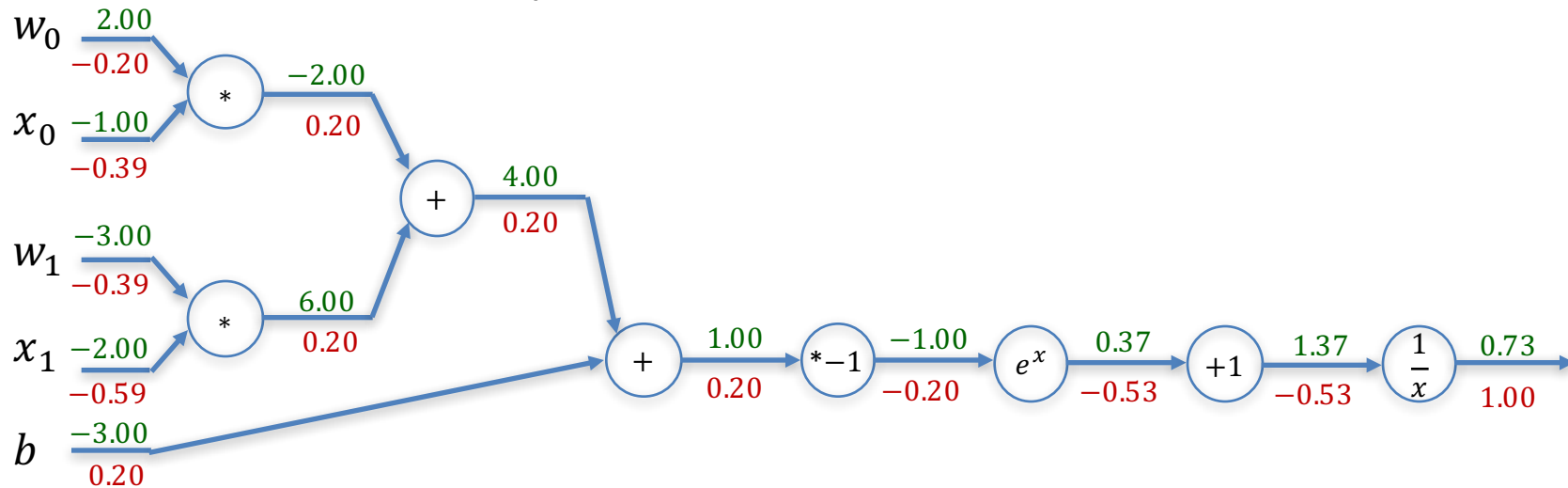
$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_a(x) = ax \quad \rightarrow \quad \frac{\partial f_a}{\partial x} = a$$

$$f_c(x) = c + x \quad \rightarrow \quad \frac{\partial f_c}{\partial x} = 1$$

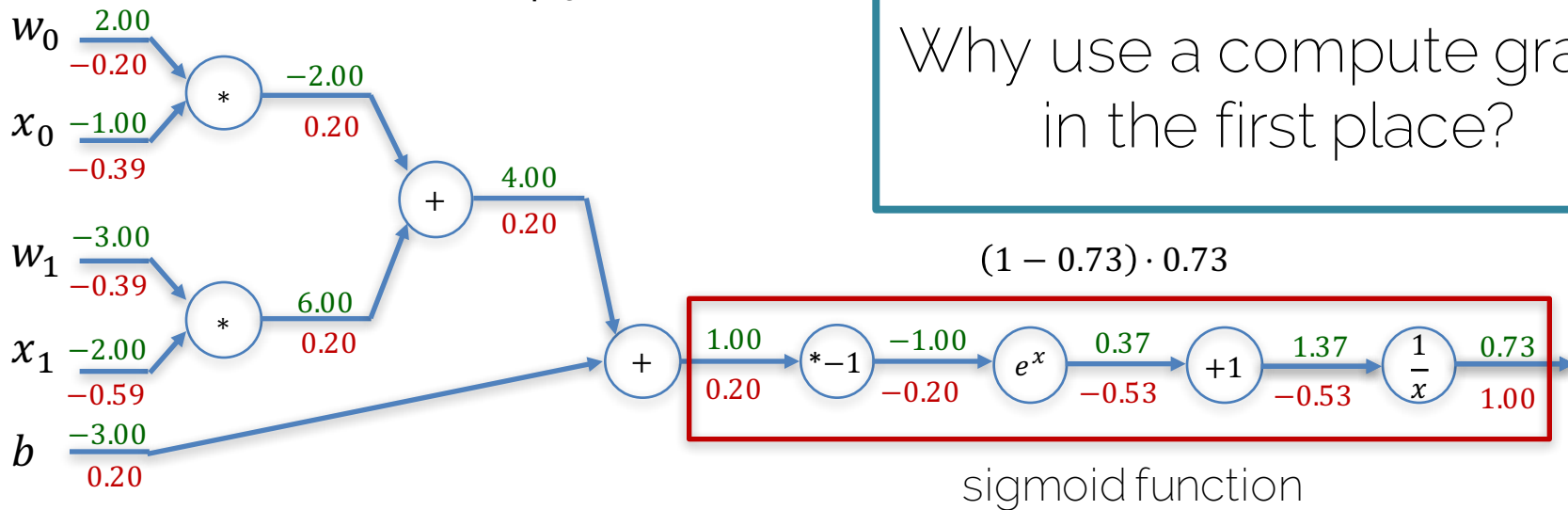
Backprop

$$f(w_0, x_0, w_1, x_1, b) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + b)}}$$



Backprop

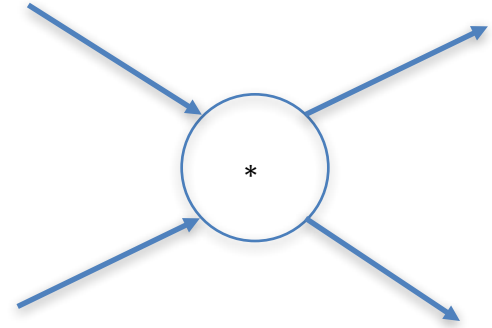
$$f(w_0, x_0, w_1, x_1, b) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + b)}}$$



$$\sigma(x) = \frac{1}{1+e^{-x}} \rightarrow \frac{\partial \sigma(x)}{\partial x} = \frac{e^{-x}}{(1+e^{-x})^2} = (1 - \sigma(x))\sigma(x)$$

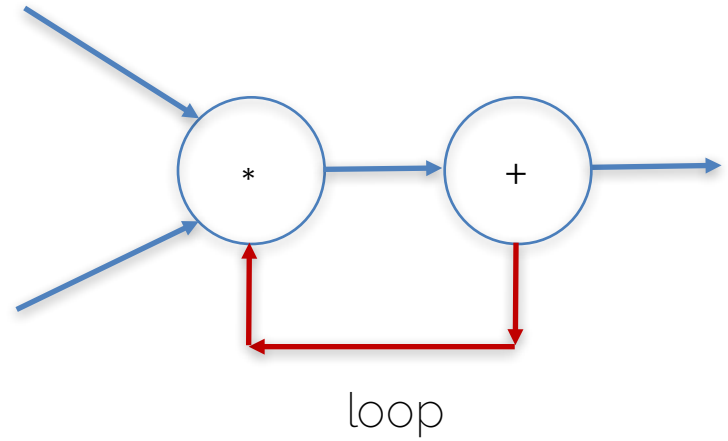
Backpropagation

What happens if there are multiple outputs in a compute node?

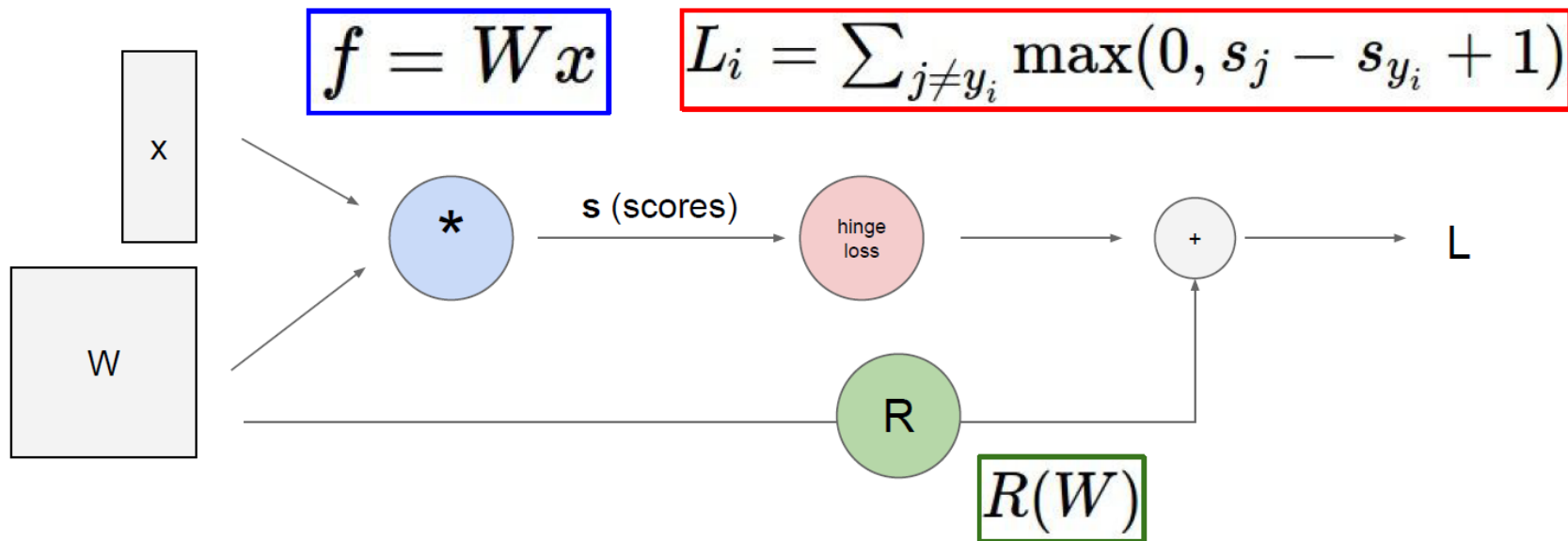


Backpropagation

What happens if there are loops in the graph?

