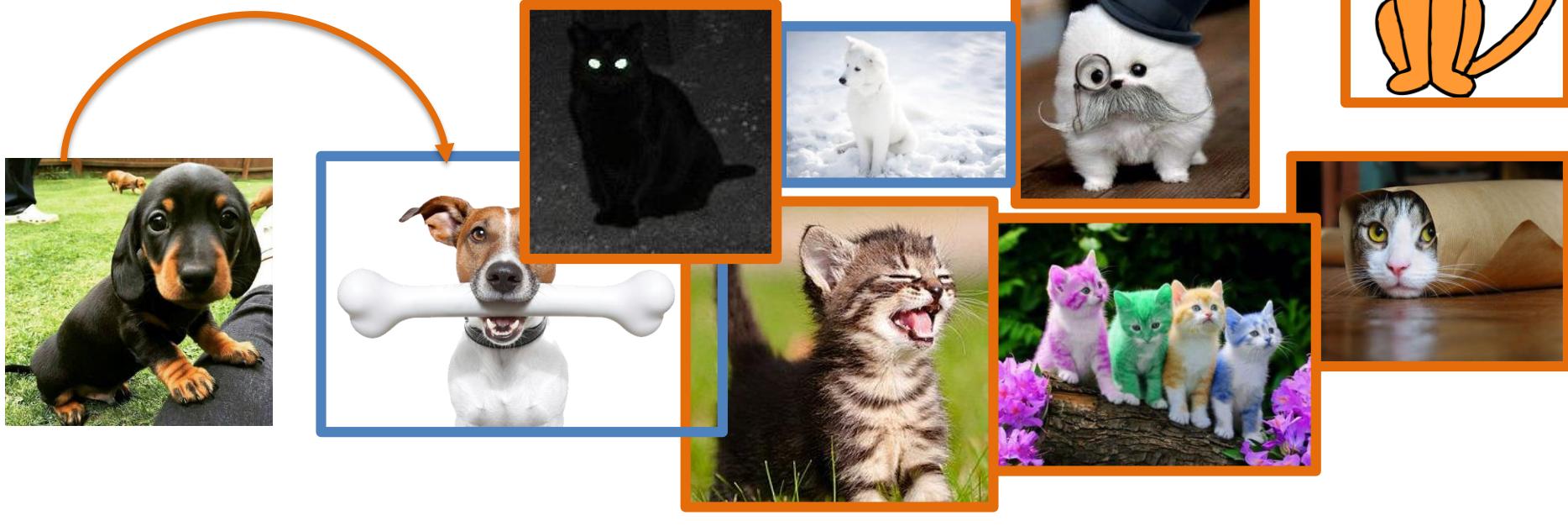


Lecture 2 Recap

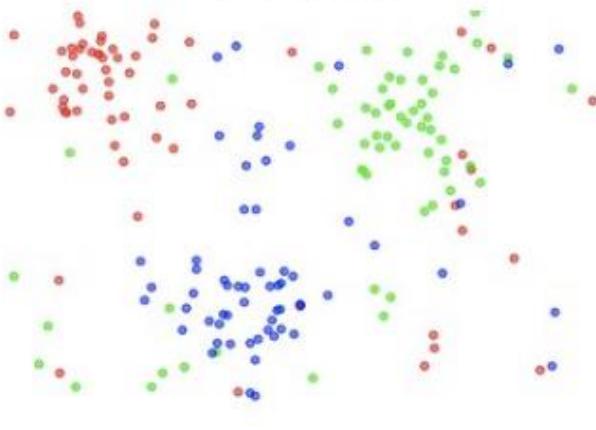
Nearest Neighbor

NN classifier = dog

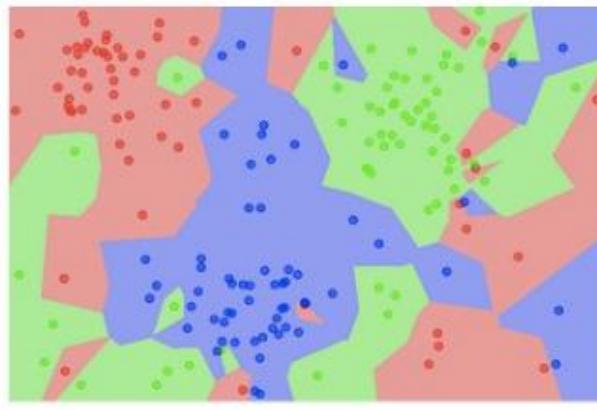


Nearest Neighbor

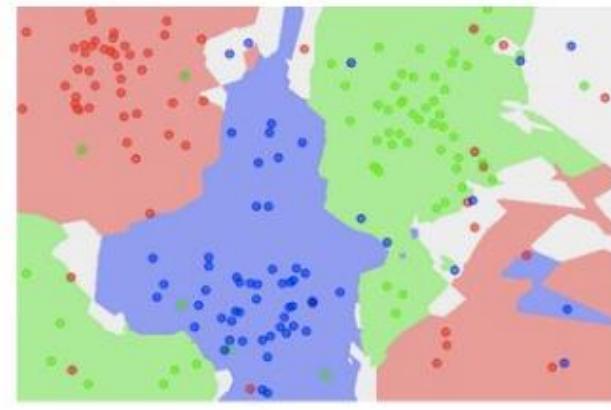
the data



NN classifier



5-NN classifier



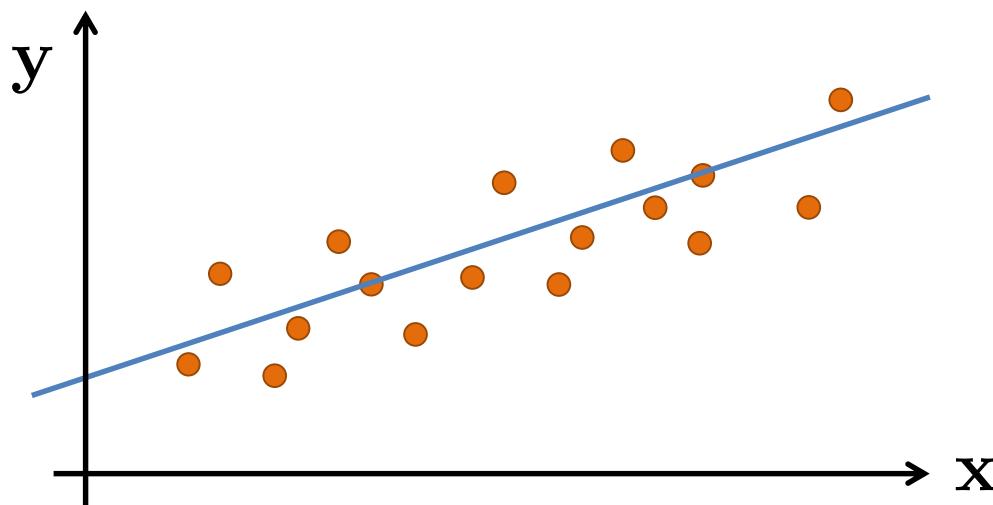
What is the performance on training data for NN classifier?

What classifier is more likely to perform best on test data?

Courtesy of Stanford course cs231n

Linear Regression

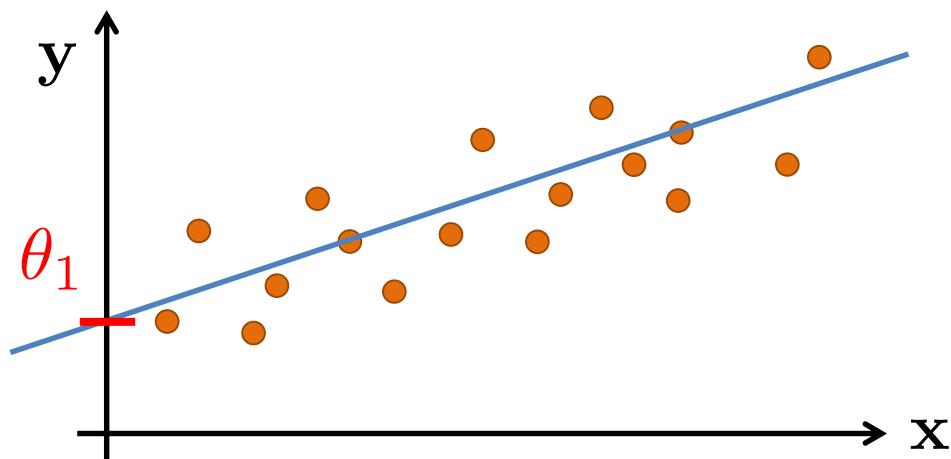
- Supervised learning
- Find a linear model that explains a target \mathbf{y} given the inputs \mathbf{X}



Linear Regression

- A linear model is expressed in the form

$$\hat{y}_i = \sum_{j=1}^d x_{ij}\theta_j = x_{i1}\theta_1 + x_{i2}\theta_2 + \cdots + x_{id}\theta_d$$



Logistic Regression

- Loss function

$$\mathcal{L}(\hat{y}_i, y_i) = y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)$$

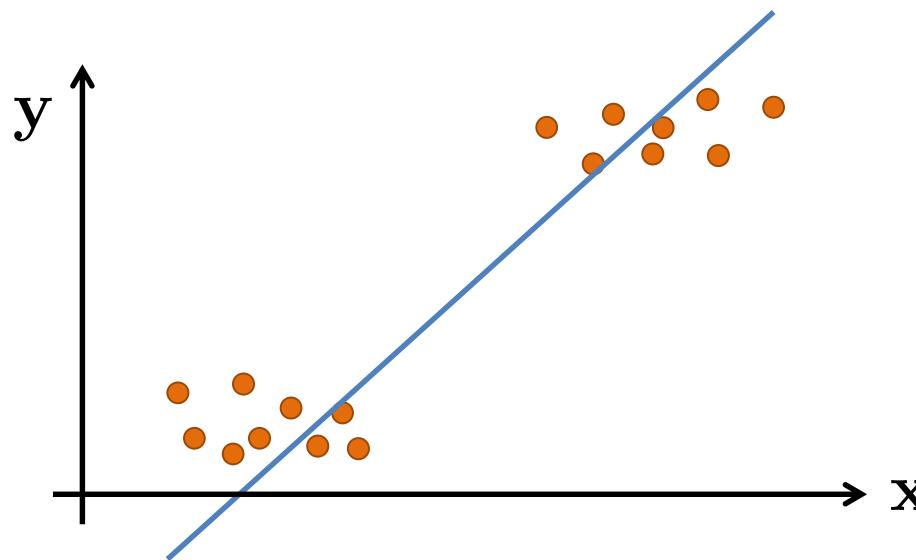
- Cost function

$$C(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^n y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)$$

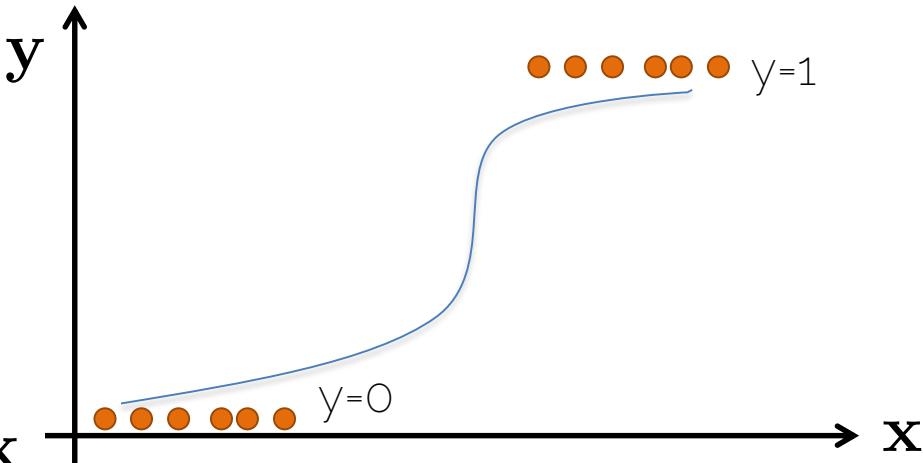
Minimization

The diagram shows two orange arrows originating from the term $\log \hat{y}_i$ in the cost function equation. One arrow points downwards to the word "Minimization". The other arrow points to the right to the equation $\hat{y}_i = \sigma(\mathbf{x}_i \boldsymbol{\theta})$.

Linear vs Logistic Regression



Prediction can exceed
training sample range
-> in case of classification [0;1] that's a real issue



Prediction is guaranteed
to be within 0 and 1

(Vanilla) Gradient Descent

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



Learning rate

$$\nabla_{\mathbf{x}} f(\mathbf{x})$$

\mathbf{x}



$$\mathbf{x}^* = \arg \min f(\mathbf{x})$$

Introduction to Neural Networks

Neural Network

- Linear score function $f = Wx$



On CIFAR-10



Credit: Li/Karpathy/Johnson

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

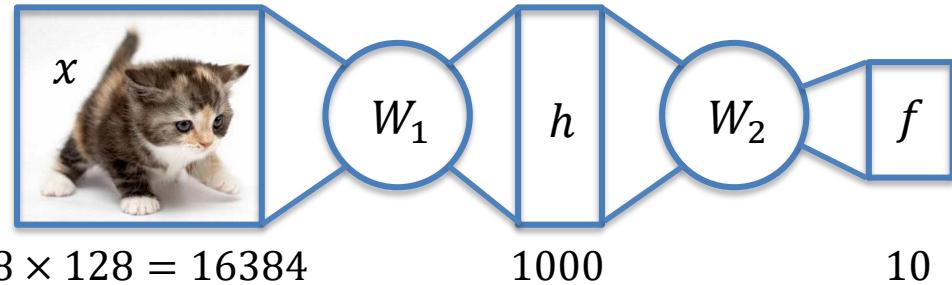
1-layer network: $f = Wx$



$$128 \times 128 = 16384$$

$$10$$

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

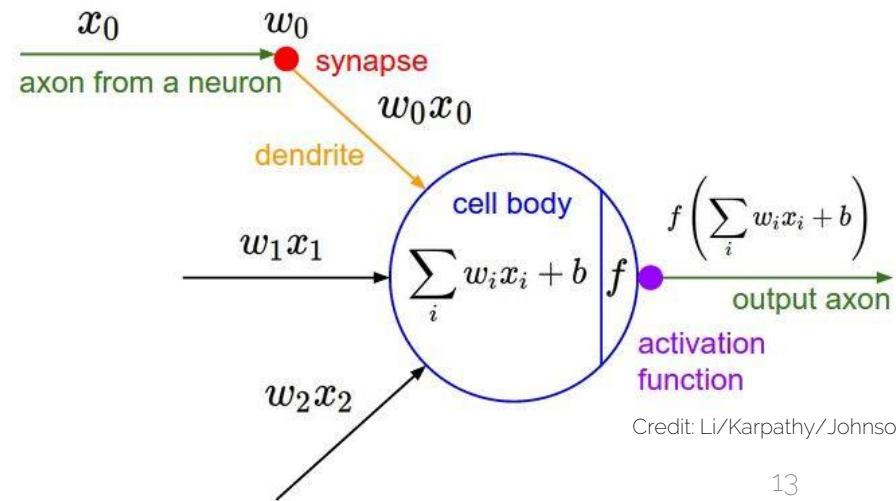
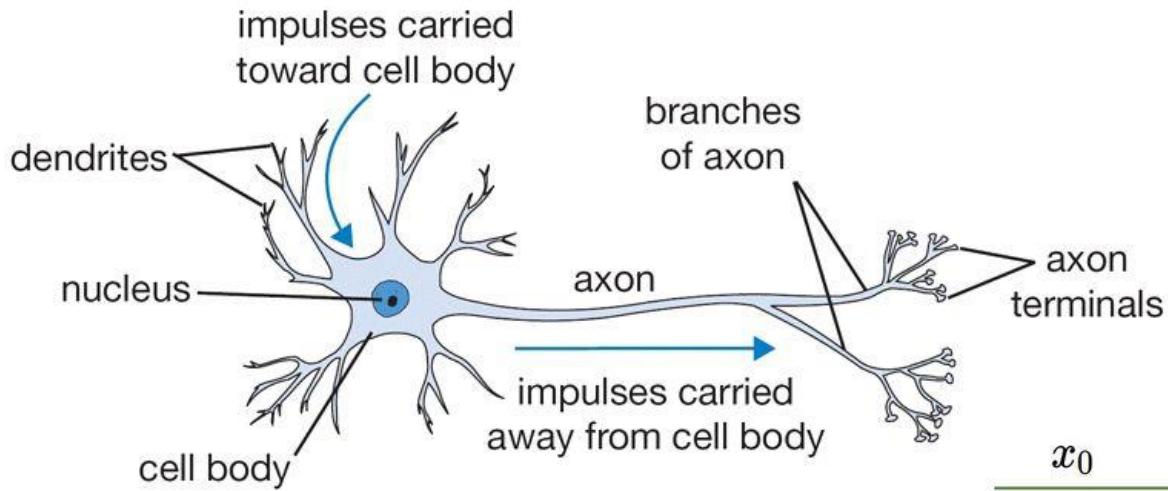