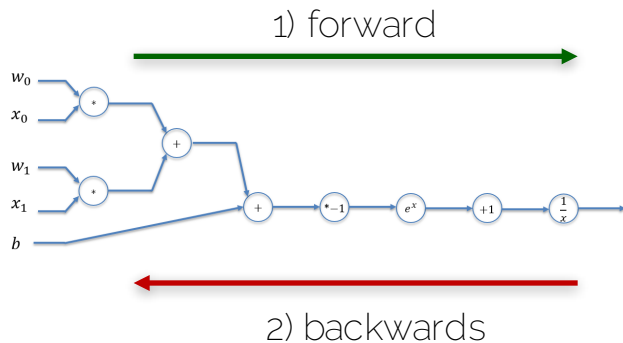


Computational Graph



Combining nodes:
Linear activation node + hinge loss + regularization

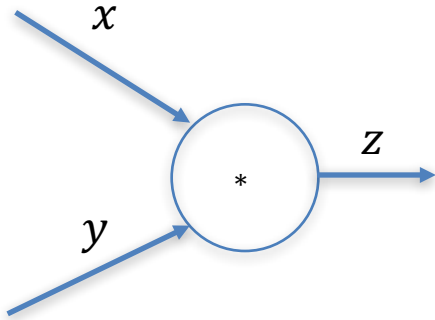
Implementation of Compute Graph



```
class ComputationalGraph(object):
    #...
    def forward(inputs):
        #1. [pass inputs to input nodes]
        #2. forward traverse the computational graph
        for node in self.graph.nodes_topologically_sorted():
            node.forward()
            # forward intermediates / loss
        return loss # final node returns loss
    def backward():
        for node in self.graph.nodes_topologically_sorted_reverse():
            node.backward() #apply chainrule
            # backward intermediate derivatives
        return inputs_gradients
```

Implementation of Nodes

- Forward and backward pass of MulNode



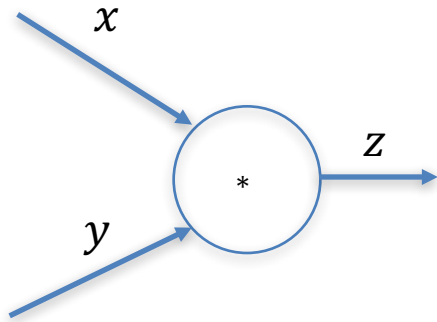
all values are scalars

```
class MulNode(object):
    def forward(x,y):
        z=x*y
        return z
    def backward(dz, x, y):
        dx = y*dz # [dz/dx * dL/dz]
        dy = x*dz # [dz/dy * dL/dz]
        return [dx, dy]
```

Issue?

Implementation of Nodes

- Forward and backward pass of MulNode



all values are scalars

```
class MulNode(object):
    def forward(x,y):
        z = x*y
        self.x = x
        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]
```

Cache results of forward pass
-> faster runtime for backward pass

Torch: Layers (GitHub)

LayerNormalization.lua	Adding Layer Normalization	2 months ago
LeakyReLU.lua	Add THNN conversion of {ELU, LeakyReLU, LogSigmoid, LogSoftMax, Looku...	a year ago
Linear.lua	Fix shared function override for specific modules	4 months ago
Log.lua	add nn.Inc & nn.Scale	a year ago
LogSigmoid.lua	lazy init	a year ago
LogSoftMax.lua	Add THNN conversion of {ELU, LeakyReLU, LogSigmoid, LogSoftMax, Looku...	a year ago
LookupTable.lua	Fix shared function override for specific modules	4 months ago
MM.lua		
MSECriterion.lua		
MV.lua		
MapTable.lua		
MarginCriterion.lua		
MarginRankingCriterion.lua		
MaskedSelect.lua		
Max.lua		
Maxout.lua		
Mean.lua		
Min.lua		
MixtureTable.lua		
Module.lua		
Mul.lua		
MulConstant.lua		
MultiCriterion.lua		
Reshape.lua	Better __tostring__ and cleans formatting	a year ago
Select.lua	Adds negative dim arguments	11 months ago
SelectTable.lua	allow SelectTable to accept input that contains tables of things that...	2 months ago
Sequential.lua	Improve error handling	a year ago
Sigmoid.lua	Add THNN conversion of {RReLU, Sigmoid, SmoothL1Criterion,SoftMax, So...	a year ago
SmoothL1Criterion.lua	Add THNN conversion of {RReLU, Sigmoid, SmoothL1Criterion,SoftMax, So...	a year ago
SoftMarginCriterion.lua	SoftMarginCriterion	a year ago
SoftMax.lua	Add THNN conversion of {RReLU, Sigmoid, SmoothL1Criterion,SoftMax, So...	a year ago
SoftMin.lua	nn.clearState	a year ago
SoftPlus.lua	Add THNN conversion of {RReLU, Sigmoid, SmoothL1Criterion,SoftMax, So...	a year ago
SoftShrink.lua	Add THNN conversion of {oftShrink, Sqrt, Square, Tanh, Threshold}	a year ago
SoftSign.lua	nn.clearState	a year ago
SparseJacobian.lua	Fix various unused variables in nn	3 years ago
SparseLinear.lua	Fixing sparse linear race condition	a year ago
SpatialAdaptiveAveragePooling.lua	Add SpatialAdaptiveAveragePooling.	4 months ago
SpatialAdaptiveMaxPooling.lua	Indices for nn.	7 months ago
SpatialAutoCropMSECriterion.lua	fix local / global var leaks	4 months ago

Torch: MulConstant

```
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
```

Init()

```
14
15 function MulConstant:updateOutput(input)
16     if self.inplace then
17         input:mul(self.constant_scalar)
18         self.output:set(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()

$$f(x) = aX$$

```
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:set(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()

Caffee: Layers (GitHub)

absval_layer.cpp	dismantle layer headers	2 years ago
absval_layer.cu	dismantle layer headers	2 years ago
accuracy_layer.cpp		
argmax_layer.cpp		
base_conv_layer.cpp		
base_data_layer.cpp		
base_data_layer.cu		
batch_norm_layer.cpp		
batch_norm_layer.cu		
batch_reindex_layer.cpp		
batch_reindex_layer.cu		
bnll_layer.cpp		
bnll_layer.cu		
concat_layer.cpp		
concat_layer.cu		
contrastive_loss_layer.cpp		
contrastive_loss_layer.cu		
conv_layer.cpp		
conv_layer.cu		
cudnn_conv_layer.cpp		
cudnn_conv_layer.cu		
cudnn_lcn_layer.cpp		
cudnn_lcn_layer.cu		
cudnn_lm_layer.cpp		
cudnn_lm_layer.cu		
cudnn_pooling_layer.cpp		
cudnn_pooling_layer.cu		
cudnn_relu_layer.cpp		
cudnn_relu_layer.cu		
cudnn_sigmoid_layer.cpp		
cudnn_sigmoid_layer.cu		
cudnn_softmax_layer.cpp		
cudnn_softmax_layer.cu		
cudnn_tanh_layer.cpp		
pooling_layer.cpp	dismantle layer headers	2 years ago
pooling_layer.cu	dismantle layer headers	2 years ago
power_layer.cpp	dismantle layer headers	2 years ago
power_layer.cu	dismantle layer headers	2 years ago
prelu_layer.cpp	dismantle layer headers	2 years ago
prelu_layer.cu	dismantle layer headers	2 years ago
reduction_layer.cpp	dismantle layer headers	2 years ago
reduction_layer.cu	dismantle layer headers	2 years ago
relu_layer.cpp	dismantle layer headers	2 years ago
relu_layer.cu	dismantle layer headers	2 years ago
reshape_layer.cpp	dismantle layer headers	2 years ago
sigmoid_cross_entropy_loss_layer.cpp	dismantle layer headers	2 years ago
sigmoid_cross_entropy_loss_layer.cu	dismantle layer headers	2 years ago
sigmoid_layer.cpp	dismantle layer headers	2 years ago
sigmoid_layer.cu	dismantle layer headers	2 years ago
silence_layer.cpp	dismantle layer headers	2 years ago
silence_layer.cu	dismantle layer headers	2 years ago
slice_layer.cpp	dismantle layer headers	2 years ago

Caffe: Sigmoid_Layer

```
#include "caffe/layers/sigmoid_layer.hpp"
```

```
namespace caffe {  
  
template <typename Dtype>  
inline Dtype sigmoid(Dtype x) {  
    return 1. / (1. + exp(-x));  
}
```

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

```
template <typename Dtype>  
void SigmoidLayer<Dtype>::Forward_cpu(const vector<Blob<Dtype>*>& bottom,  
    const vector<Blob<Dtype>*>& top) {  
    const Dtype* bottom_data = bottom[0]->cpu_data();  
    Dtype* top_data = top[0]->mutable_cpu_data();  
    const int count = bottom[0]->count();  
    for (int i = 0; i < count; ++i) {  
        top_data[i] = sigmoid(bottom_data[i]);  
    }  
}
```

Forward()

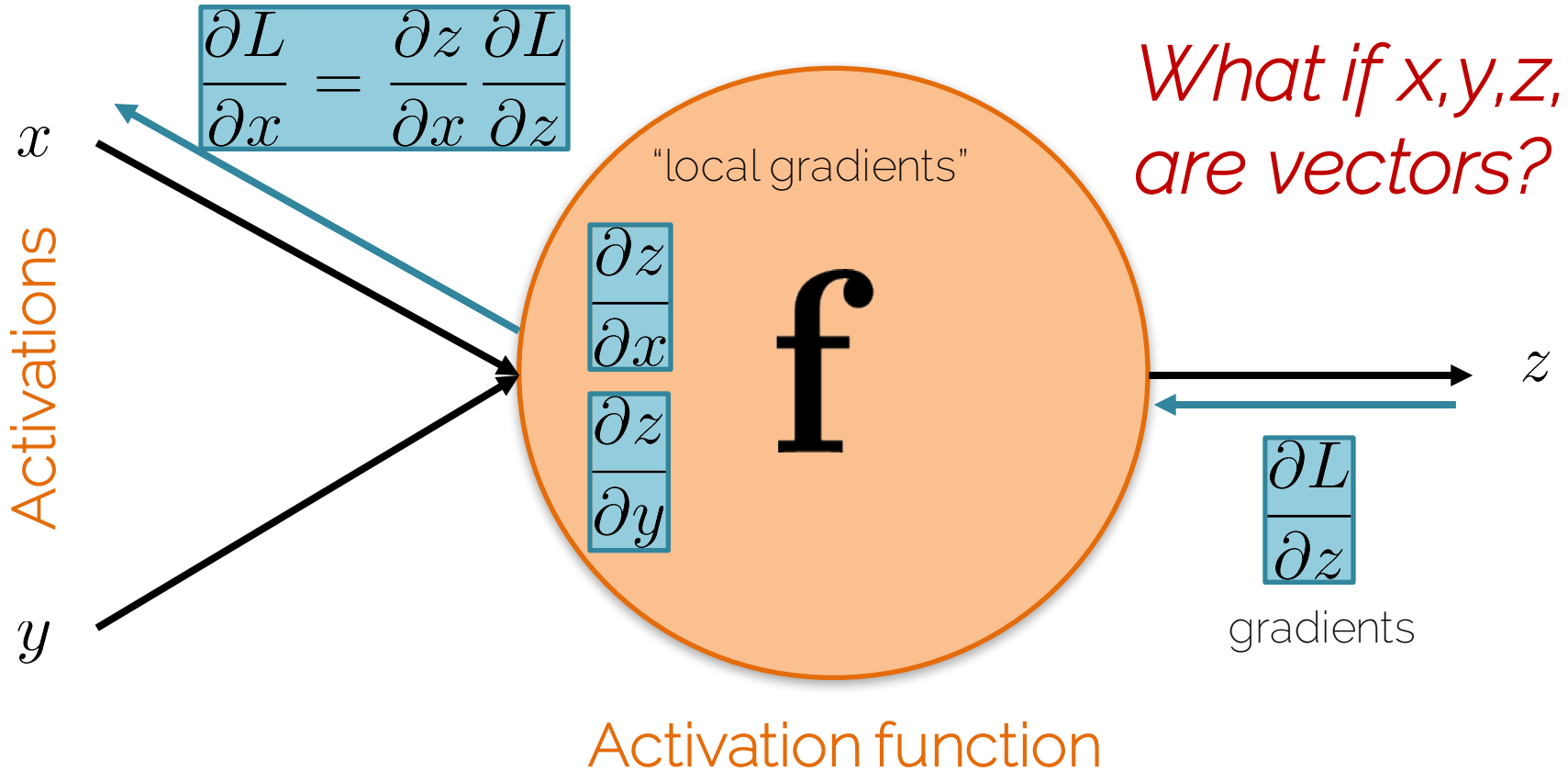
```
template <typename Dtype>  
void SigmoidLayer<Dtype>::Backward_cpu(const vector<Blob<Dtype>*>& top,  
    const vector<bool>& propagate_down,  
    const vector<Blob<Dtype>*>& bottom) {  
    if (propagate_down[0]) {  
        const Dtype* top_data = top[0]->cpu_data();  
        const Dtype* top_diff = top[0]->cpu_diff();  
        Dtype* bottom_diff = bottom[0]->mutable_cpu_diff();  
        const int count = bottom[0]->count();  
        for (int i = 0; i < count; ++i) {  
            const Dtype sigmoid_x = top_data[i];  
            bottom_diff[i] = top_diff[i] * sigmoid_x * (1. - sigmoid_x);  
        }  
    }  
}
```

Backward()

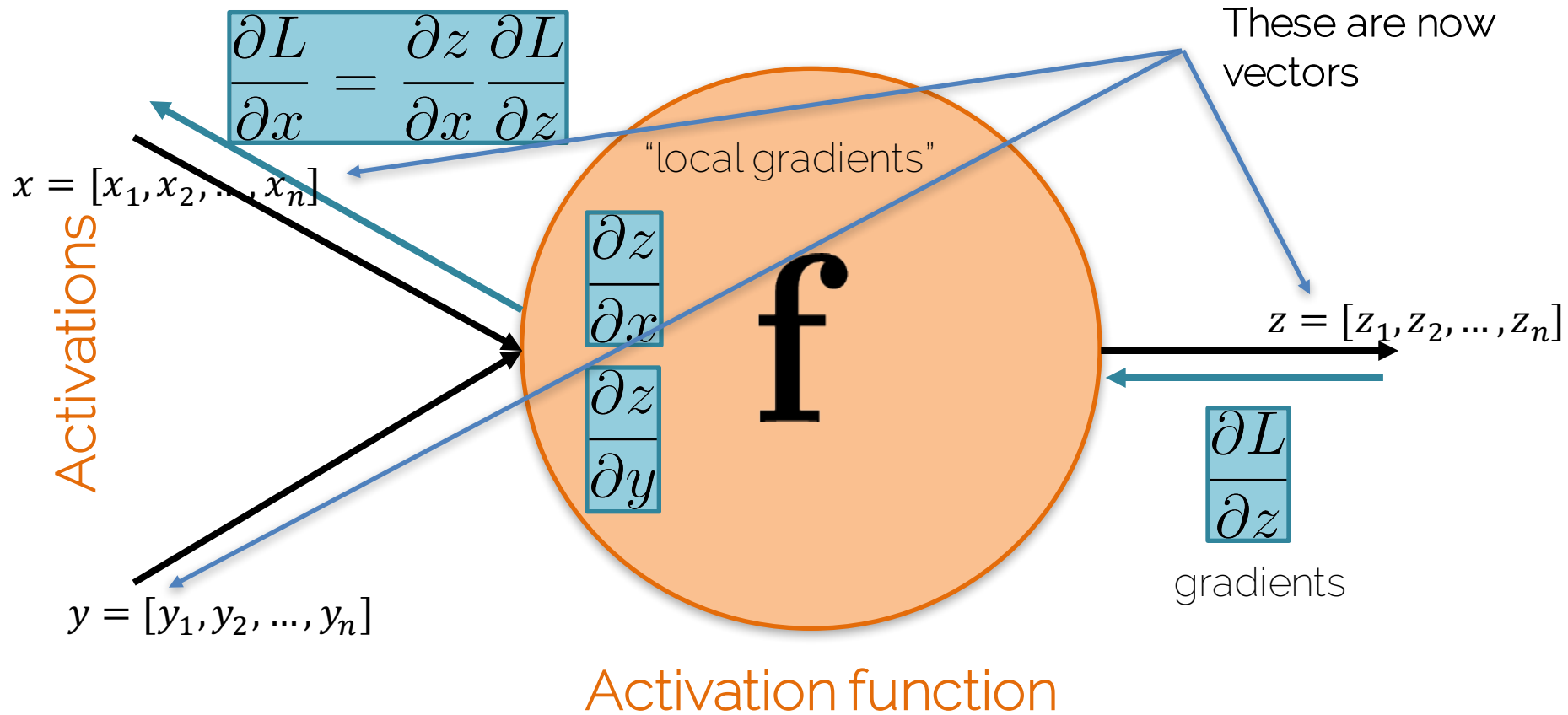
$$\sigma'(x) = (1 - \sigma(x))\sigma(x)$$

```
37 }  
38 }  
39
```