$$128 \times 128 = 16384$$

$$1000$$

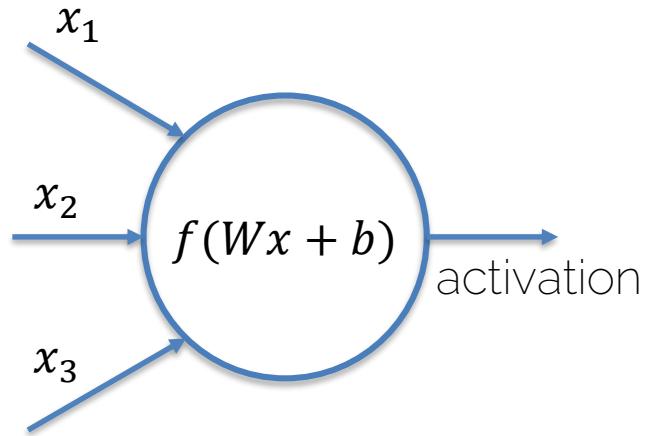
$$10$$

Neurons



Credit: Li/Karpathy/Johnson

Neurons

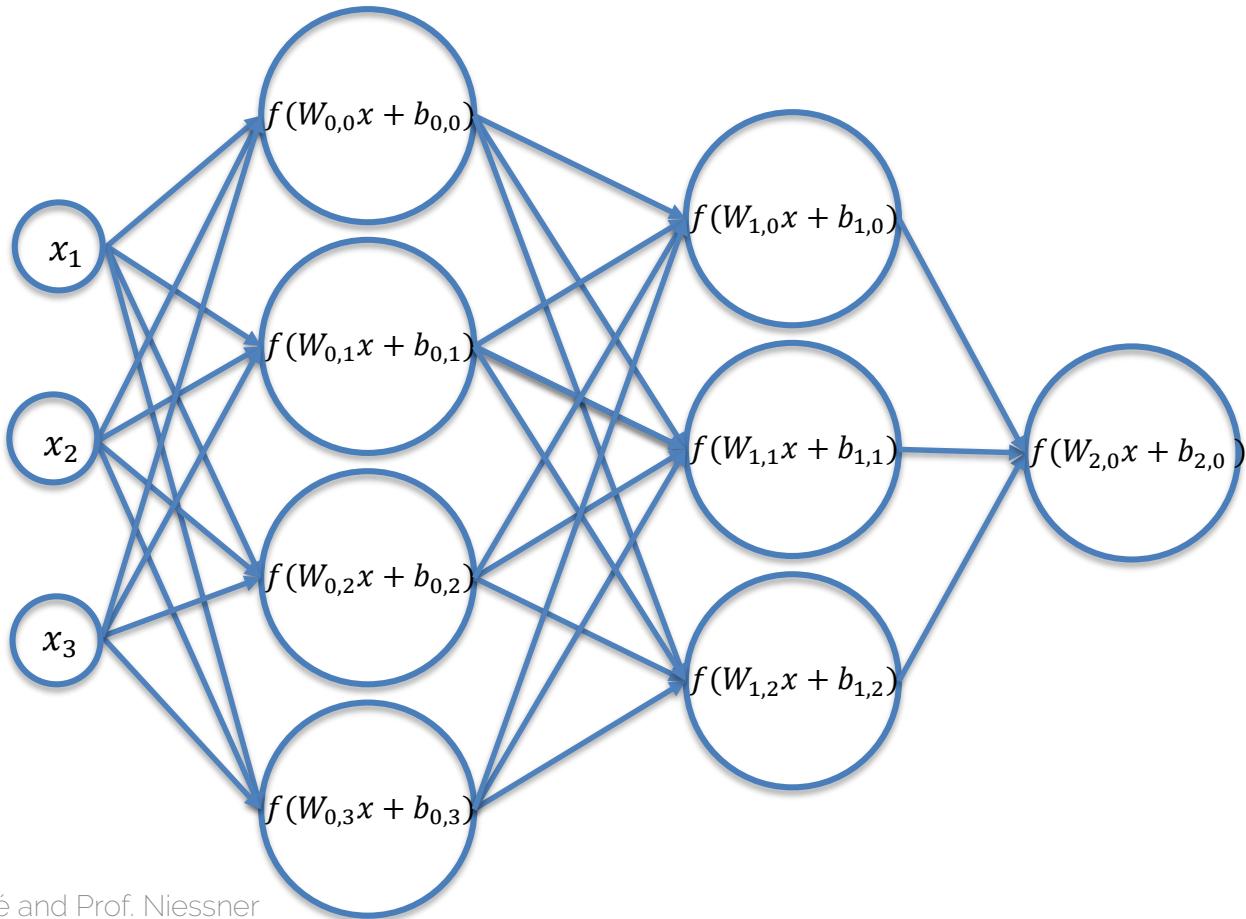


Linear function: $Wx + b$

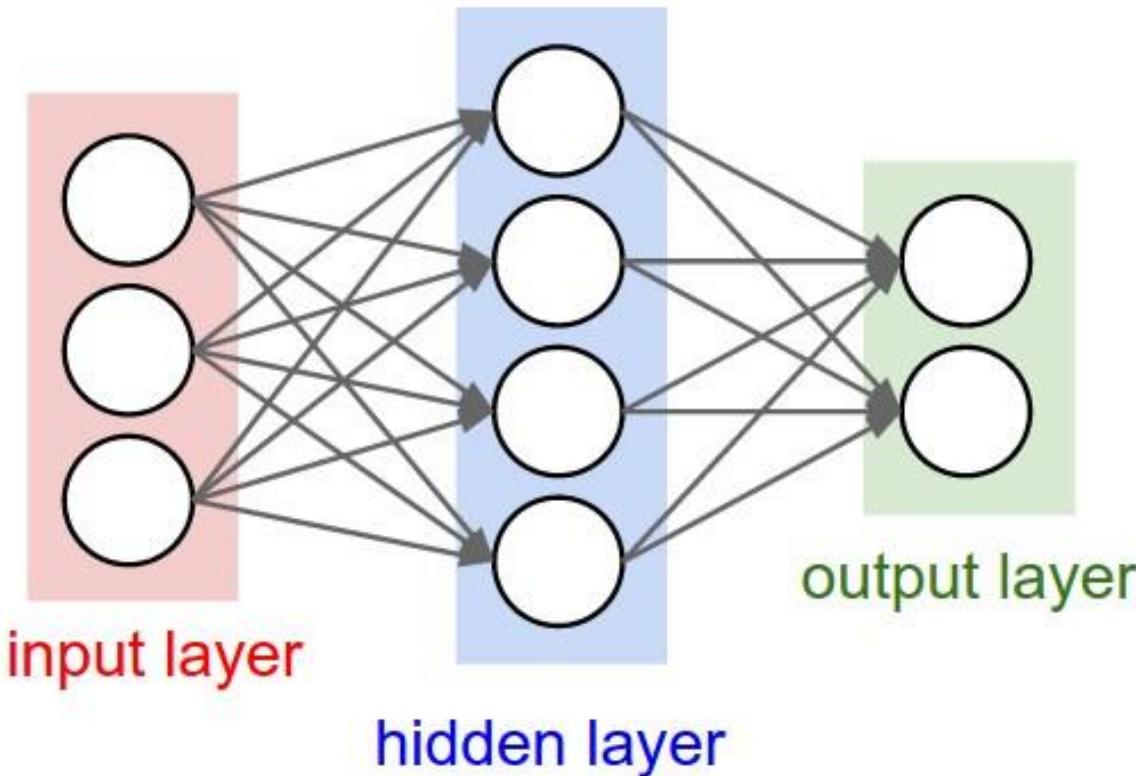
Non-linearity: activation: $f(x)$

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

Net of Neurons

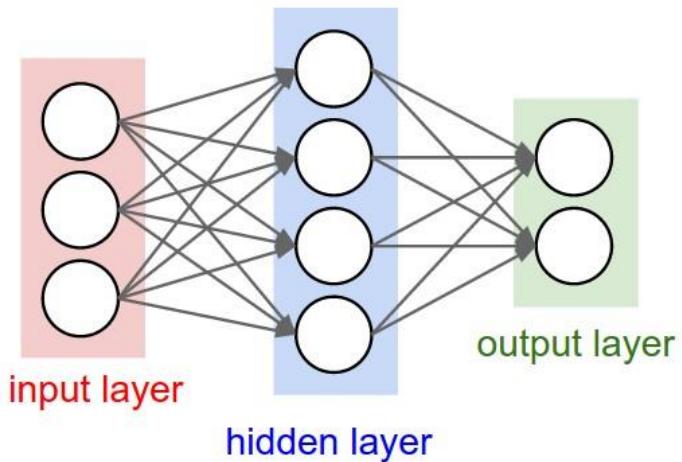


Neural Network

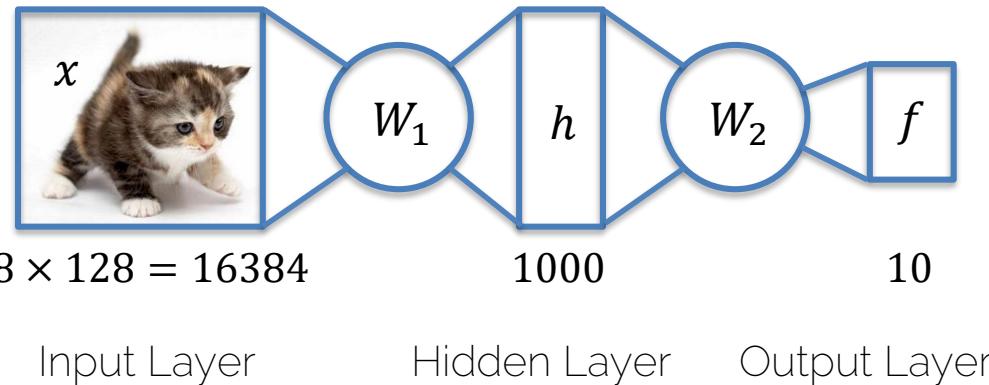


Credit: Li/Karpathy/Johnson

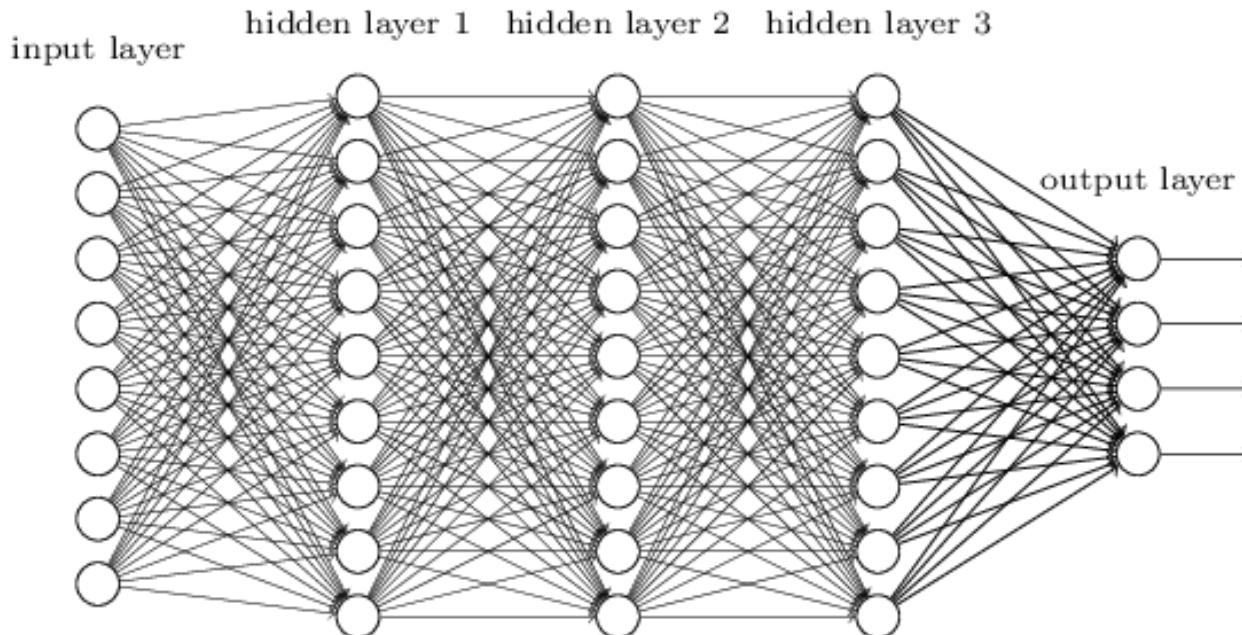
Neural Network



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



Neural Network



Neural Network

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

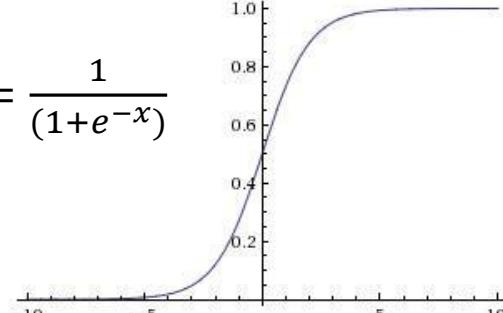
Why activation functions?

Why not just concatenate?

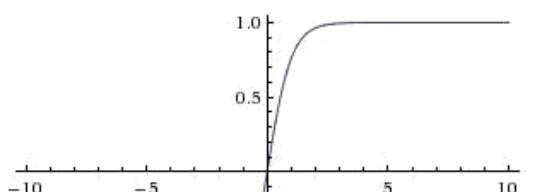
Would be much cheaper to
compute....

Activation Functions

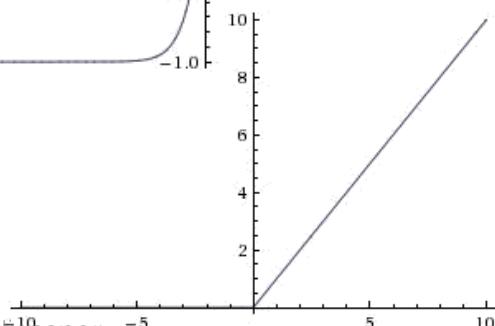
Sigmoid: $\sigma(x) = \frac{1}{(1+e^{-x})}$



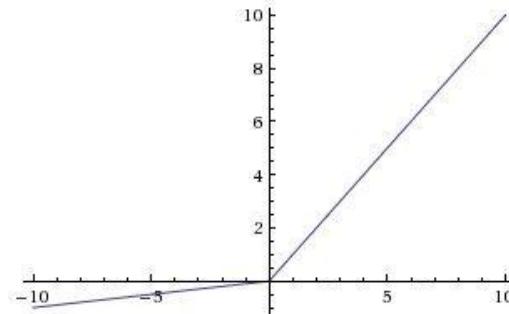
tanh: $\tanh(x)$



ReLU: $\max(0, x)$



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



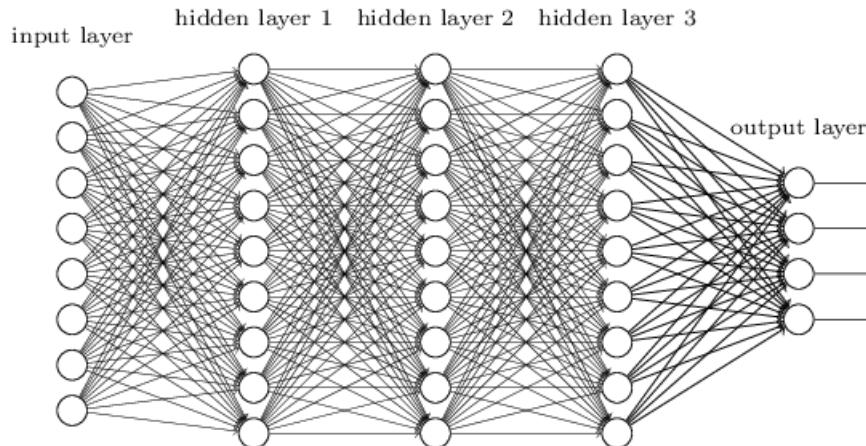
Parametric ReLU: $\max(\alpha x, x)$

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 $[x_i; y_i]$,
 - 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 ☺

Artificial Neural Network vs Brain



U.S.

World

Opinion

Politics

Entertainment

Business

Lifestyle

TV

Radio

More :



Login

Wa

Hot Topics

'Unprecedented' NoKo threat

| Another ESPN controversy

| Serial killer scare

GOOGLE · 3 days ago

Google's artificial intelligence computer 'no longer constrained by limits of human knowledge'



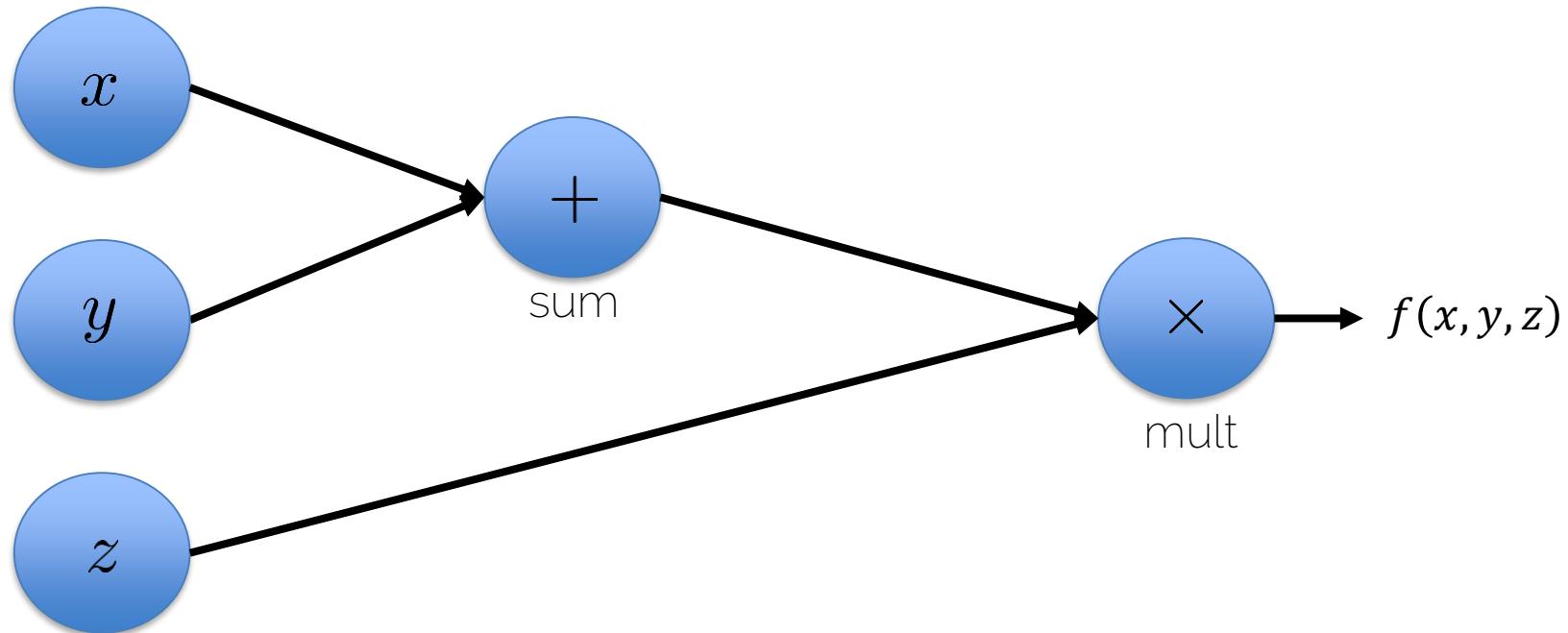
Computational Graphs

Computational Graphs

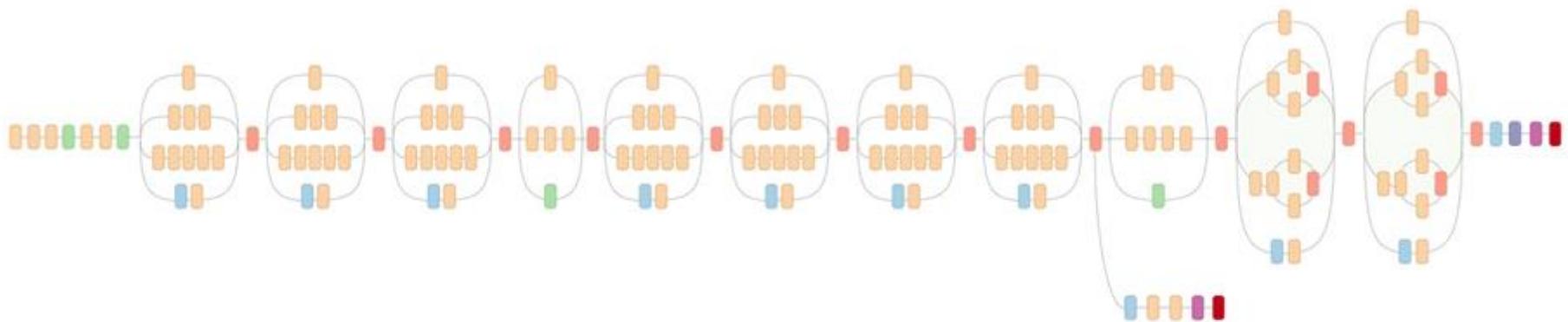
- Neural network is a computational graph
 - It has compute nodes
 - It has edges that connect nodes
 - It is directional
 - It is organized in 'layers'

Computational Graphs

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



Computational Graphs

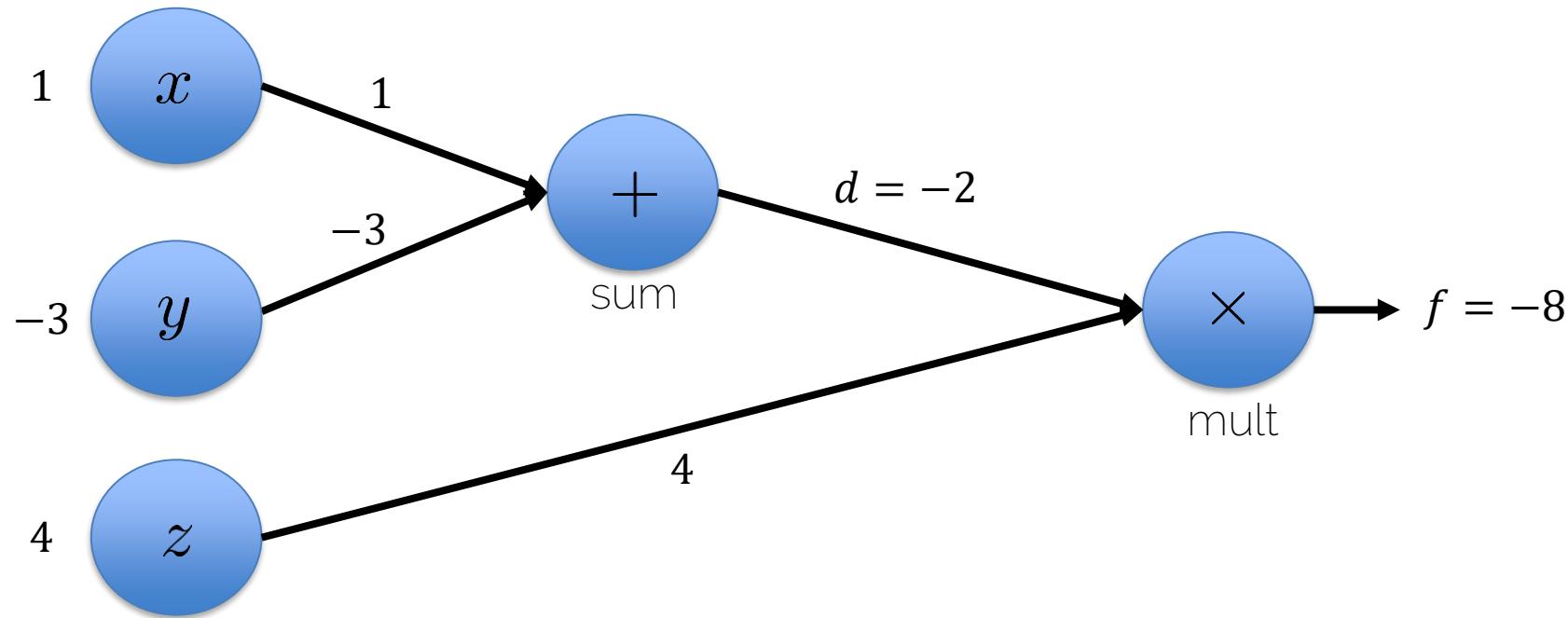


- Convolution
- AvgPool
- MaxPool
- Concat
- Dropout
- Fully connected
- Softmax

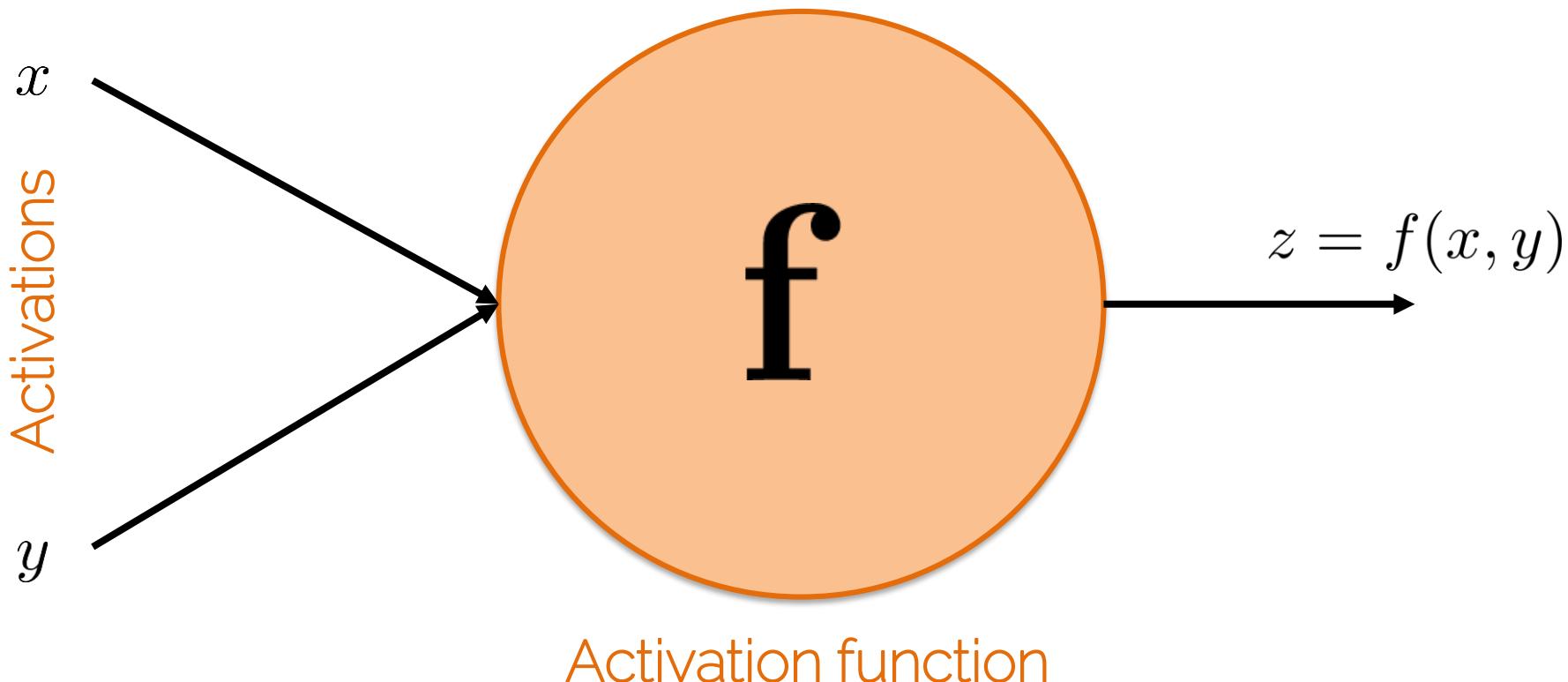
Another view of GoogLeNet's architecture.