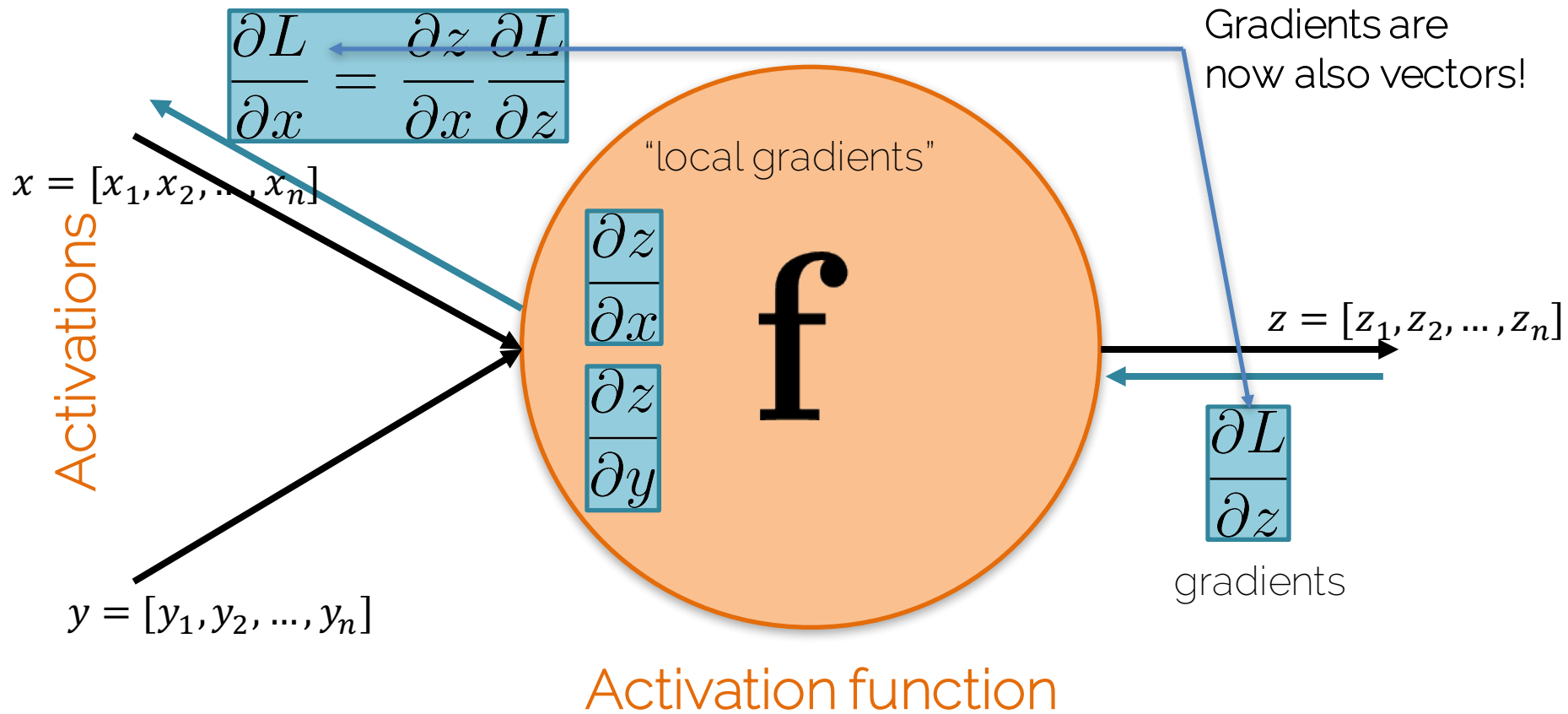
Vectorized Operations



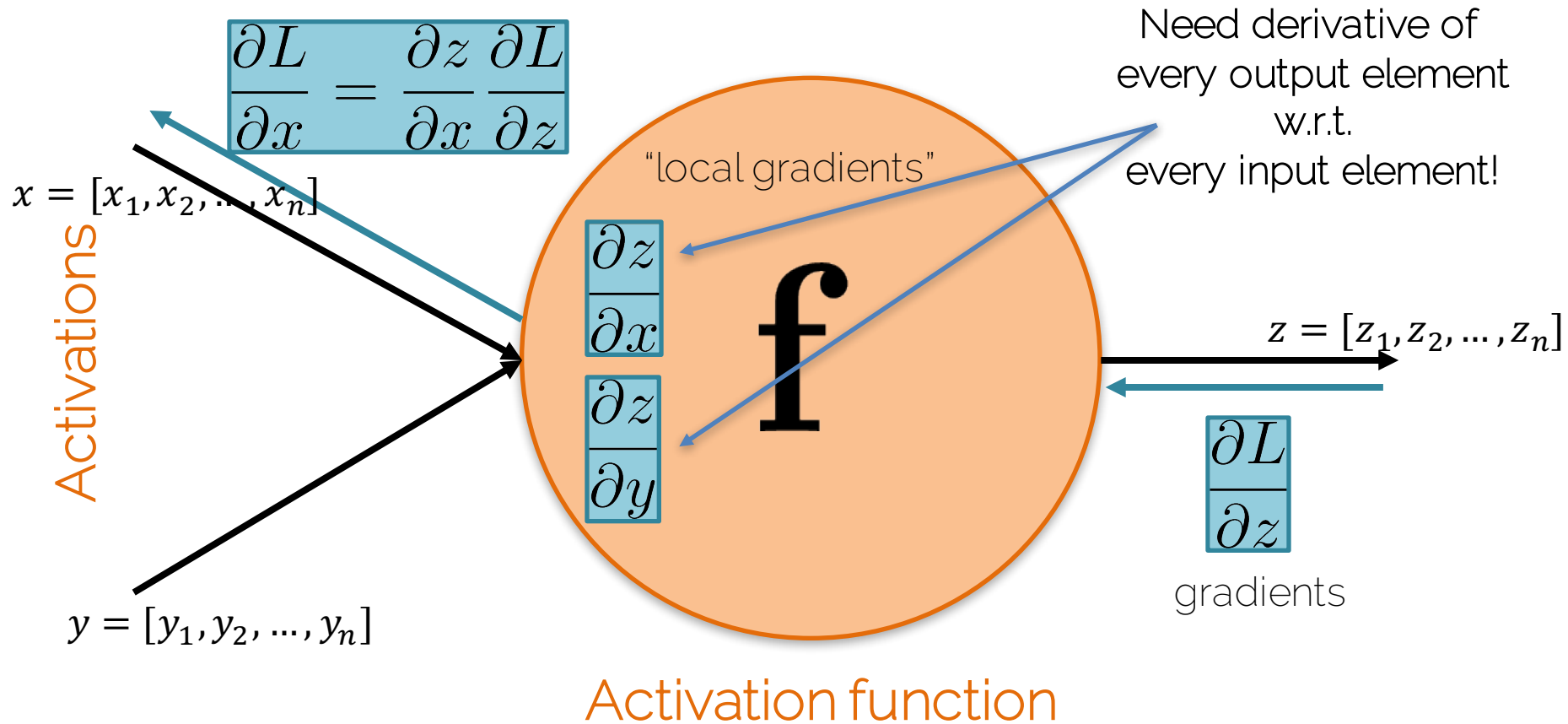
Vectorized Operations



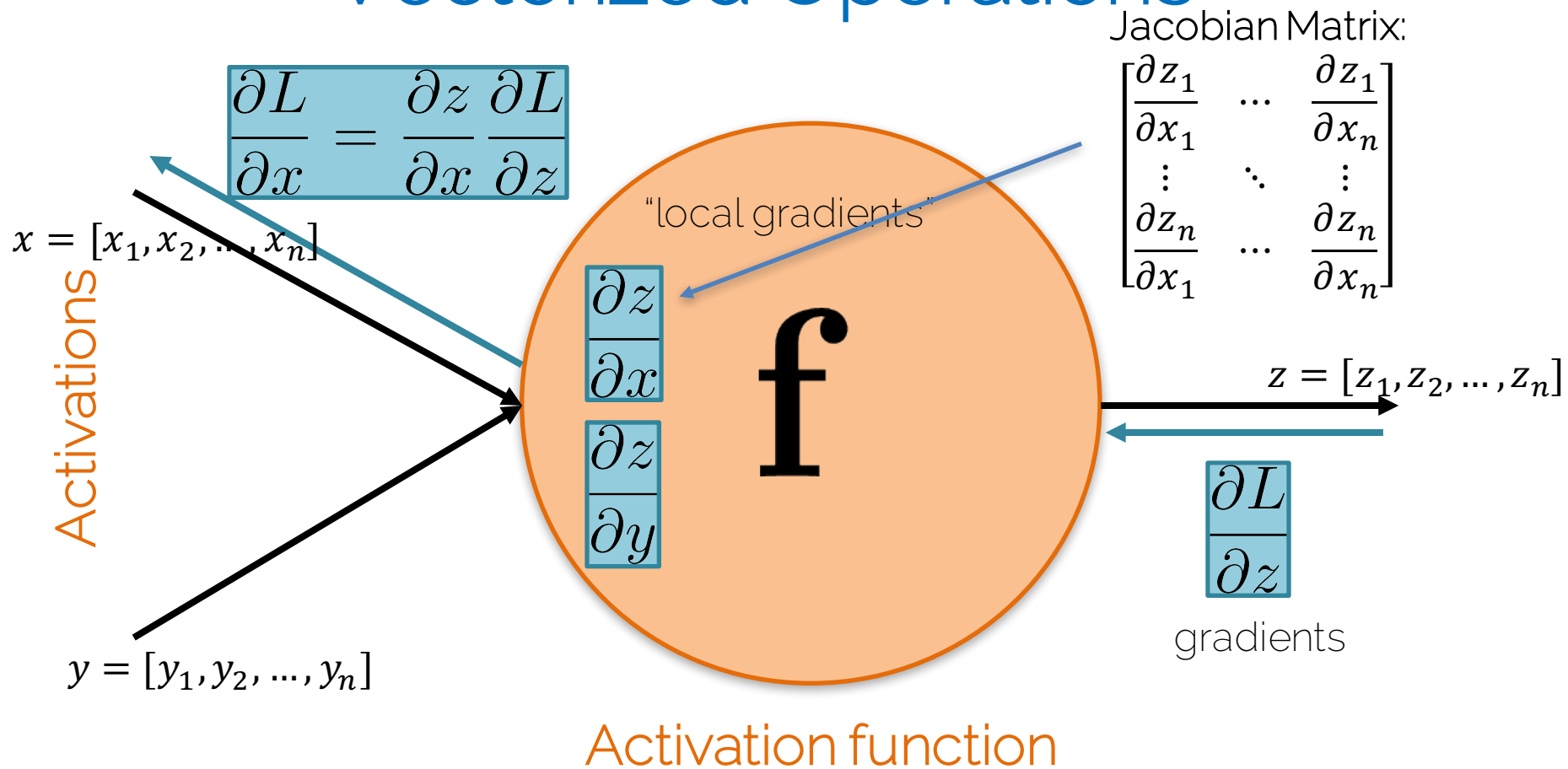
Vectorized Operations



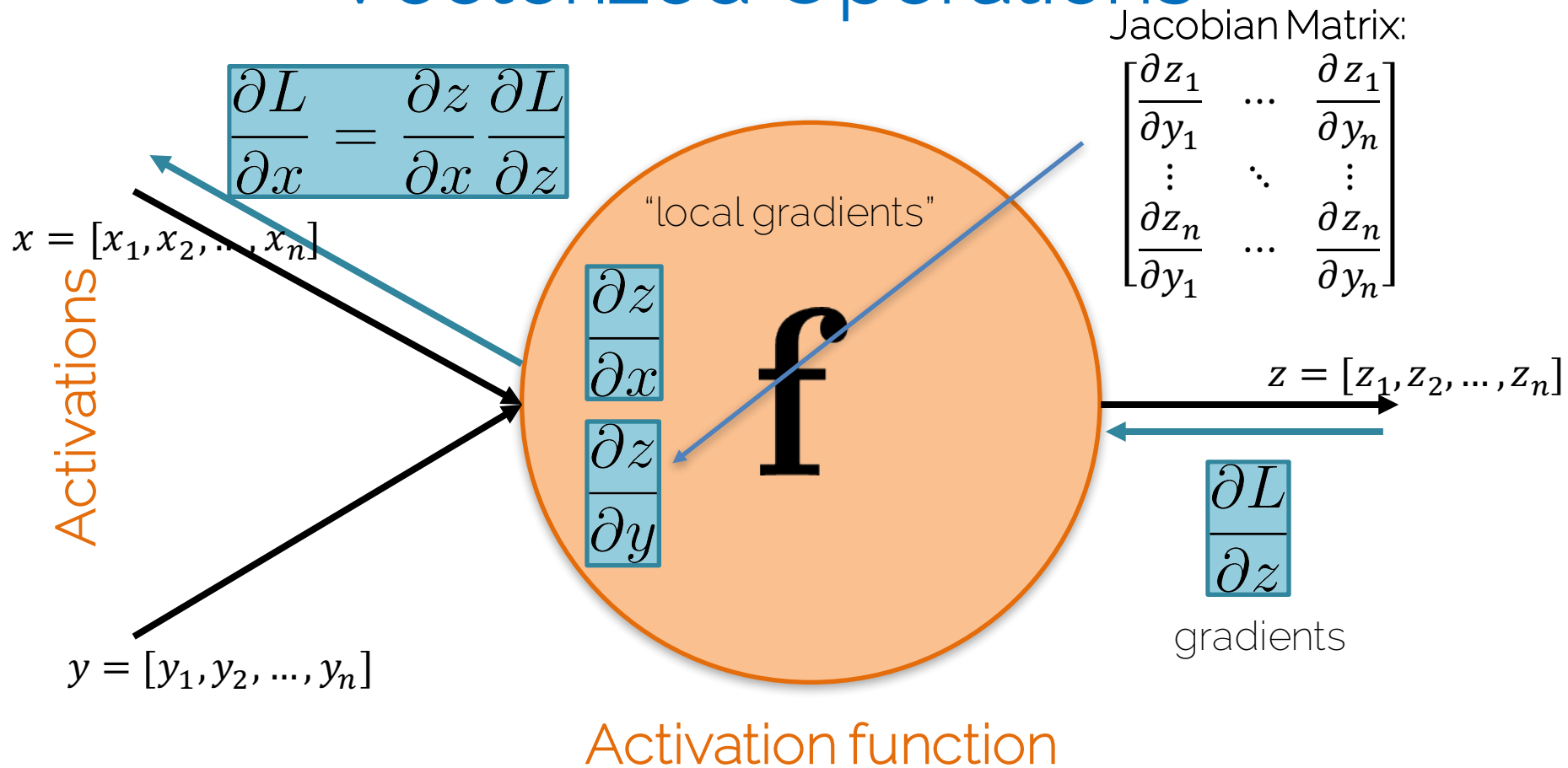
Vectorized Operations



Vectorized Operations



Vectorized Operations



Vectorized Operations

Assuming input and output $\in \mathbb{R}^{4096}$

Jacobian Matrix:

$$\begin{bmatrix} \frac{\partial z_1}{\partial y_1} & \dots & \frac{\partial z_1}{\partial y_n} \\ \frac{\partial z_2}{\partial y_1} & \dots & \frac{\partial z_2}{\partial y_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_n}{\partial y_1} & \dots & \frac{\partial z_n}{\partial y_n} \end{bmatrix}$$