Evaluation: Forward Pass

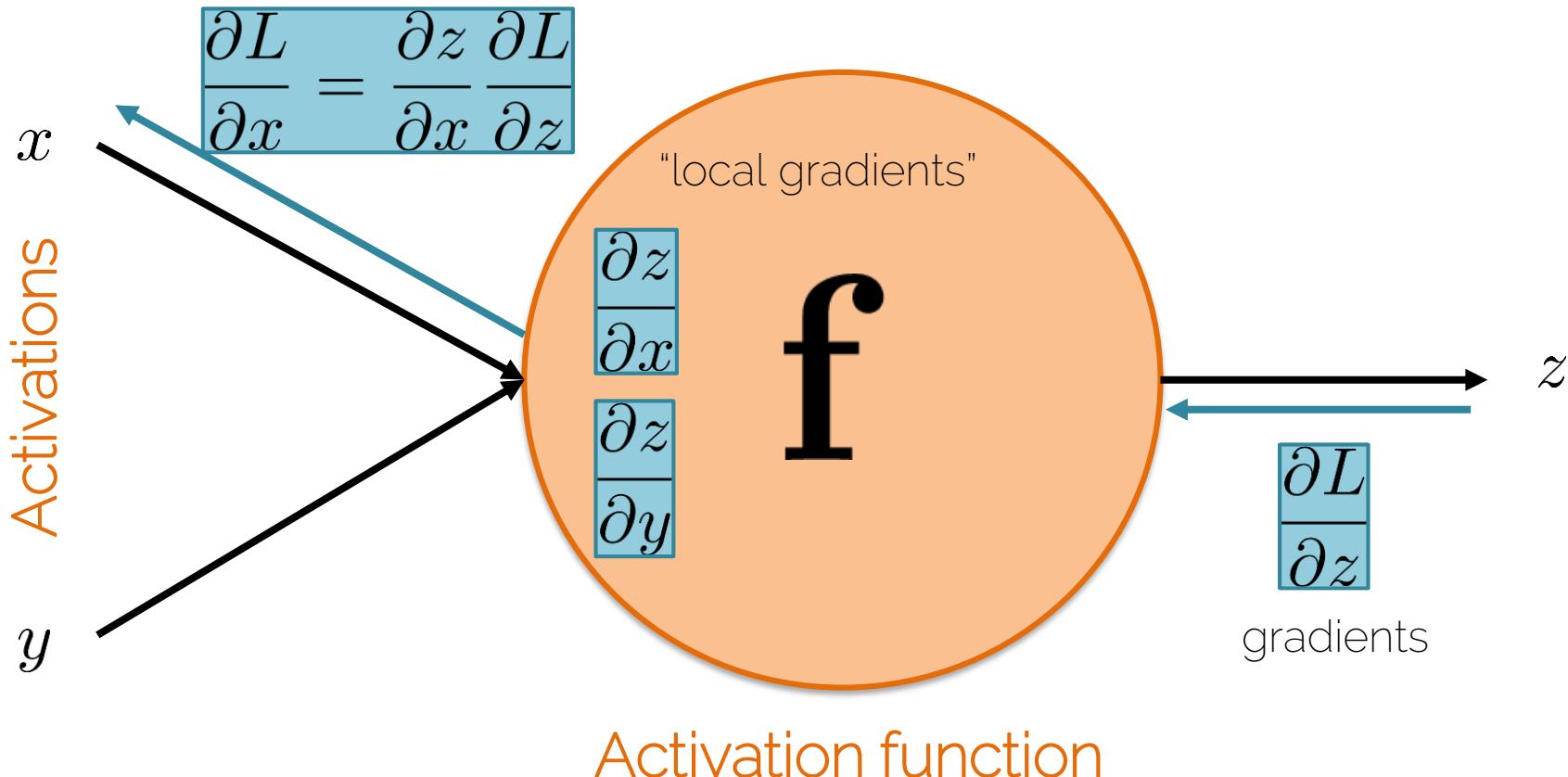
- $f(x, y, z) = (x + y) \cdot z$ Initialization $x = 1, y = -3, z = 4$