$$x = [x_1, x_2, \dots, x_n]$$

$$x \in \mathbb{R}^{4096}$$

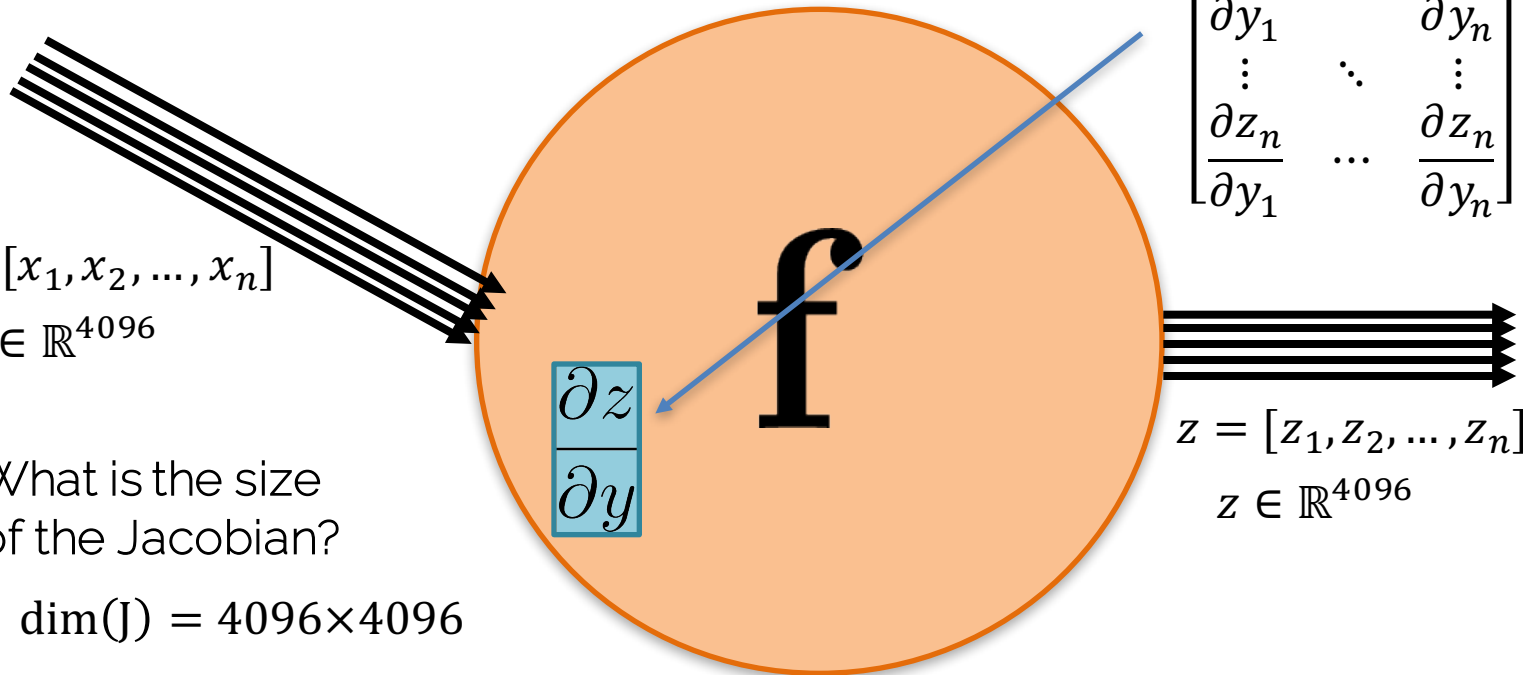
What is the size
of the Jacobian?

$$\dim(J) = 4096 \times 4096$$

$$\begin{bmatrix} \partial z \\ \partial y \end{bmatrix}$$

$$z = [z_1, z_2, \dots, z_n]$$

$$z \in \mathbb{R}^{4096}$$



Vectorized Operations

How efficient is that:

- **$\dim(\mathbf{J}) = 4096 \times 4096 = 16.78 \text{ mio}$**
- Assuming floats (i.e., 4 bytes / elem)
- $\rightarrow 64 \text{ MB}$

Typically, networks are run in batches:

- Assuming mini-batch size of 16
- $\rightarrow \mathbf{\dim(\mathbf{J}) = (16 \cdot 4096) \times (16 \cdot 4096) = 4295 \text{ mio}}$
- $\rightarrow 16.384 \text{ MB} = 16 \text{ GB}$

Jacobian Matrix:

$$\begin{bmatrix} \frac{\partial z_1}{\partial y_1} & \dots & \frac{\partial z_1}{\partial y_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_n}{\partial y_1} & \dots & \frac{\partial z_n}{\partial y_n} \end{bmatrix}$$

How to handle this?

Loss Functions

Naïve Losses

$$L^2 \text{ Loss: } L^2 = \sum_{i=1}^n (y_i - f(x_i))^2$$

- Sum of squared differences (SSD)
- Prone to outliers
- Compute-efficient (optimization)
- Optimum is the mean

$$L^1 \text{ Loss: } L^1 = \sum_{i=1}^n |y_i - f(x_i)|$$

- Sum of absolute differences
- Robust
- Costly to compute
- Optimum is the median

12	24	42	23
34	32	5	2
12	31	12	31
31	64	5	13

x_i

15	20	40	25
34	32	5	2
12	31	12	31
31	64	5	13

y_i

$$L^2(x, y) = 9 + 16 + 4 + 4 + 0 + \dots + 0 = 66$$

$$L^1(x, y) = 3 + 4 + 2 + 2 + 0 + \dots + 0 = 15$$

Hinge Loss (SVM Loss)

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

Hinge Loss (SVM Loss)

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

Given a function with weights W ,
Training pairs $[\mathbf{x}_i; \mathbf{y}_i]$ (input and labels)

Score function $s = f(\mathbf{x}_i, W)$

e.g., $f(\mathbf{x}_i, W) = W \cdot [x_0, x_1, \dots, x_N]^T$

Hinge Loss (SVM Loss)

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

Given a function with weights W ,
Training pairs $[x_i; y_i]$ (input and labels)

Score function $s = f(x_i, W)$
e.g., $f(x_i, W) = W \cdot [x_0, x_1, \dots, x_N]^T$

Suppose: 3 training examples and 3 classes



Hinge Loss (SVM Loss)

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

Given a function with weights W ,
Training pairs $[x_i; y_i]$ (input and labels)

Score function $s = f(x_i, W)$
e.g., $f(x_i, W) = W \cdot [x_0, x_1, \dots, x_N]^T$

Suppose: 3 training examples and 3 classes



scores	cat	3.2	1.3	2.2
	chair	5.1	4.9	2.5
	"car"	-1.7	2.0	-3.1

Hinge Loss (SVM Loss)

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

Given a function with weights W ,
Training pairs $[x_i; y_i]$ (input and labels)

Score function $s = f(x_i, W)$
e.g., $f(x_i, W) = W \cdot [x_0, x_1, \dots, x_N]^T$

Suppose: 3 training examples and 3 classes



scores	cat	3.2	1.3	2.2
	chair	5.1	4.9	2.5
	"car"	-1.7	2.0	-3.1

Loss 2.9

$$\begin{aligned} L_1 &= \\ &= \max(0, 5.1 - 3.2 + 1) + \\ &= \max(0, -1.7 - 3.2 + 1) = \\ &= \max(0, 2.9) + \max(0, -3.9) = \\ &= 2.9 + 0 = \\ &= \mathbf{2.9} \end{aligned}$$

Hinge Loss (SVM Loss)

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

Given a function with weights W ,
Training pairs $[x_i; y_i]$ (input and labels)

Score function $s = f(x_i, W)$
e.g., $f(x_i, W) = W \cdot [x_0, x_1, \dots, x_N]^T$

Suppose: 3 training examples and 3 classes



scores	cat	3.2	1.3	2.2
	chair	5.1	4.9	2.5
	"car"	-1.7	2.0	-3.1

Loss	2.9	0
------	-----	---

$$\begin{aligned} L_2 &= \\ &= \max(0, 1.3 - 4.9 + 1) + \\ &= \max(0, 2.0 - 4.9 + 1) = \\ &= \max(0, -2.6) + \max(0, -1.9) = \\ &= 0 + 0 = \\ &= 0 \end{aligned}$$

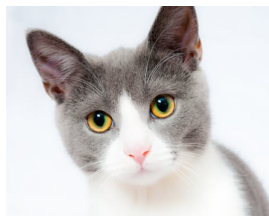
Hinge Loss (SVM Loss)

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

Given a function with weights W ,
Training pairs $[x_i; y_i]$ (input and labels)

Score function $s = f(x_i, W)$
e.g., $f(x_i, W) = W \cdot [x_0, x_1, \dots, x_N]^T$

Suppose: 3 training examples and 3 classes



scores	cat	3.2	1.3	2.2
	chair	5.1	4.9	2.5
	"car"	-1.7	2.0	-3.1

Loss	2.9	0	10.9
------	-----	---	------

$$\begin{aligned} L_3 &= \\ &= \max(0, 2.2 - (-3.1) + 1) + \\ &= \max(0, 2.5 - (-3.1) + 1) = \\ &= \max(0, 5.3) + \max(0, 5.6) = \\ &= 5.3 + 5.6 = \\ &= \mathbf{10.9} \end{aligned}$$

Hinge Loss (SVM Loss)

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

Given a function with weights W ,
Training pairs $[x_i; y_i]$ (input and labels)

Score function $s = f(x_i, W)$
e.g., $f(x_i, W) = W \cdot [x_0, x_1, \dots, x_N]^T$

Suppose: 3 training examples and 3 classes



scores	cat	3.2	1.3	2.2
	chair	5.1	4.9	2.5
	"car"	-1.7	2.0	-3.1
<hr/>				
	Loss	2.9	0	10.9

Full Loss (over all pairs):

$$\begin{aligned}
 L &= \frac{1}{N} \sum_{i=1}^N L_i = \\
 &= \frac{L_1 + L_2 + L_3}{3} = \\
 &= \frac{2.9 + 0 + 10.9}{3} = \\
 &= \mathbf{4.6}
 \end{aligned}$$

Weight Regularization & SVM Loss

Multiclass SVM loss $L_i = \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1)$

Full loss $L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \lambda R(W)$

$$L^1\text{-reg: } R^1(W) = \sum_{i=1}^N \sum_{j \neq y_i} |w_i|$$

$$L^2\text{-reg: } R^2(W) = \sum_{i=1}^N \sum_{j \neq y_i} w_i^2$$

Weight Regularization & SVM Loss

Multiclass SVM loss $L_i = \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1)$

Full loss $L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \lambda R(W)$

$$L^1\text{-reg: } R^1(W) = \sum_{i=1}^N \sum_{j \neq y_i} |w_i|$$

$$L^2\text{-reg: } R^2(W) = \sum_{i=1}^N \sum_{j \neq y_i} w_i^2$$

Example:

$$x = [1, 1, 1, 1]$$

$$w_1 = [1, 0, 0, 0]$$

$$R^2(w_1) = 1$$

$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

$$R^2(w_2) = 0.25^2 + 0.25^2 + 0.25^2 + 0.25^2 = 0.25$$

$$w_1^T x = w_2^T x = 1$$

Cross-Entropy (Softmax)

Softmax
$$L_i = -\log\left(\frac{e^{s y_i}}{\sum_j e^{s_j}}\right)$$

Score function
$$s = f(x_i, W)$$

e.g.,
$$f(x_i, W) = W \cdot [x_0, x_1, \dots, x_N]^T$$

Given a function with weights W ,
Training pairs $[x_i; y_i]$ (input and labels)