The Flow of Gradients



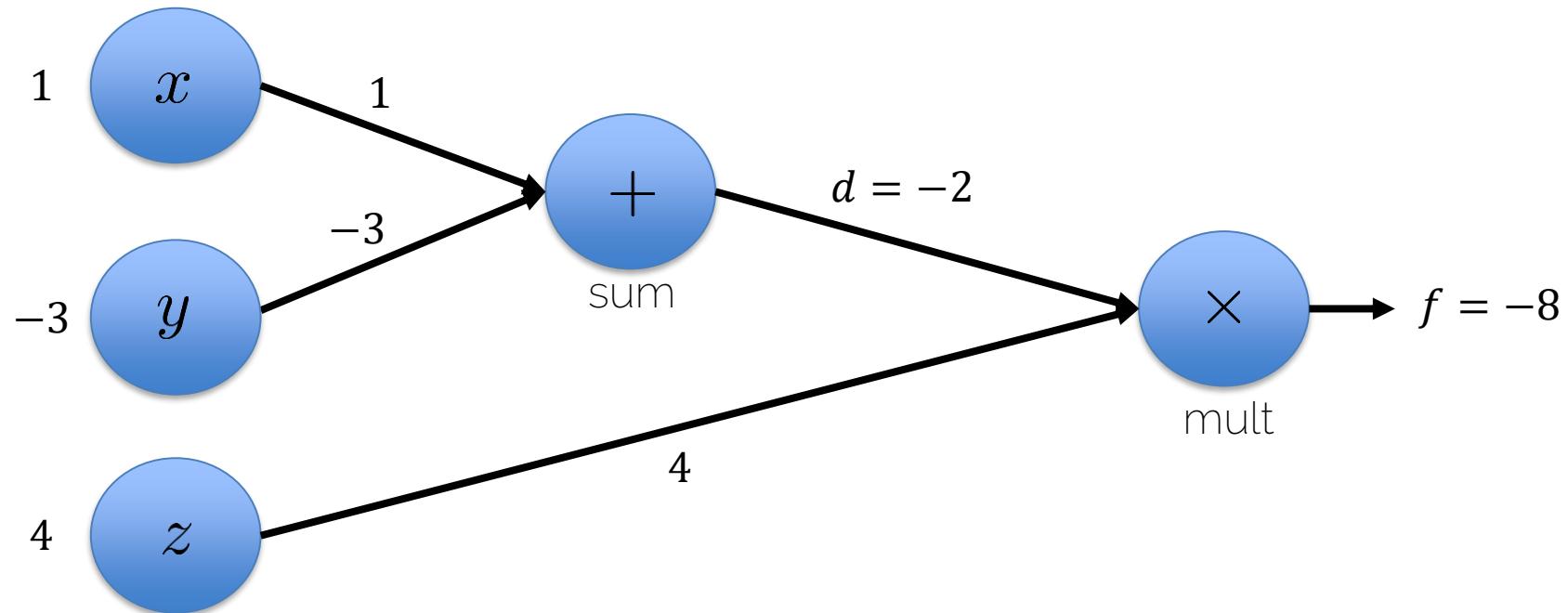
The Flow of Gradients



Backpropagation

Backprop: Forward Pass

- $f(x, y, z) = (x + y) \cdot z$ Initialization $x = 1, y = -3, z = 4$



Backprop: Backward Pass

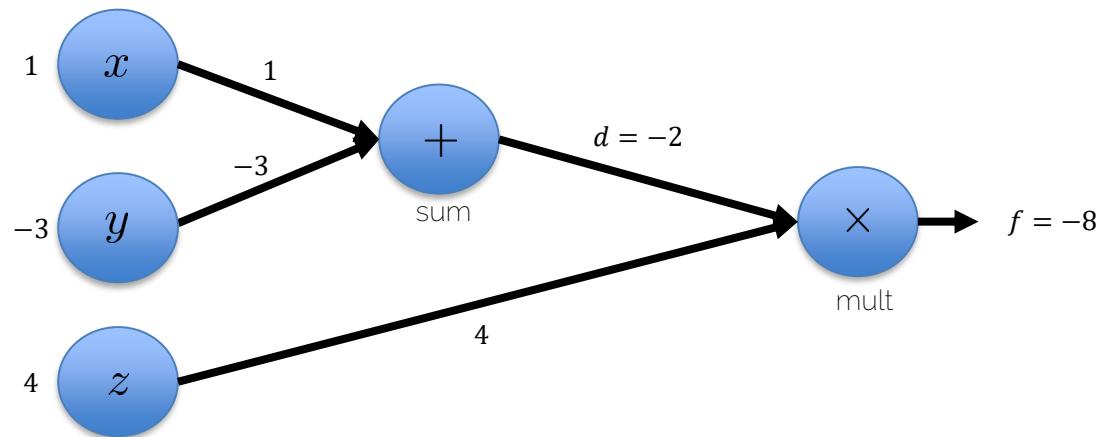
$$f(x, y, z) = (x + y) \cdot z$$

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

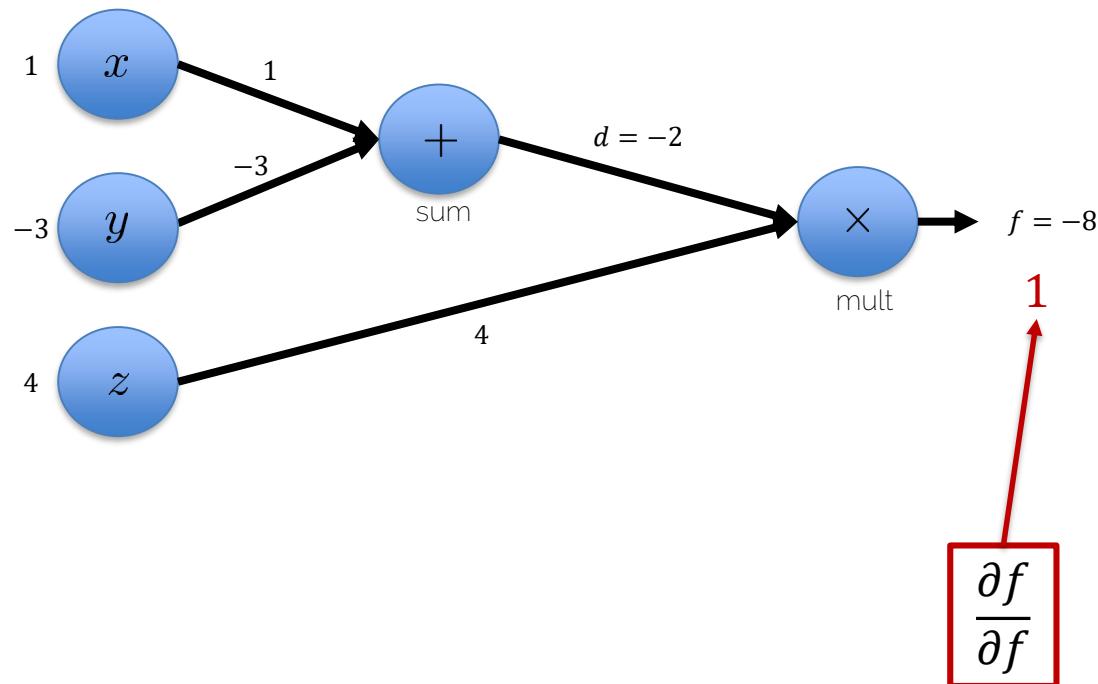
$$f(x, y, z) = (x + y) \cdot z$$

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

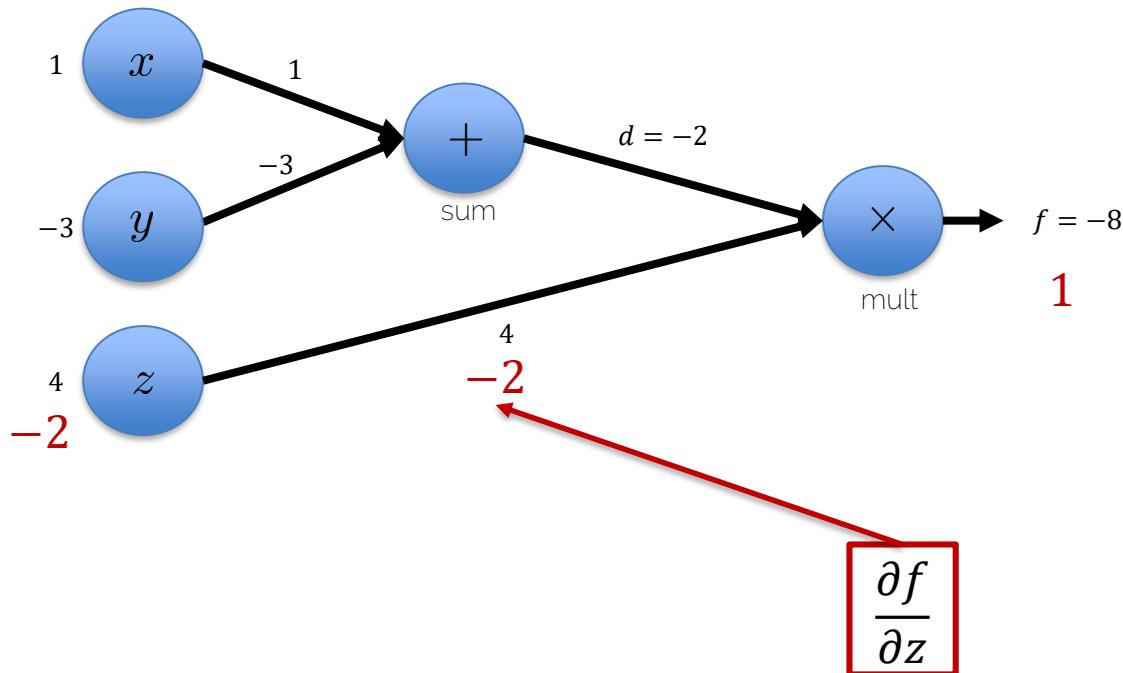
$$f(x, y, z) = (x + y) \cdot z$$

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

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

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

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



Backprop: Backward Pass

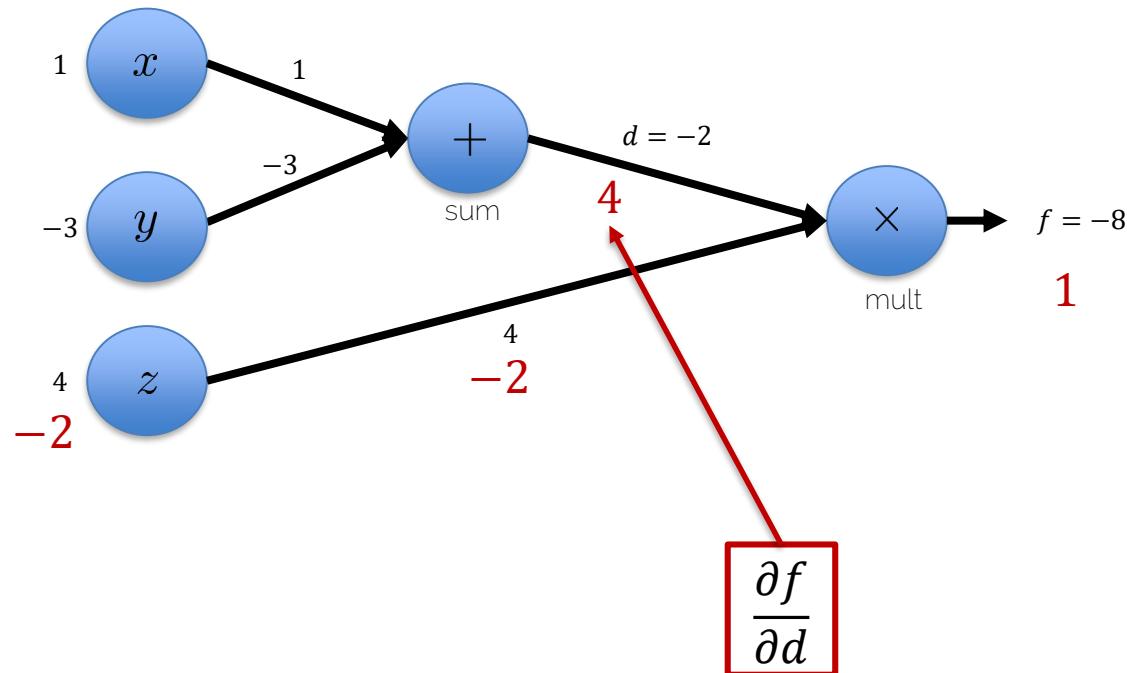
$$f(x, y, z) = (x + y) \cdot z$$

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

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

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

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



Backprop: Backward Pass

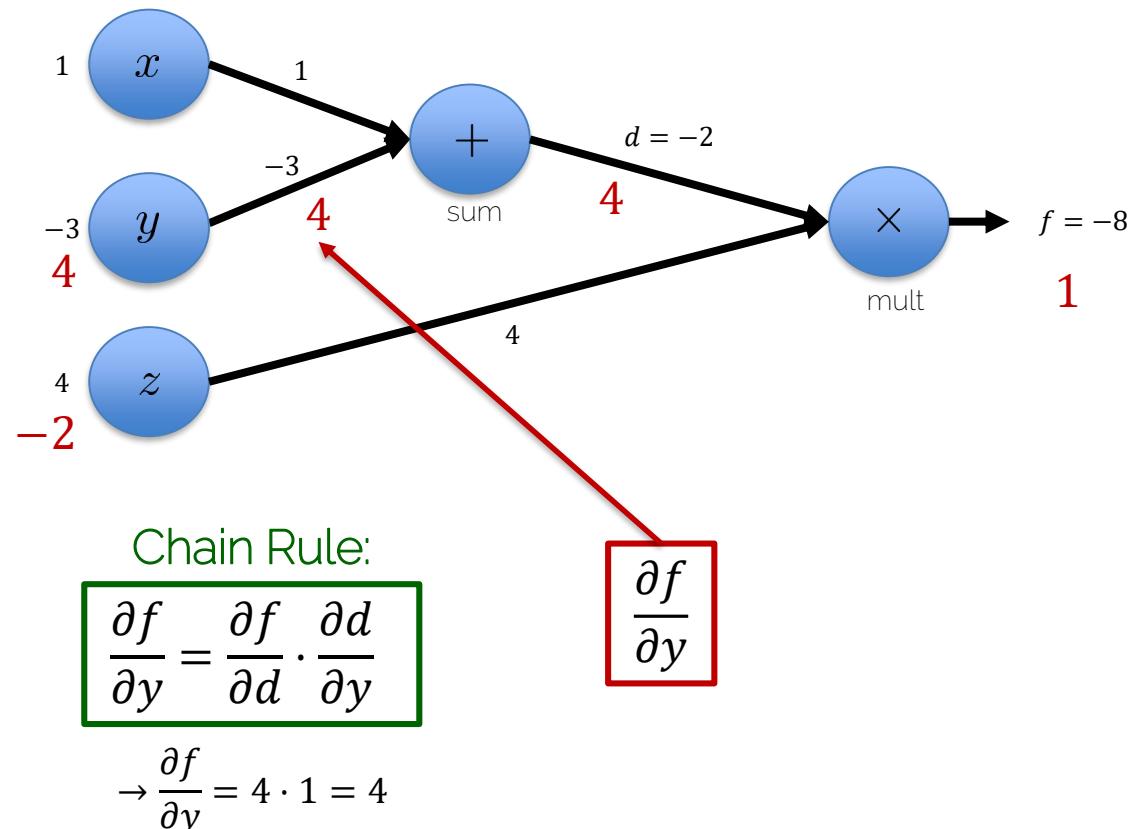
$$f(x, y, z) = (x + y) \cdot z$$

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

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

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

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



Backprop: Backward Pass

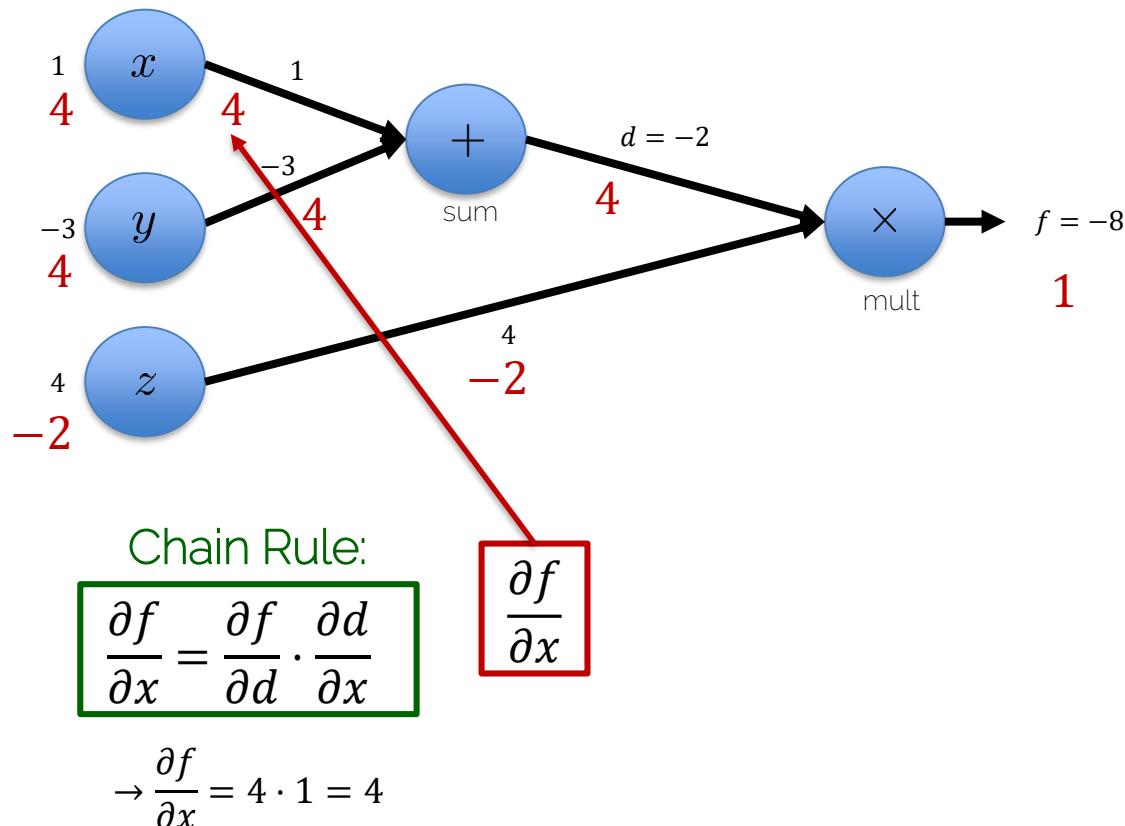
$$f(x, y, z) = (x + y) \cdot z$$

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

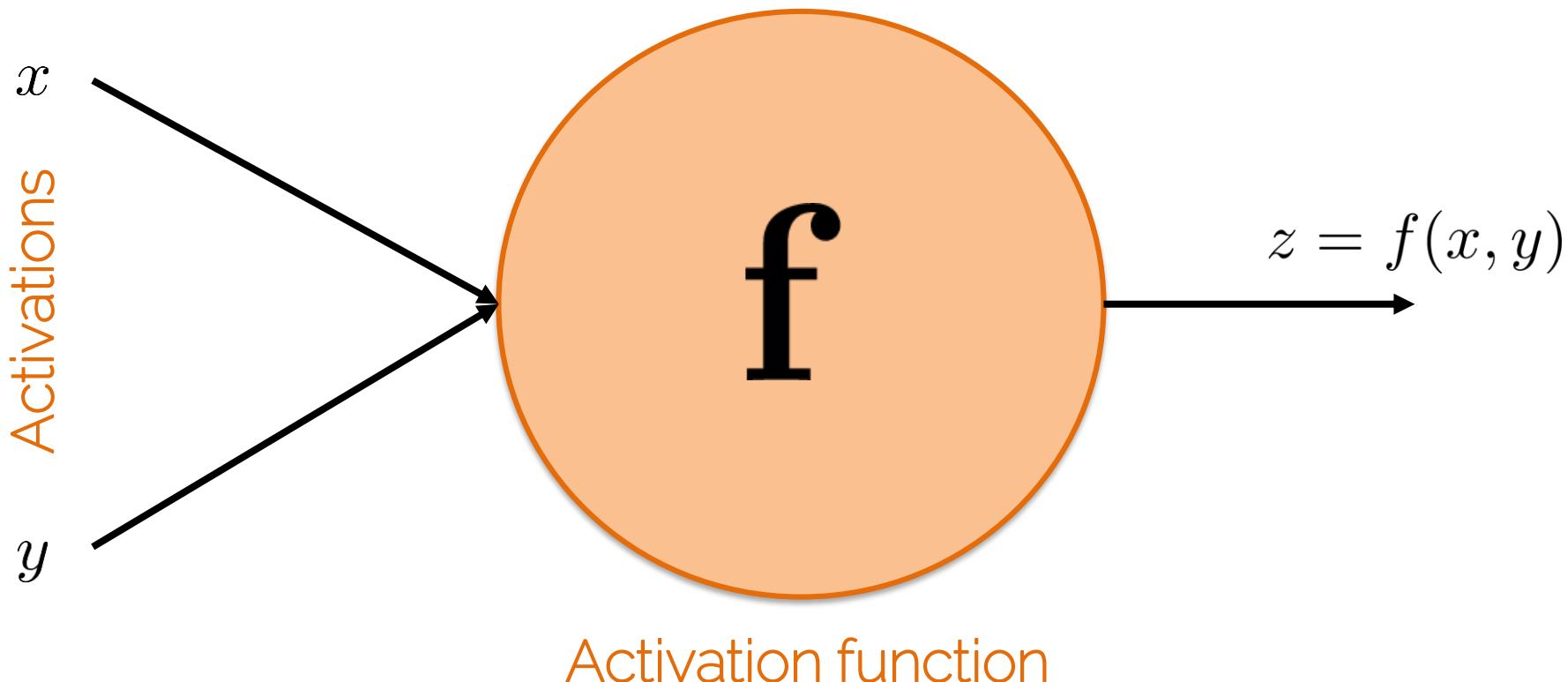
$$d = x + y \quad \boxed{\frac{\partial d}{\partial x} = 1} \quad \frac{\partial d}{\partial y} = 1$$

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

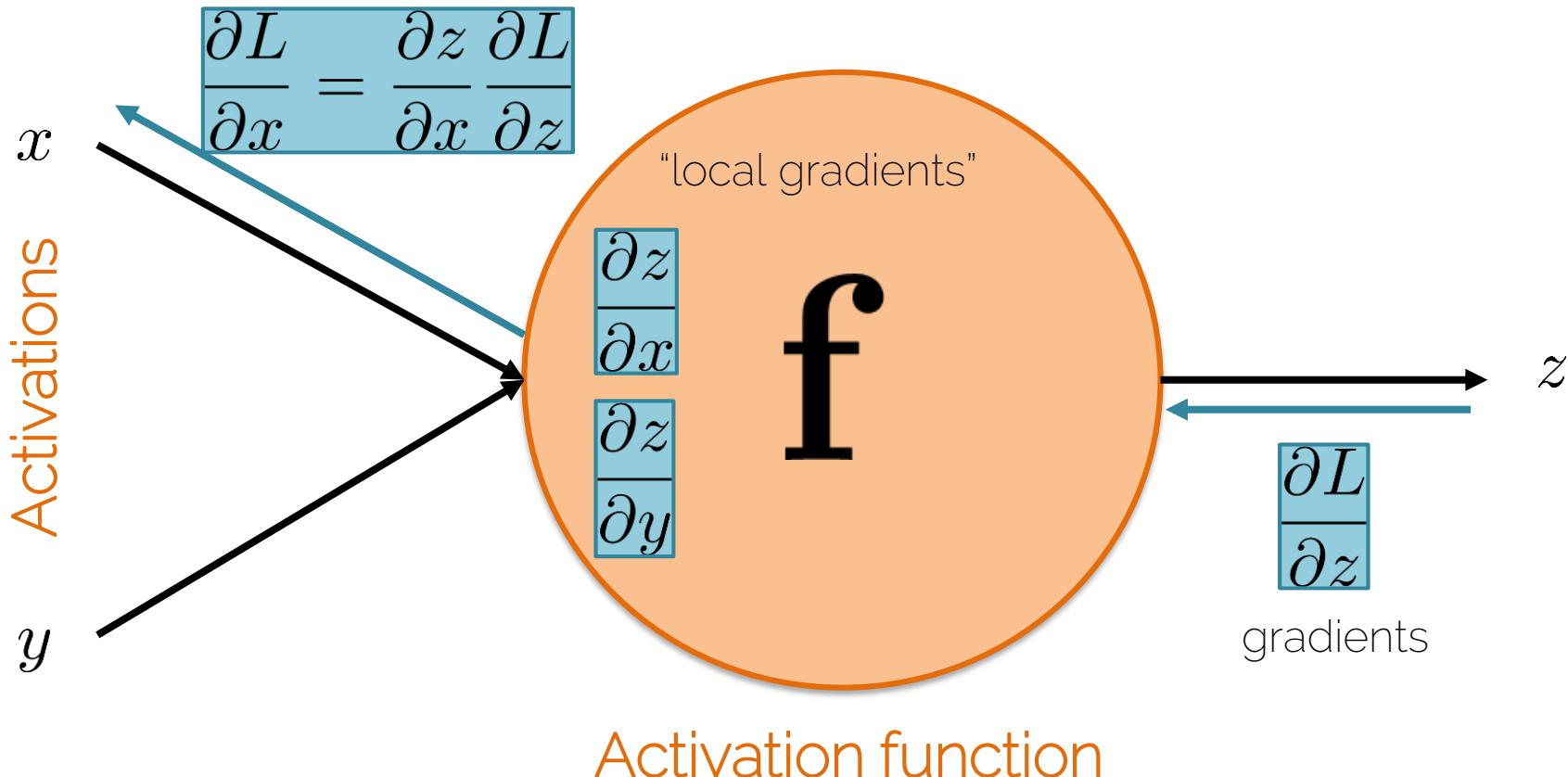
What is $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$?