Suppose: 3 training examples and 3 classes



scores	cat	3.2	1.3	2.2
	chair	5.1	4.9	2.5
	"car"	-1.7	2.0	-3.1

Loss

Cross-Entropy (Softmax)

Softmax
$$L_i = -\log\left(\frac{e^{s y_i}}{\sum_j e^{s_j}}\right)$$

Score function
$$s = f(x_i, W)$$

e.g.,
$$f(x_i, W) = W \cdot [x_0, x_1, \dots, x_N]^T$$

Given a function with weights W ,
Training pairs $[x_i; y_i]$ (input and labels)

Suppose: 3 training examples and 3 classes



scores	cat	3.2	1.3	2.2
	chair	5.1	4.9	2.5
	"car"	-1.7	2.0	-3.1

3.2
5.1
-1.7

Loss

Cross-Entropy (Softmax)

Softmax
$$L_i = -\log\left(\frac{e^{s y_i}}{\sum_j e^{s_j}}\right)$$

Score function
$$s = f(x_i, W)$$

e.g.,
$$f(x_i, W) = W \cdot [x_0, x_1, \dots, x_N]^T$$

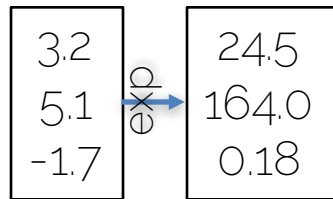
Given a function with weights W ,
Training pairs $[x_i; y_i]$ (input and labels)

Suppose: 3 training examples and 3 classes



scores	cat	3.2	1.3	2.2
	chair	5.1	4.9	2.5
	"car"	-1.7	2.0	-3.1

Loss



Cross-Entropy (Softmax)

Softmax
$$L_i = -\log\left(\frac{e^{s y_i}}{\sum_j e^{s_j}}\right)$$

Score function
$$s = f(x_i, W)$$

e.g.,
$$f(x_i, W) = W \cdot [x_0, x_1, \dots, x_N]^T$$

Suppose: 3 training examples and 3 classes



scores	cat	3.2	1.3	2.2
	chair	5.1	4.9	2.5
	"car"	-1.7	2.0	-3.1

Loss

Given a function with weights W ,
Training pairs $[x_i; y_i]$ (input and labels)



Cross-Entropy (Softmax)

Softmax
$$L_i = -\log\left(\frac{e^{s y_i}}{\sum_j e^{s_j}}\right)$$

Score function
$$s = f(x_i, W)$$

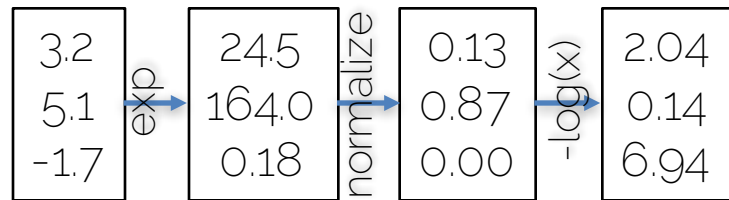
 e.g.,
$$f(x_i, W) = W \cdot [x_0, x_1, \dots, x_N]^T$$

Suppose: 3 training examples and 3 classes



scores	cat	3.2	1.3	2.2
	chair	5.1	4.9	2.5
	"car"	-1.7	2.0	-3.1
Loss		2.04	0.14	6.94

Given a function with weights W ,
 Training pairs $[x_i; y_i]$ (input and labels)



Cross-Entropy (Softmax)

Softmax $L_i = -\log\left(\frac{e^{s y_i}}{\sum_j e^{s_j}}\right)$

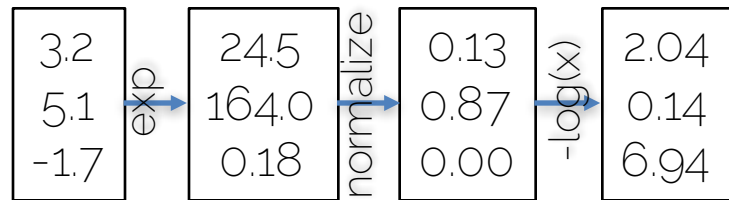
Score function $s = f(x_i, W)$
 e.g., $f(x_i, W) = W \cdot [x_0, x_1, \dots, x_N]^T$

Suppose: 3 training examples and 3 classes



scores	cat	3.2	1.3	2.2
	chair	5.1	4.9	2.5
	"car"	-1.7	2.0	-3.1
Loss		2.04	0.14	6.94

Given a function with weights W ,
 Training pairs $[x_i; y_i]$ (input and labels)



$$\begin{aligned}
 L &= \frac{1}{N} \sum_{i=1}^N L_i = \\
 &= \frac{L_1 + L_2 + L_3}{3} = \\
 &= \frac{2.04 + 0.14 + 6.94}{3} = \\
 &= \mathbf{9.12}
 \end{aligned}$$

Hinge Loss vs Softmax

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

Softmax: $L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$

Hinge Loss vs Softmax

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

$$\text{Softmax: } L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

Given the following scores:

$$s = [5, -3, 2]$$

$$s = [5, 10, 10]$$

$$s = [5, -20, -20]$$

$$y_i = 0$$

Hinge Loss vs Softmax

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

Softmax: $L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$

Given the following scores:

$$s = [5, -3, 2]$$

Hinge loss:

$$\begin{aligned} &\max(0, -3 - 5 + 1) + \\ &\max(0, 2 - 5 + 1) = 0 \end{aligned}$$

$$s = [5, 10, 10]$$

$$\begin{aligned} &\max(0, 10 - 5 + 1) + \\ &\max(0, 10 - 5 + 1) = 12 \end{aligned}$$

$$s = [5, -20, -20]$$

$$\begin{aligned} &\max(0, -20 - 5 + 1) + \\ &\max(0, -20 - 5 + 1) = 0 \end{aligned}$$

$$y_i = 0$$

Hinge Loss vs Softmax

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

Softmax: $L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$

Given the following scores:

$$s = [5, -3, 2]$$

Hinge loss:

$$\max(0, -3 - 5 + 1) + \max(0, 2 - 5 + 1) = 0$$

Softmax loss:

Google...
 $-\ln(e^5 / (e^5 + e^{-3} + e^2)) = 0.05$

$$s = [5, 10, 10]$$

$$\max(0, 10 - 5 + 1) + \max(0, 10 - 5 + 1) = 12$$

Google...
 $-\ln(e^5 / (e^5 + e^{10} + e^{10})) = 5.70$

$$s = [5, -20, -20]$$

$$\max(0, -20 - 5 + 1) + \max(0, -20 - 5 + 1) = 0$$

Google...
 $-\ln(e^5 / (e^5 + e^{-20} + e^{-20})) = 2.e-11$

$y_i = 0$

Hinge Loss vs Softmax

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

Softmax: $L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$

Given the following scores:

$$s = [5, -3, 2]$$

Hinge loss:

$$\max(0, -3 - 5 + 1) + \max(0, 2 - 5 + 1) = 0$$

Softmax loss:

Google...
 $-\ln(e^5 / (e^5 + e^{-3} + e^2)) = 0.05$

$$s = [5, 10, 10]$$

$$\max(0, 10 - 5 + 1) + \max(0, 10 - 5 + 1) = 12$$

Google...
 $-\ln(e^5 / (e^5 + e^{10} + e^{10})) = 5.70$

$$s = [5, -20, -20]$$

$$\max(0, -20 - 5 + 1) + \max(0, -20 - 5 + 1) = 0$$

Google...
 $-\ln(e^5 / (e^5 + e^{-20} + e^{-20})) = 2.e-11$

$y_i = 0$

Softmax *always* wants to improve!

Hinge Loss saturates

Loss in Compute Graph

Score function $s = f(x_i, W)$
e.g., $f(x_i, W) = W \cdot [x_0, x_1, \dots, x_N]^T$

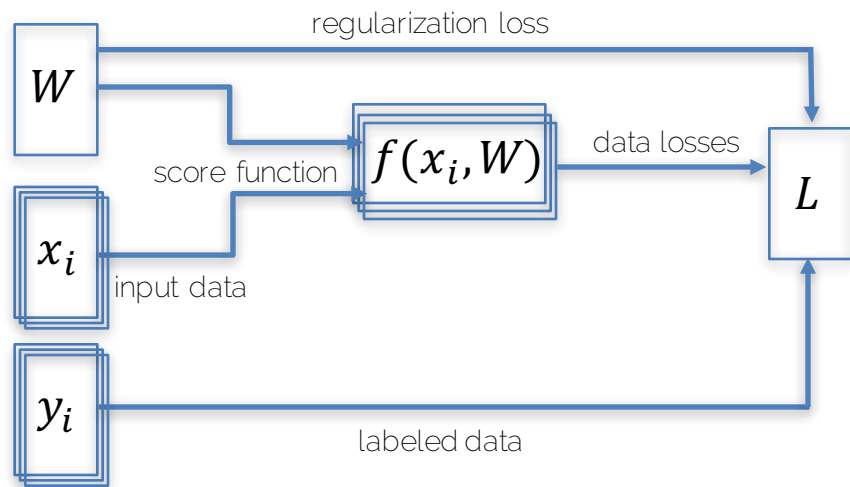
Given a function with weights W ,
Training pairs $[x_i; y_i]$ (input and labels)

Softmax $L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$

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

e.g., L^2 -reg: $R^2(W) = \sum_{i=1}^N w_i^2$

Full Loss $L = \frac{1}{N} \sum_{i=1}^N L_i + R^2(W)$



Compute Graphs

Score function $s = f(x_i, W)$
e.g., $f(x_i, W) = W \cdot [x_0, x_1, \dots, x_N]^T$

Softmax $L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$

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

e.g., L^2 -reg: $R^2(W) = \sum_{i=1}^N w_i^2$

Full Loss $L = \frac{1}{N} \sum_{i=1}^N L_i + R^2(W)$

Compute Graphs

Score function $s = f(x_i, W)$
e.g., $f(x_i, W) = W \cdot [x_0, x_1, \dots, x_N]^T$

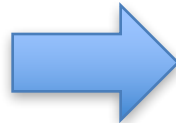
Softmax $L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$

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

e.g., L^2 -reg: $R^2(W) = \sum_{i=1}^N w_i^2$

Full Loss $L = \frac{1}{N} \sum_{i=1}^N L_i + R^2(W)$

Want to find optimal W . I.e., weights are unknowns of optimization problem



Compute gradient w.r.t. W .

Gradient $\nabla_W L$ is computed via backpropagation

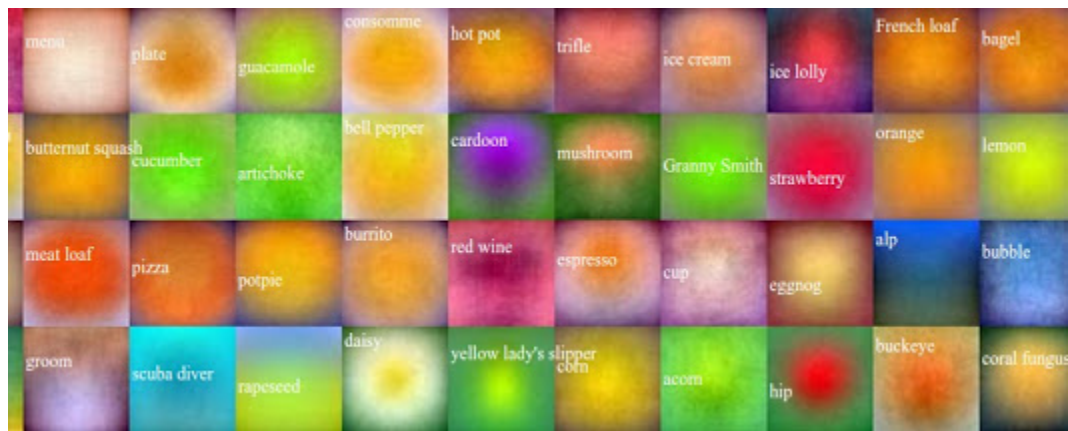
Introduction to Neural Networks

Neural Network

- Linear score function $f = Wx$



On CIFAR-10



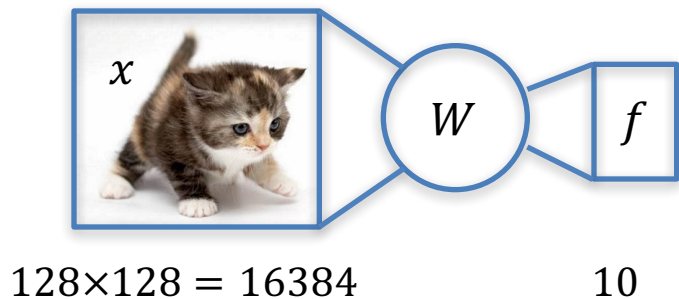
On ImageNet

Neural Network

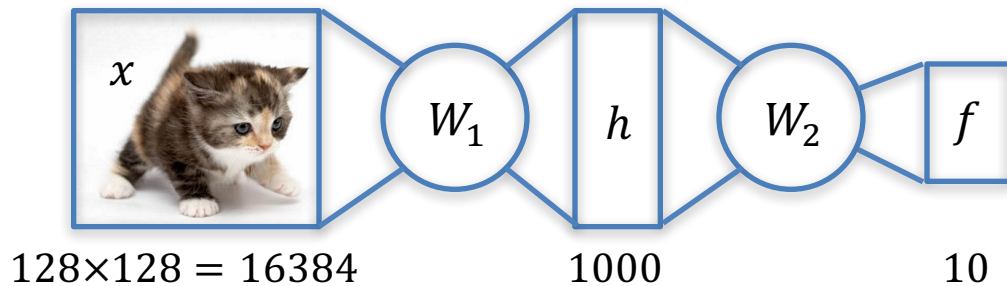
- Linear score function $f = Wx$
- Neural network is a nesting of 'functions'
 - 2-layers: $f = W_2 \max(0, W_1 x)$
 - 3-layers: $f = W_3 \max(0, W_2 \max(0, W_1 x))$
 - 4-layers: $f = W_4 \tanh(W_3, \max(0, W_2 \max(0, W_1 x)))$
 - 5-layers: $f = W_5 \sigma(W_4 \tanh(W_3, \max(0, W_2 \max(0, W_1 x))))$
 - ... up to hundreds of layers

Neural Network

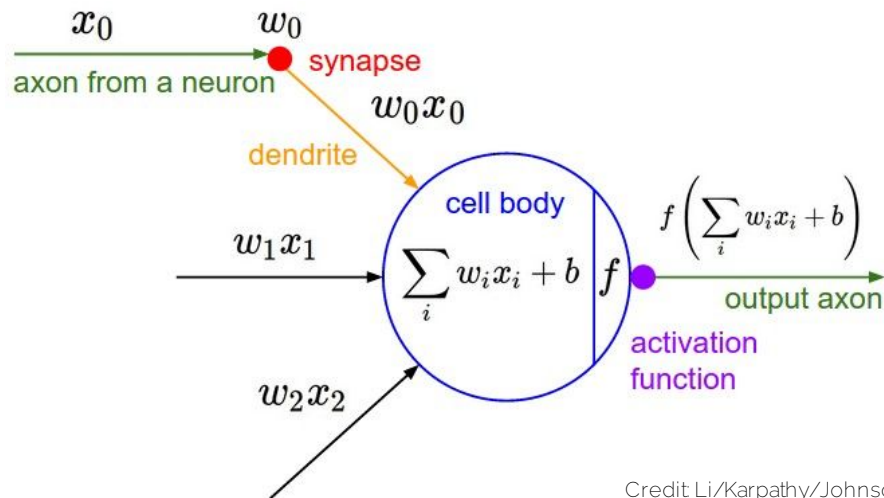
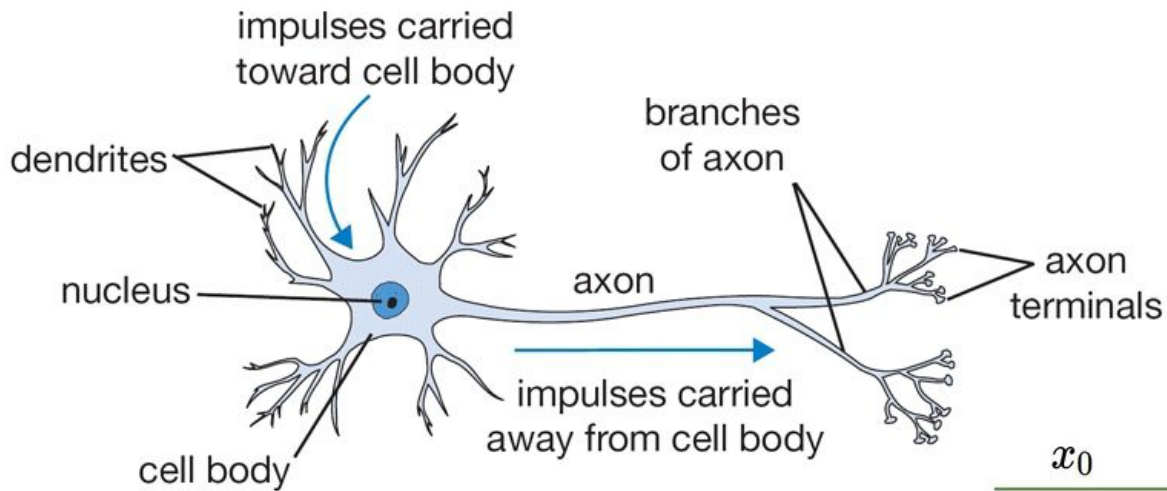
1-layer network: $f = Wx$



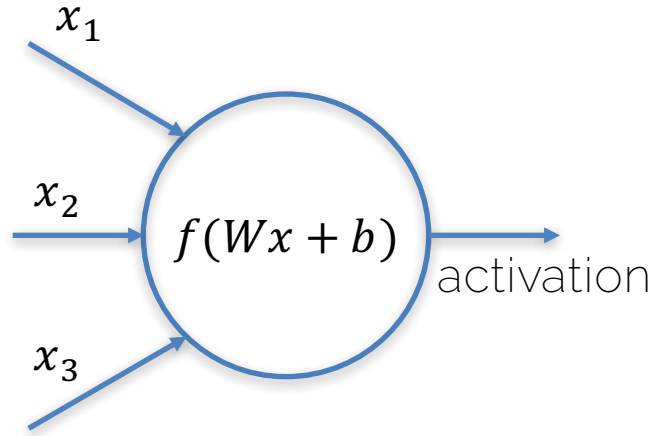
2-layer network: $f = W_2 \max(0, W_1 x)$