The Flow of Gradients

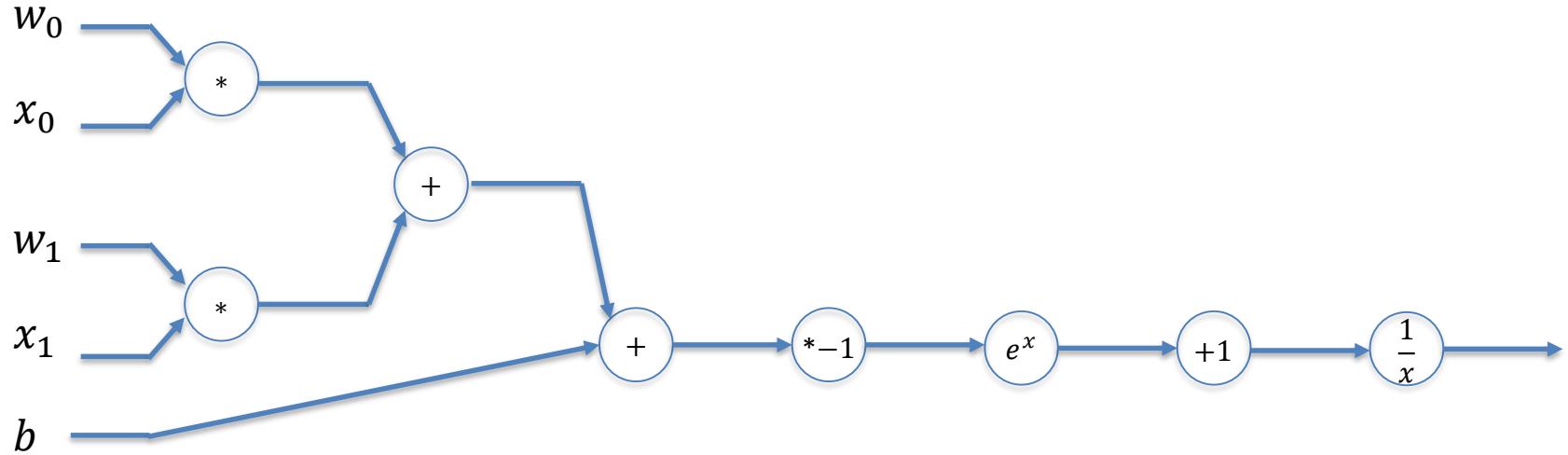


The Flow of Gradients



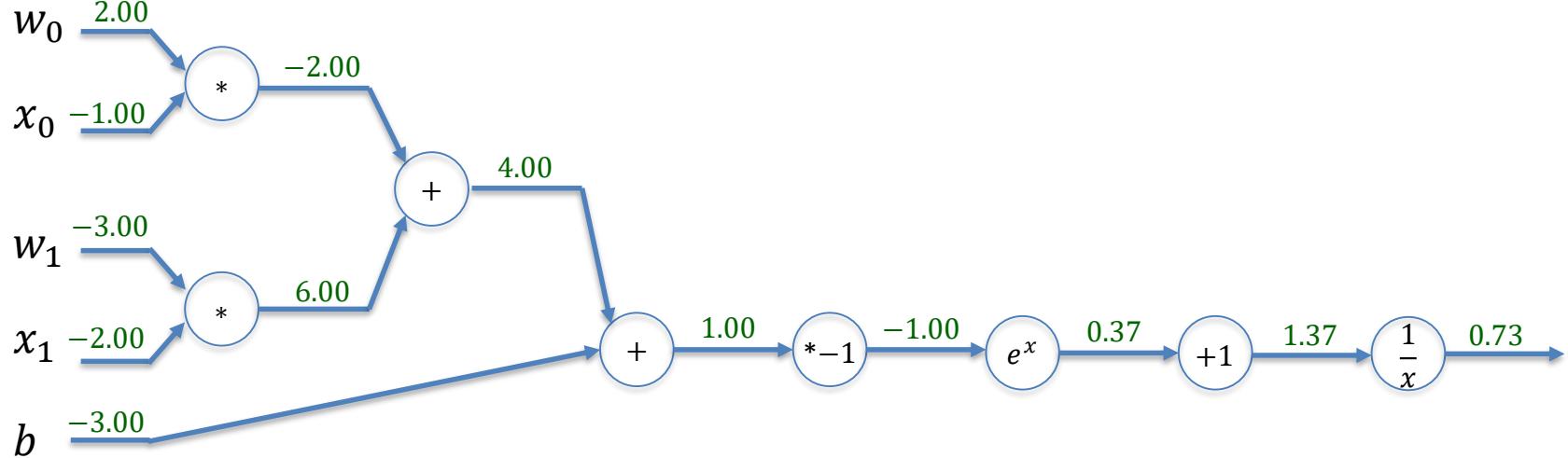
Backprop

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

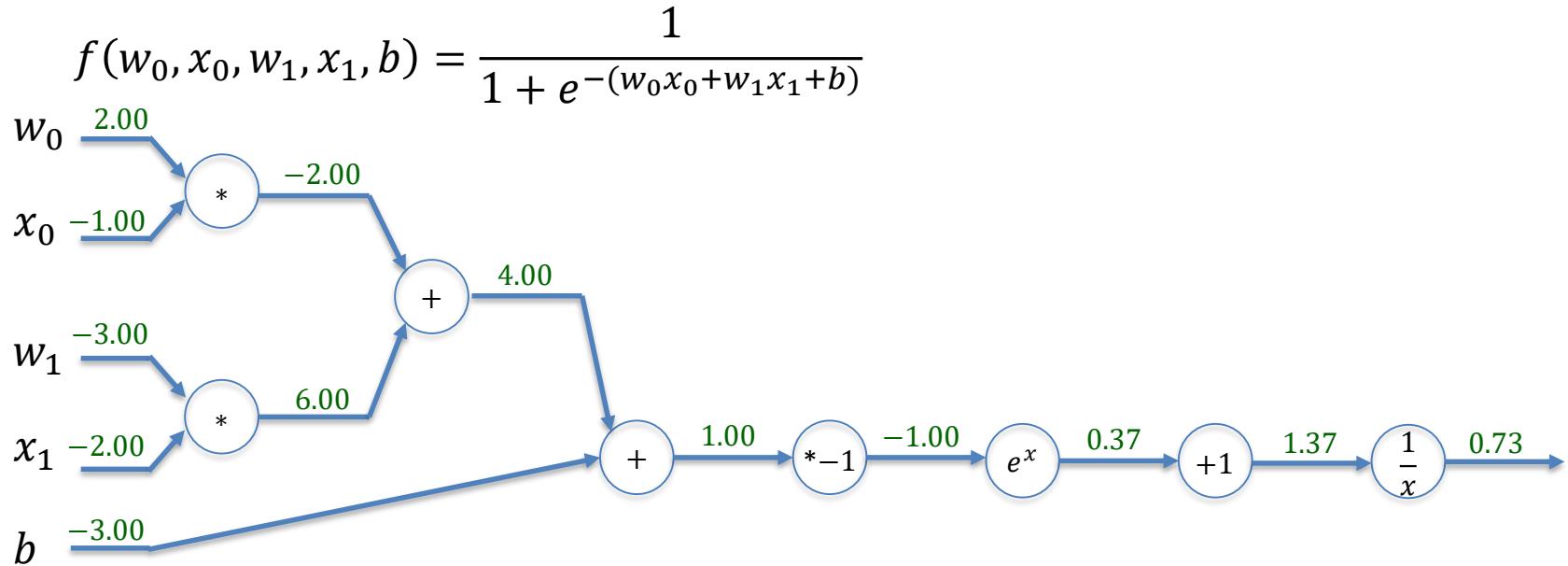


Backprop

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



Backprop



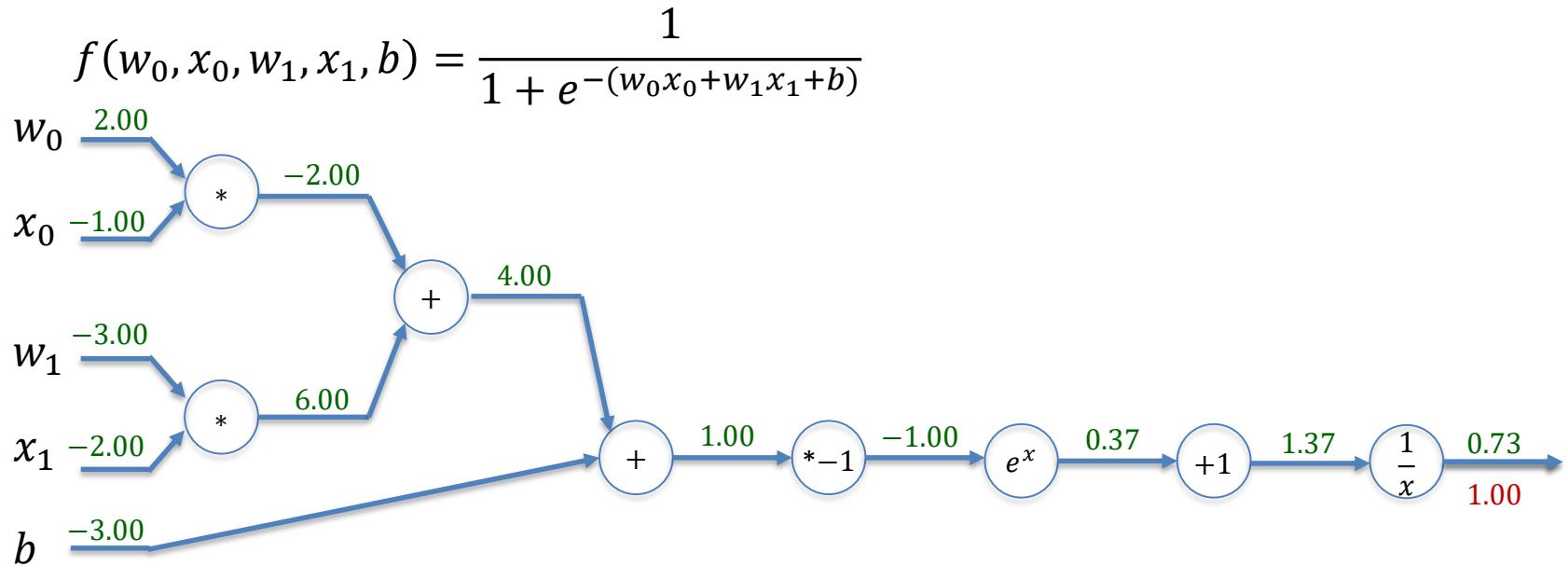
$$f(x) = e^x \quad \rightarrow \quad \frac{\partial f}{\partial x} = e^x$$

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

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

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

Backprop



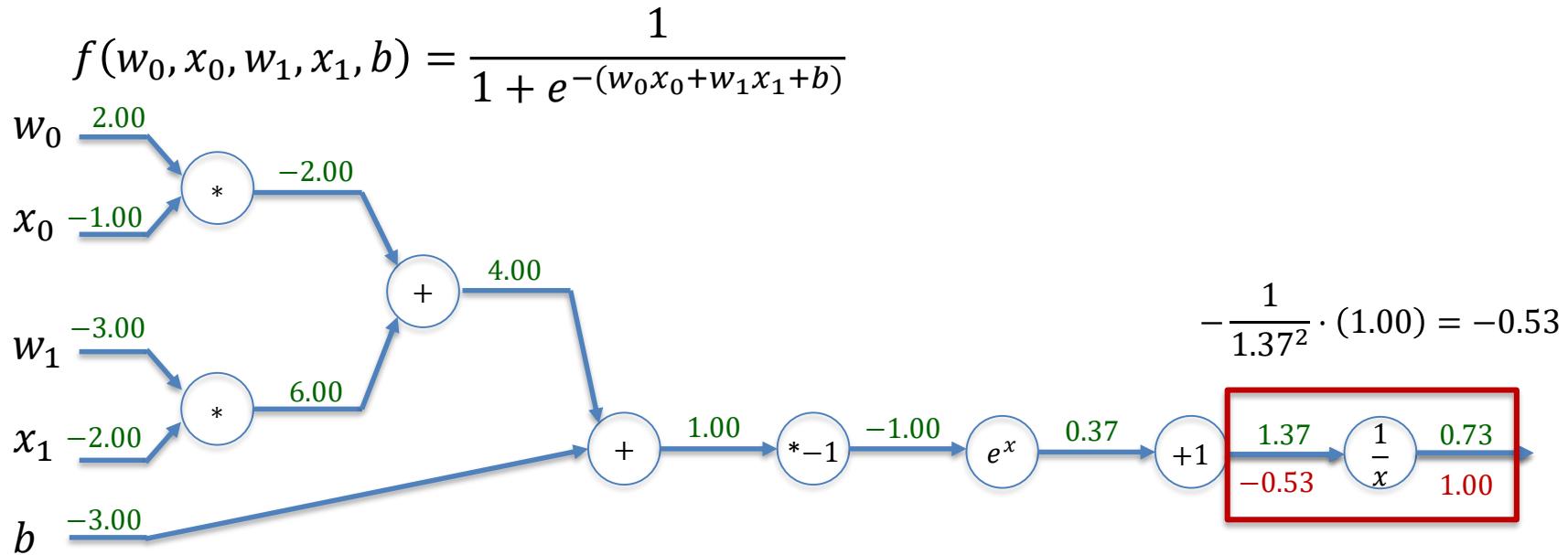
$$f(x) = e^x \quad \rightarrow \quad \frac{\partial f}{\partial x} = e^x$$

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

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

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

Backprop



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

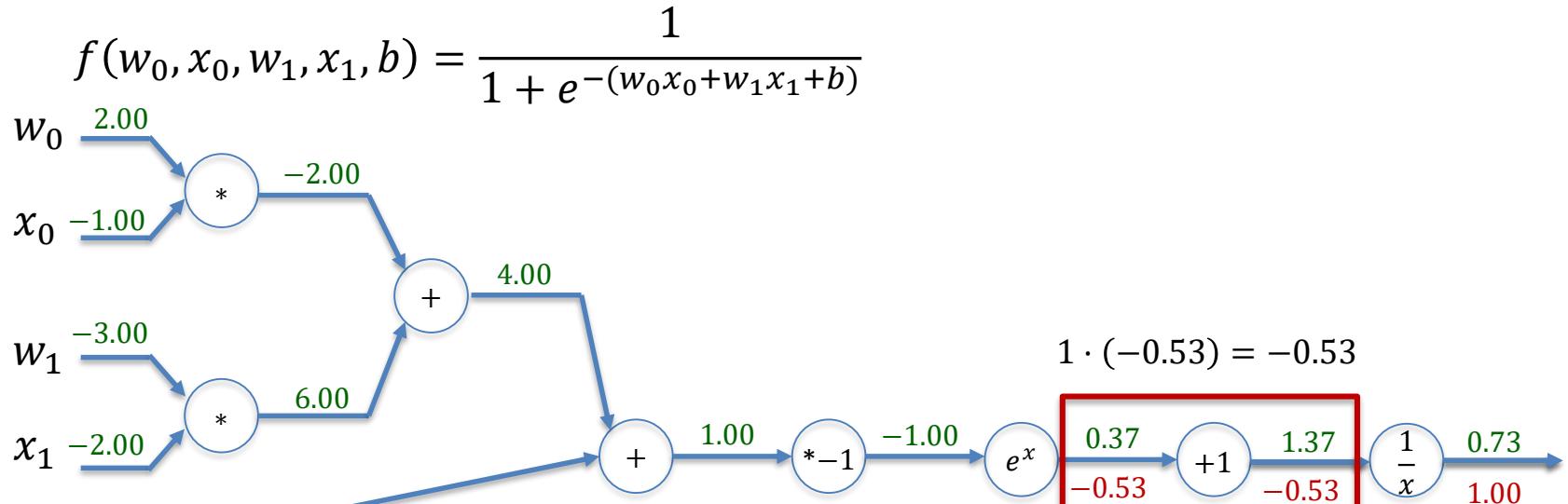
$$f_a(x) = ax \rightarrow \frac{\partial f_a}{\partial x} = a$$

Prof. Leal-Taixé and Prof. Niessner

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

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

Backprop



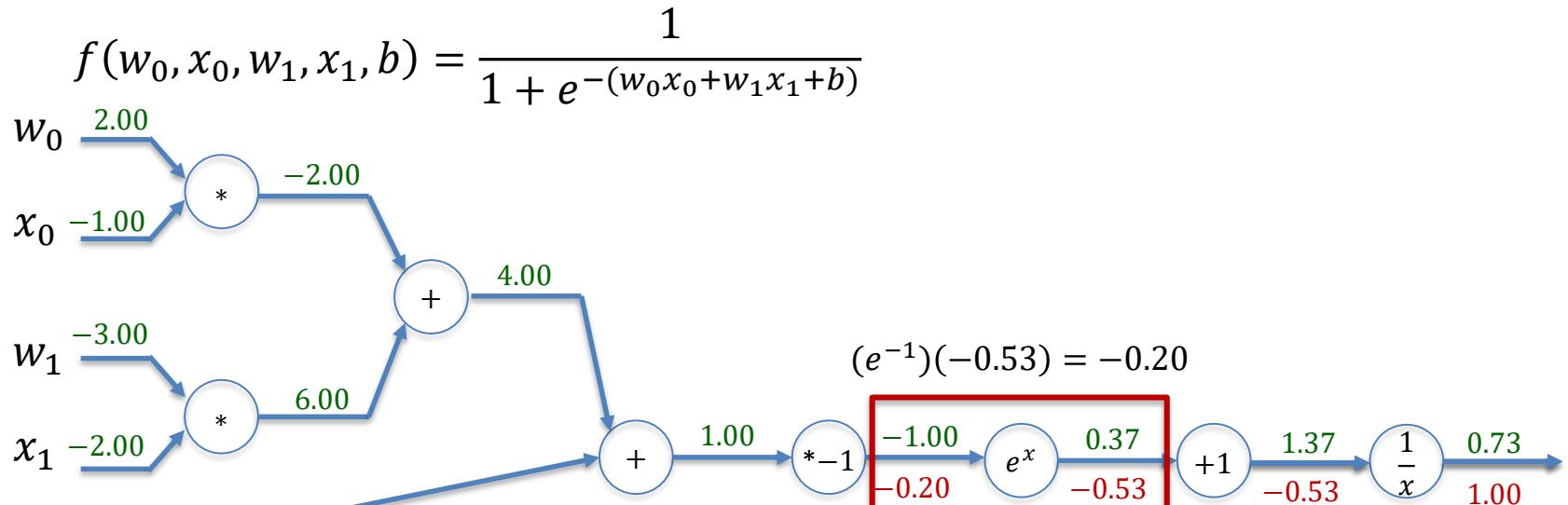
$$f(x) = e^x \quad \rightarrow \quad \frac{\partial f}{\partial x} = e^x$$

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

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

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

Backprop



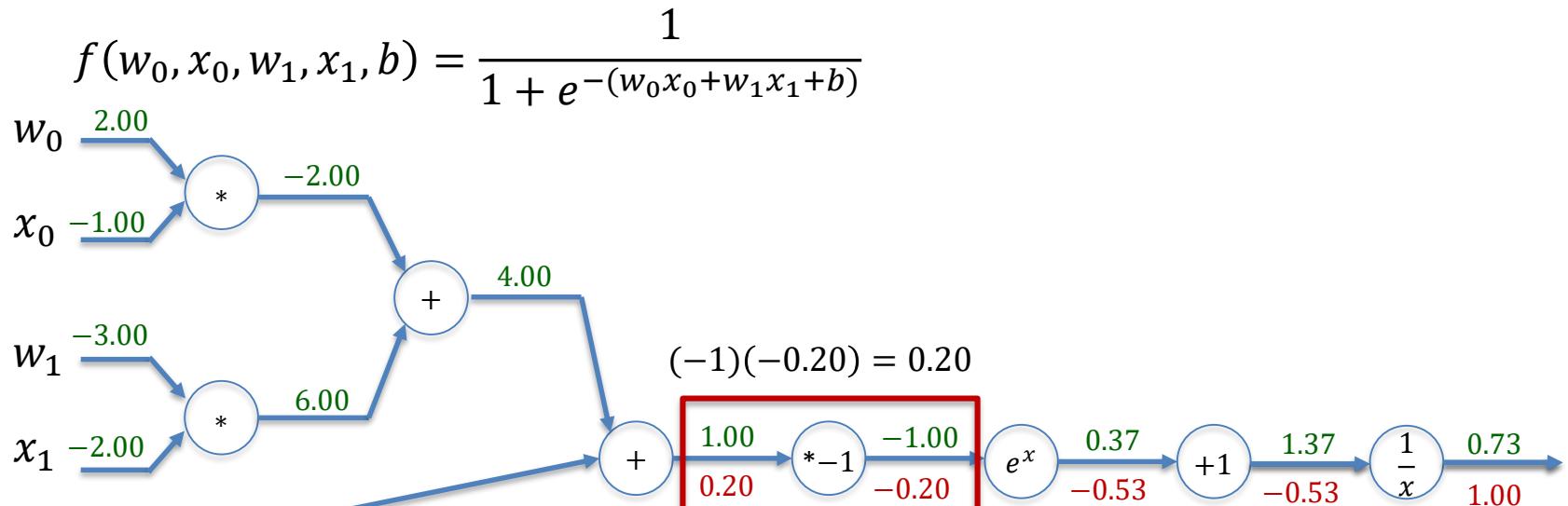
$$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$$

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

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

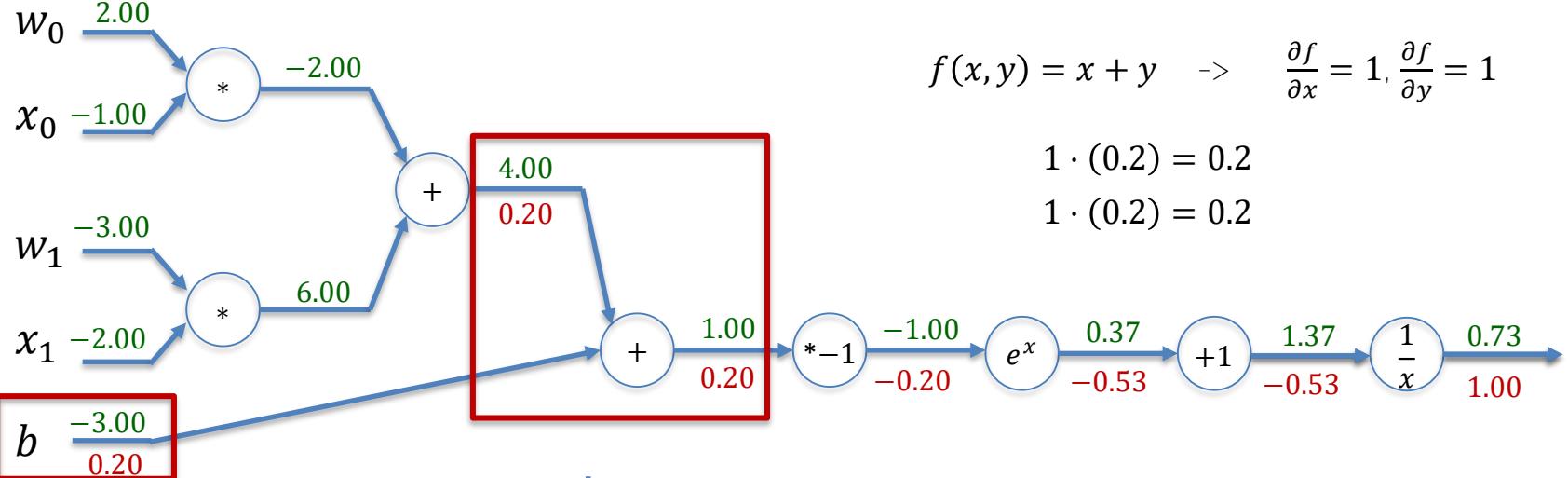
$$f_c(x) = c + x \rightarrow \frac{\partial f_c}{\partial x} = 1$$

Backprop



Backprop

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

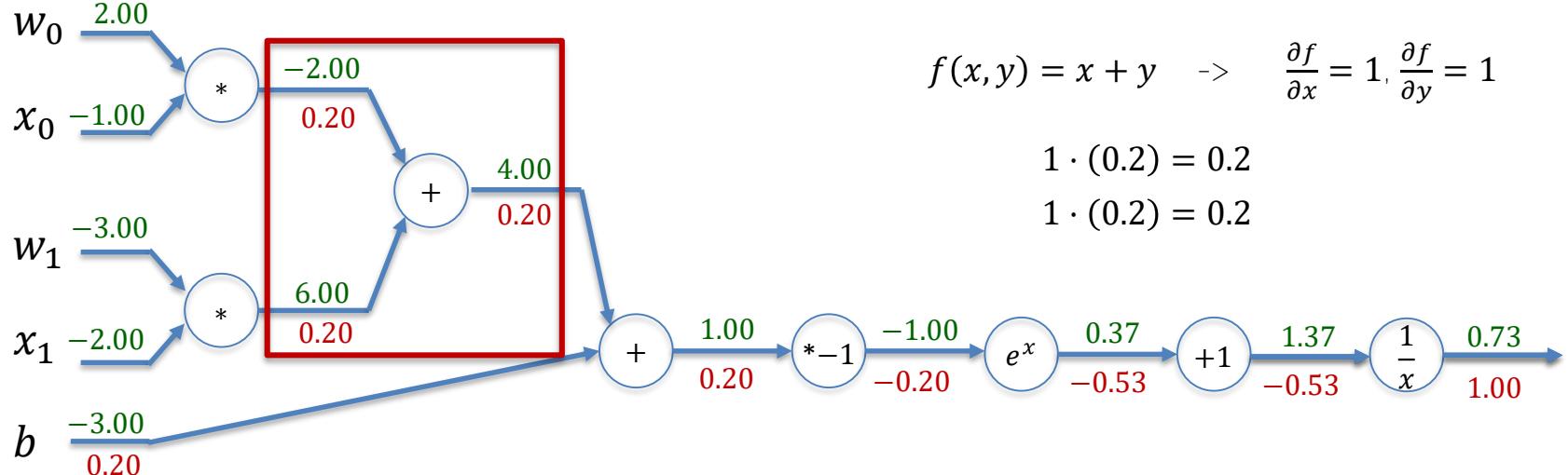


$$f(x) = e^x \rightarrow \frac{\partial f}{\partial x} = e^x$$

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

Backprop

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



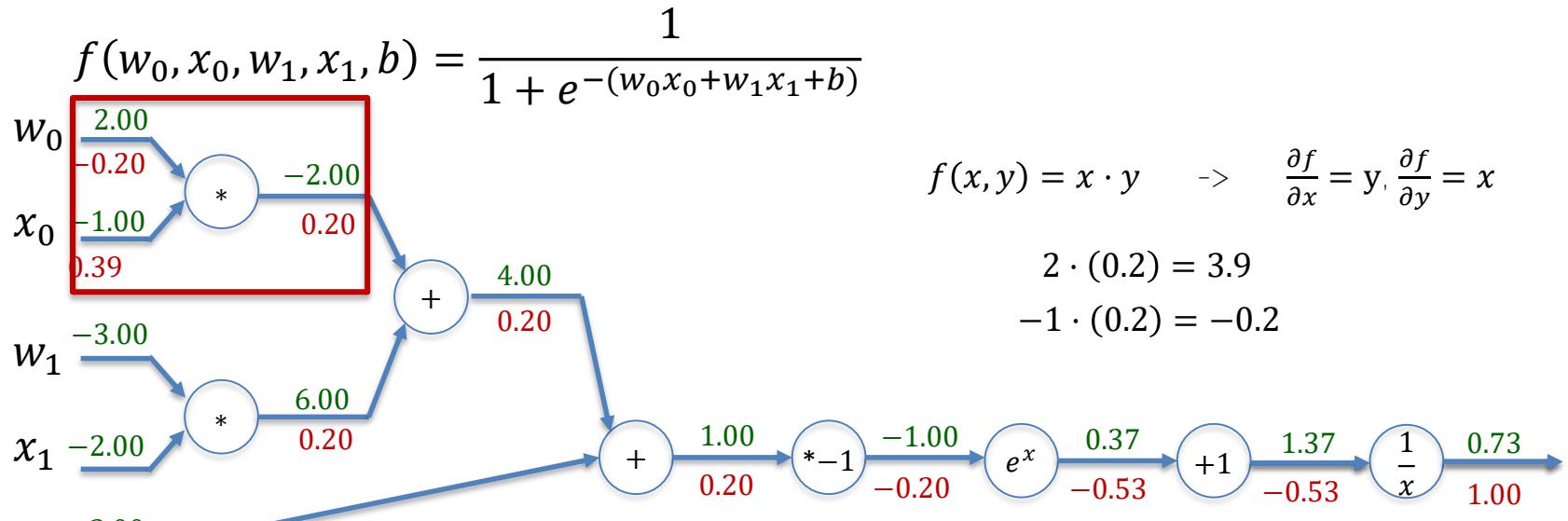
$$f(x) = e^x \quad \rightarrow \quad \frac{\partial f}{\partial x} = e^x$$