Neurons



Neurons

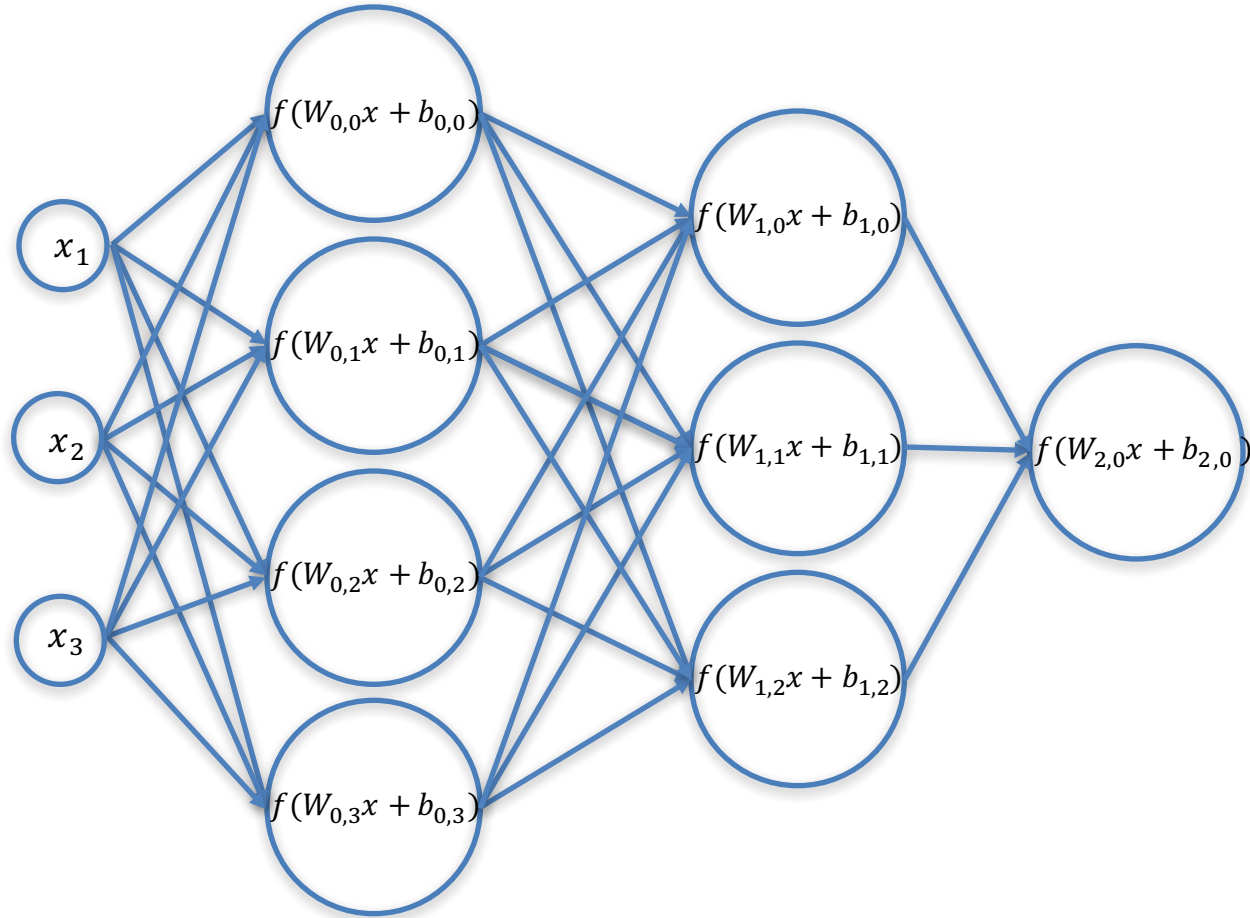


Linear function: $Wx + b$

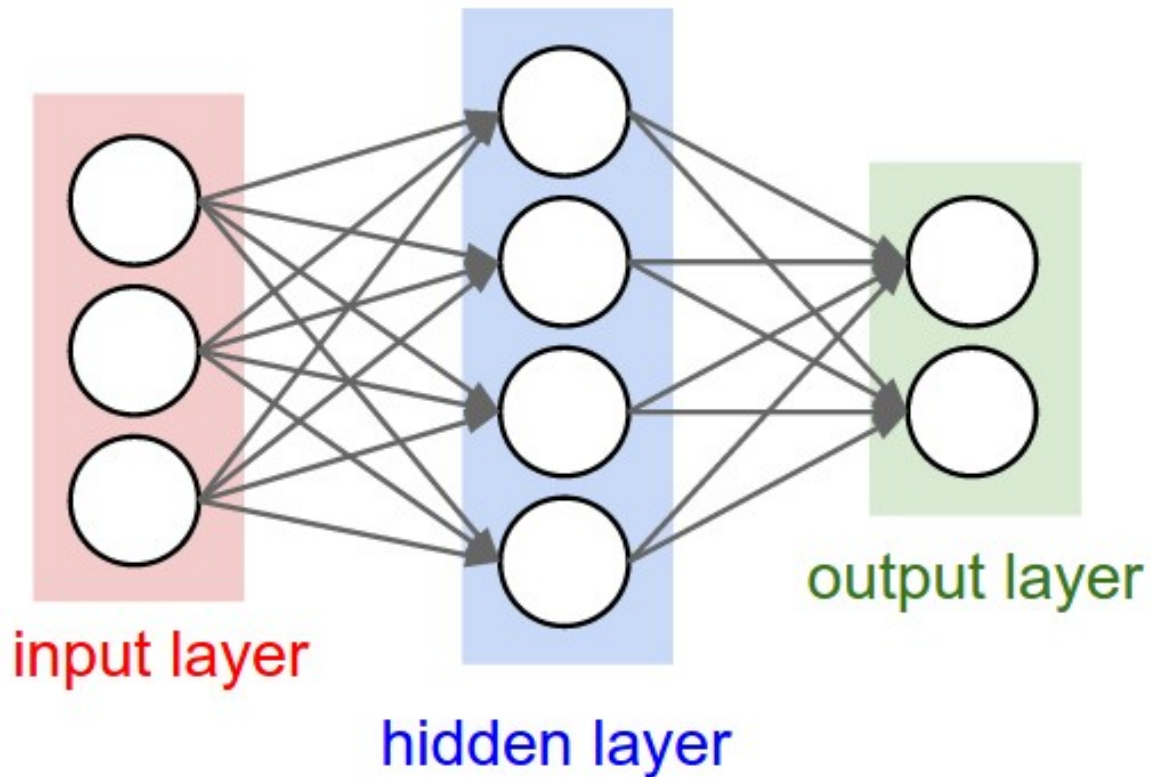
Non-linearity (activation): $f(x)$

Every neuron computes: $f(Wx + b)$

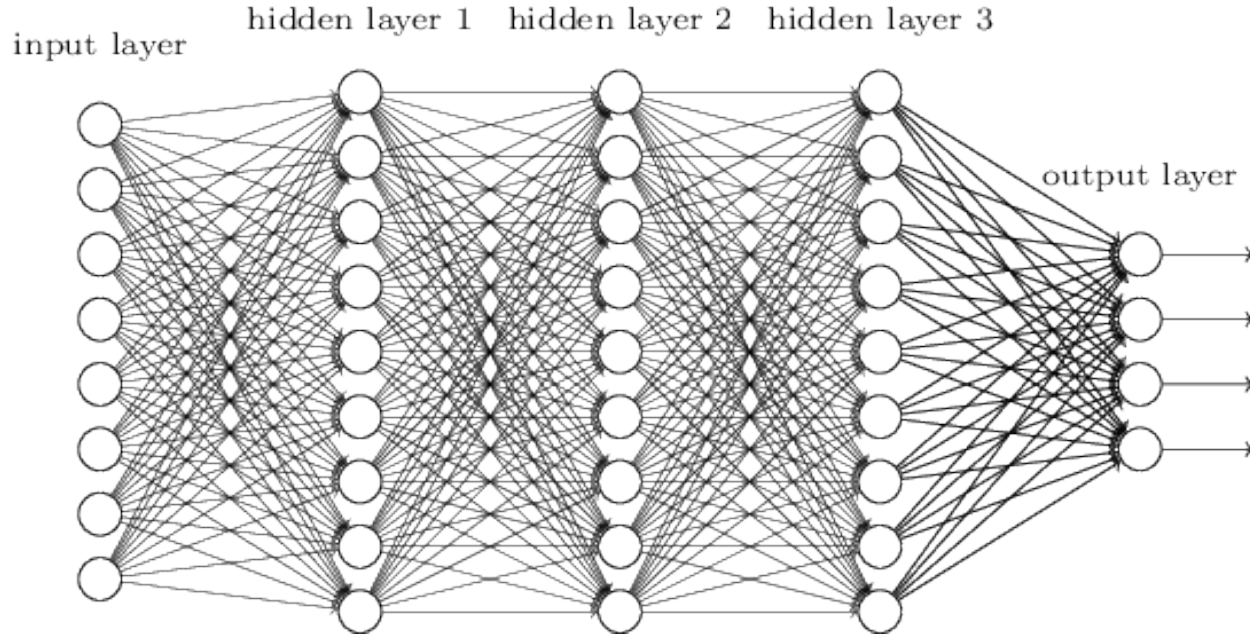
Net of Neurons



Neural Network



Neural Network



Neural Network

$$f = W_3 \cdot (W_2 \cdot (W_1 \cdot x))$$

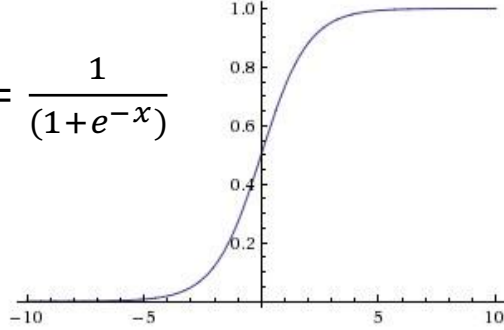
Why activation functions?

Why not just concatenate?

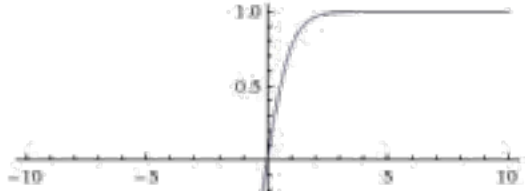
Would be much cheaper to
compute...

Activation Functions

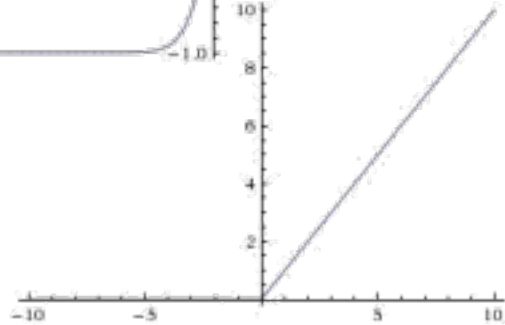
Sigmoid: $\sigma(x) = \frac{1}{1+e^{-x}}$



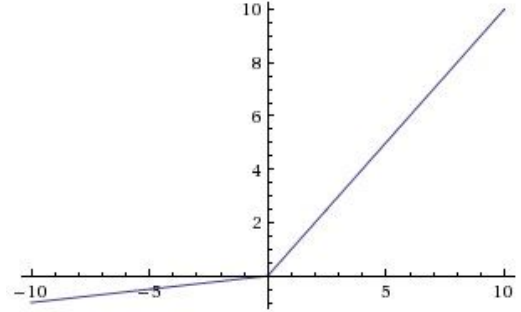
tanh: $\tanh(x)$



ReLU: $\max(0, x)$



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



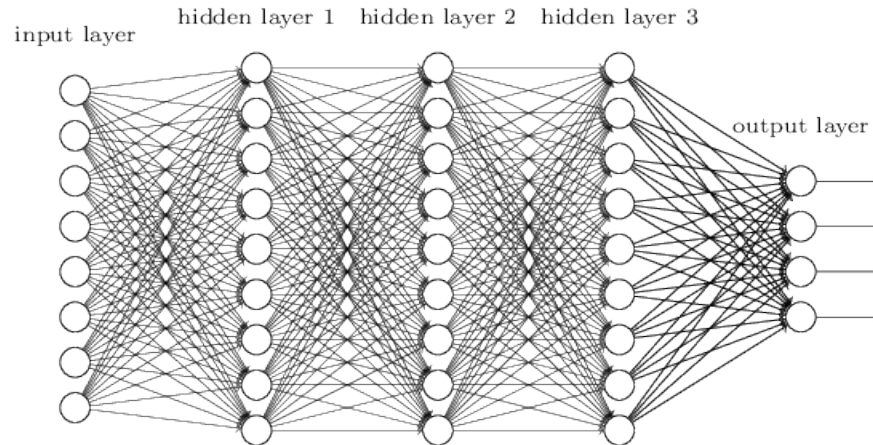
Parametric ReLU: $\max(\alpha x, x)$

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

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

Neural Network

Why organize a neural network into layers?



Neural Network

- Summary
 - Given a dataset with ground truth training pairs $[\mathbf{x}_i; \mathbf{y}_i]$,
 - Find optimal weights \mathbf{W} using stochastic gradient descent, such that the loss function is minimized
 - Compute gradients with backpropagation (use batch-mode; more later)
 - Iterate many times over training set (SGD; more later)

Artificial Neural Network vs Brain



Artificial neural networks: inspired but not even close to the brain!
It's much more complex than simple linearity + activations
Great for the media and news articles 😊

Administrative Things

- Next Thursday 18th of May:
 - More on Neural Networks 😊
 - More theory, best practices, applications
- Tomorrow: 1st Q&A session of the first assignment
- Tentative date for the exam: 18th of August
(not confirmed though)