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

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

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

Backprop



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

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

Prof. Leal-Taixé and Prof. Niessner

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

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

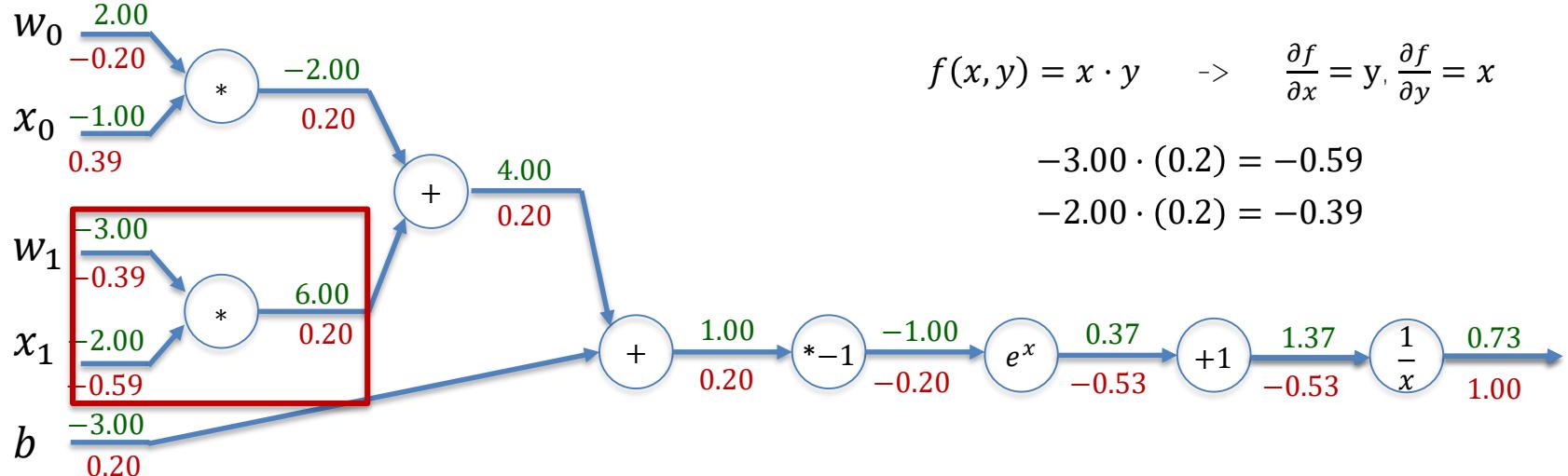
$$-1 \cdot (0.2) = -0.2$$

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

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

Backprop

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



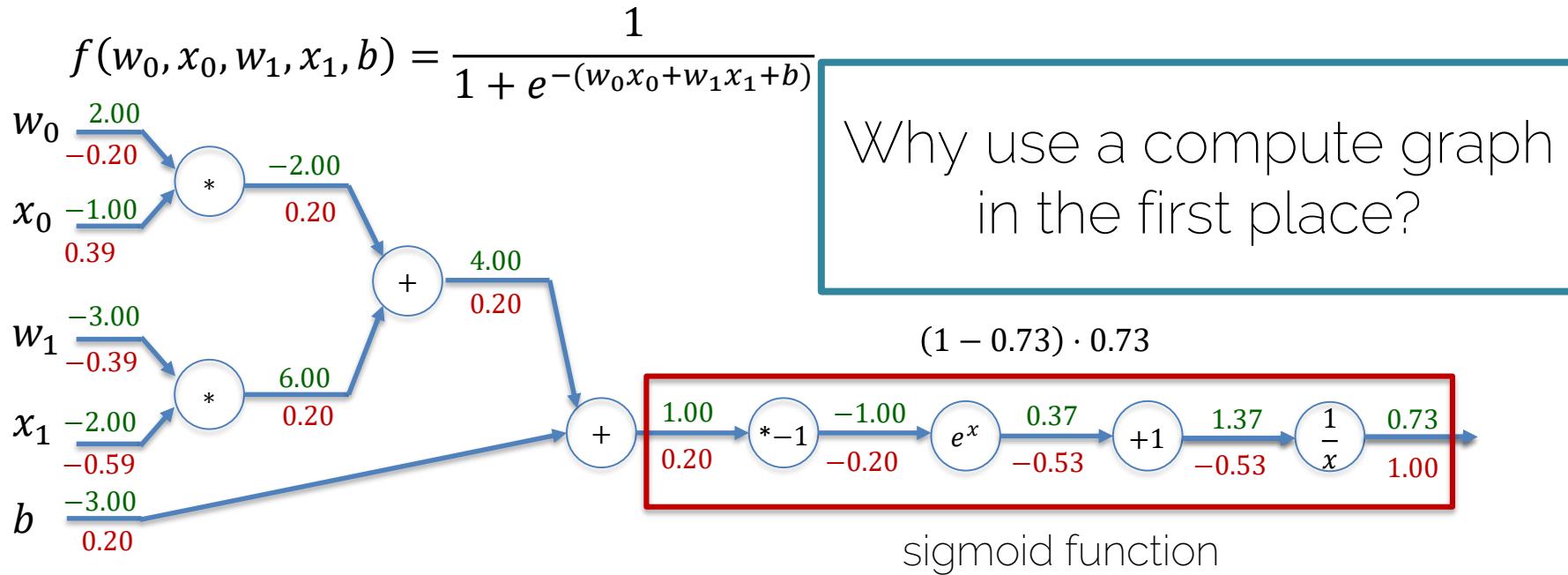
$$f(x) = e^x \quad \rightarrow \quad \frac{\partial f}{\partial x} = e^x$$

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

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

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

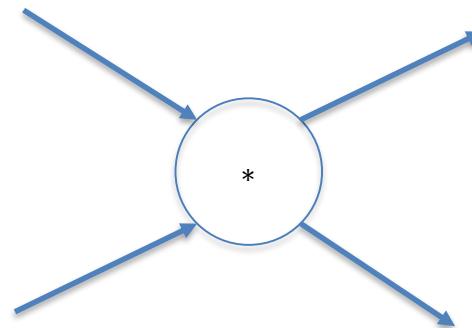
Backprop



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

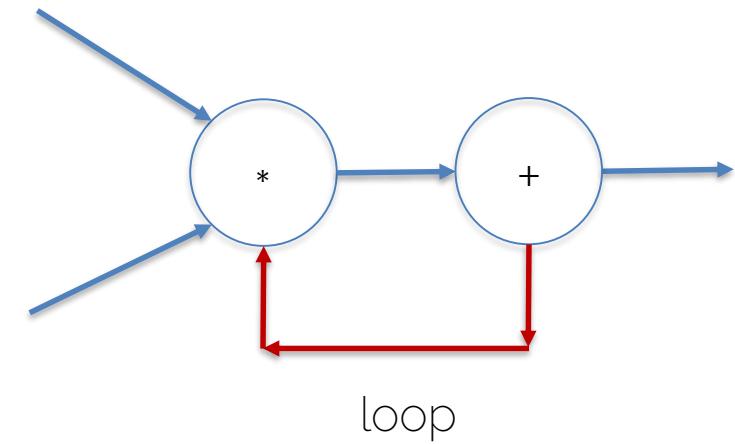
Backpropagation

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

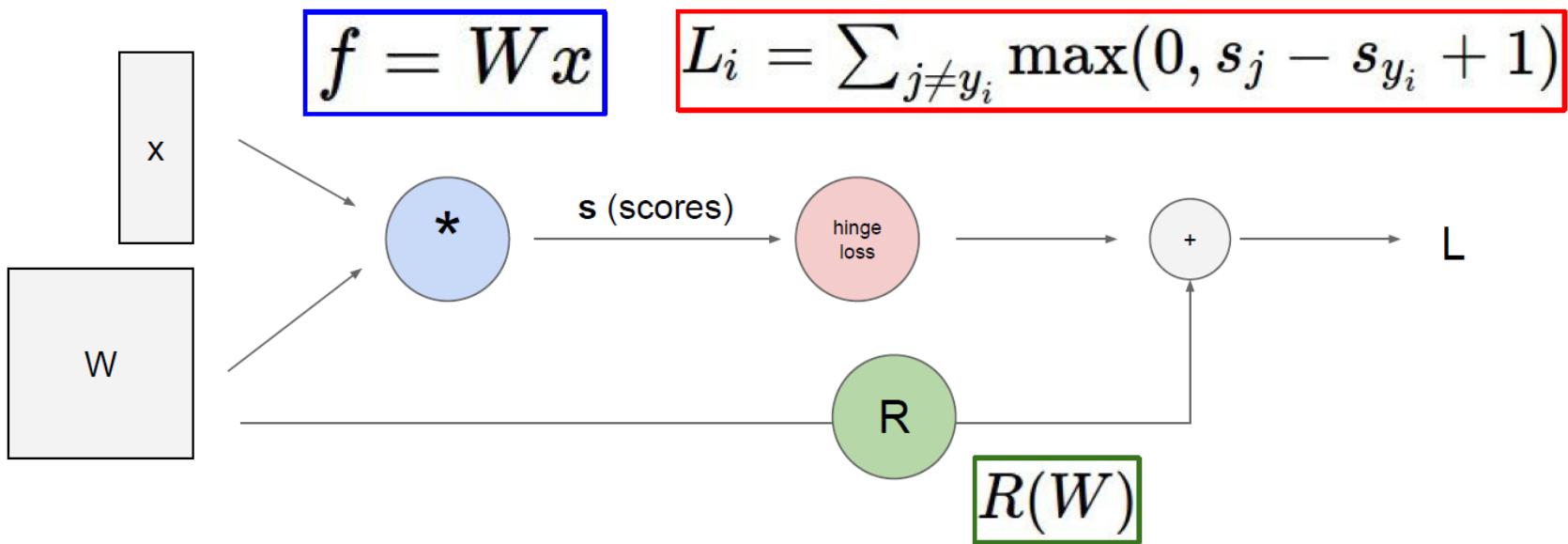


Backpropagation

What happens if there are loops in the graph?



Computational Graph



Combining nodes:
Linear activation node + hinge loss + regularization

Computational Graph: Logistic Regression

- Loss function

$$\mathcal{L}(\hat{y}_i, y_i) = y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)$$

- Cost function

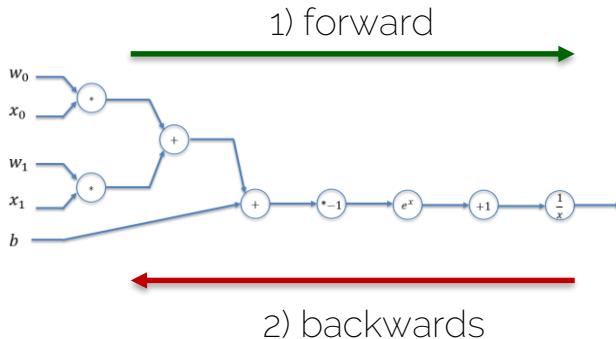
More Next Lecture!

$$C(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^n y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)$$

Minimization

The diagram shows two orange arrows originating from the term $y_i \log \hat{y}_i$ in the cost function equation. One arrow points downwards to the word "Minimization", and another arrow points to the right to the equation $\hat{y}_i = \sigma(\mathbf{x}_i \boldsymbol{\theta})$.

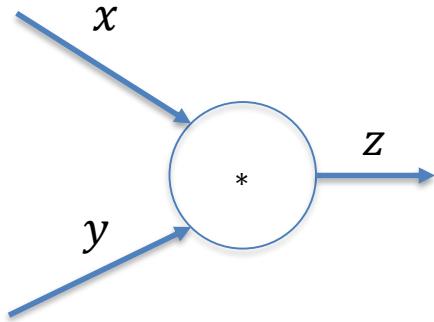
Implementation of Compute Graph



```
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



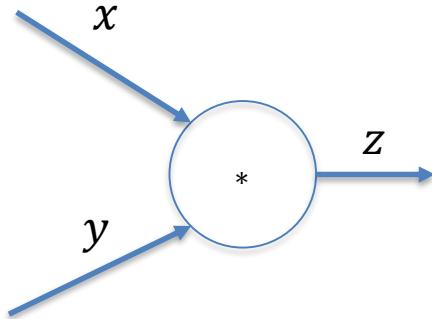
all values are scalars

```
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	Reshape.lua Better __tostring__ and cleans formatting	a year ago
MSECriterion.lua	Select.lua Adds negative dim arguments	11 months ago
MV.lua	SelectTable.lua allow SelectTable to accept input that contains tables of things that...	2 months ago
MapTable.lua	Sequential.lua Improve error handling	a year ago
MarginCriterion.lua	Sigmoid.lua Add THNN conversion of {RReLU, Sigmoid, SmoothL1Criterion, SoftMax, So...	a year ago
MarginRankingCriterion.lua	SmoothL1Criterion.lua Add THNN conversion of {RReLU, Sigmoid, SmoothL1Criterion, SoftMax, So...	a year ago
MaskedSelect.lua	SoftMarginCriterion.lua SoftMarginCriterion	a year ago
Max.lua	SoftMax.lua Add THNN conversion of {RReLU, Sigmoid, SmoothL1Criterion, SoftMax, So...	a year ago
Maxout.lua	SoftMin.lua nn.clearState	a year ago
Mean.lua	SoftPlus.lua Add THNN conversion of {RReLU, Sigmoid, SmoothL1Criterion, SoftMax, So...	a year ago
Min.lua	SoftShrink.lua Add THNN conversion of {softShrink, Sqrt, Square, Tanh, Threshold}	a year ago
MixtureTable.lua	SoftSign.lua nn.clearState	a year ago
Module.lua	SparseJacobian.lua Fix various unused variables in nn	3 years ago
Mul.lua	SparseLinear.lua Fixing sparse linear race condition	a year ago
MulConstant.lua	SpatialAdaptiveAveragePooling.lua Add SpatialAdaptiveAveragePooling.	4 months ago
MultiCriterion.lua	SpatialAdaptiveMaxPooling.lua Indices for nn.	7 months ago
Prof. Leal-Taixé and Prof. Niessner	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
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
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
```

Init()

$$f(x) = aX$$

Forward()

Backward()

Caffe: Layers (GitHub)

absval_layer.cpp	dismantle layer headers	2 years ago
absval_layer.cu	dismantle layer headers	2 years ago
accuracy_layer.cpp	dismantle layer headers	2 years ago
argmax_layer.cpp	dismantle layer headers	2 years ago
base_conv_layer.cpp	dismantle layer headers	2 years ago
base_data_layer.cpp	dismantle layer headers	2 years ago
base_data_layer.cu	dismantle layer headers	2 years ago
batch_norm_layer.cpp	dismantle layer headers	2 years ago
batch_norm_layer.cu	dismantle layer headers	2 years ago
batch_reindex_layer.cpp	dismantle layer headers	2 years ago
batch_reindex_layer.cu	dismantle layer headers	2 years ago
bnll_layer.cpp	dismantle layer headers	2 years ago
bnll_layer.cu	dismantle layer headers	2 years ago
concat_layer.cpp	dismantle layer headers	2 years ago
concat_layer.cu	dismantle layer headers	2 years ago
contrastive_loss_layer.cpp	dismantle layer headers	2 years ago
contrastive_loss_layer.cu	dismantle layer headers	2 years ago
conv_layer.cpp	dismantle layer headers	2 years ago
conv_layer.cu	dismantle layer headers	2 years ago
cudnn_conv_layer.cpp	dismantle layer headers	2 years ago
cudnn_conv_layer.cu	dismantle layer headers	2 years ago
cudnn_lcn_layer.cpp	dismantle layer headers	2 years ago
cudnn_lcn_layer.cu	dismantle layer headers	2 years ago
cudnn_lrn_layer.cpp	dismantle layer headers	2 years ago
cudnn_lrn_layer.cu	dismantle layer headers	2 years ago
cudnn_pooling_layer.cpp	dismantle layer headers	2 years ago
cudnn_pooling_layer.cu	dismantle layer headers	2 years ago
cudnn_relu_layer.cpp	dismantle layer headers	2 years ago
cudnn_relu_layer.cu	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

```
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
```

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

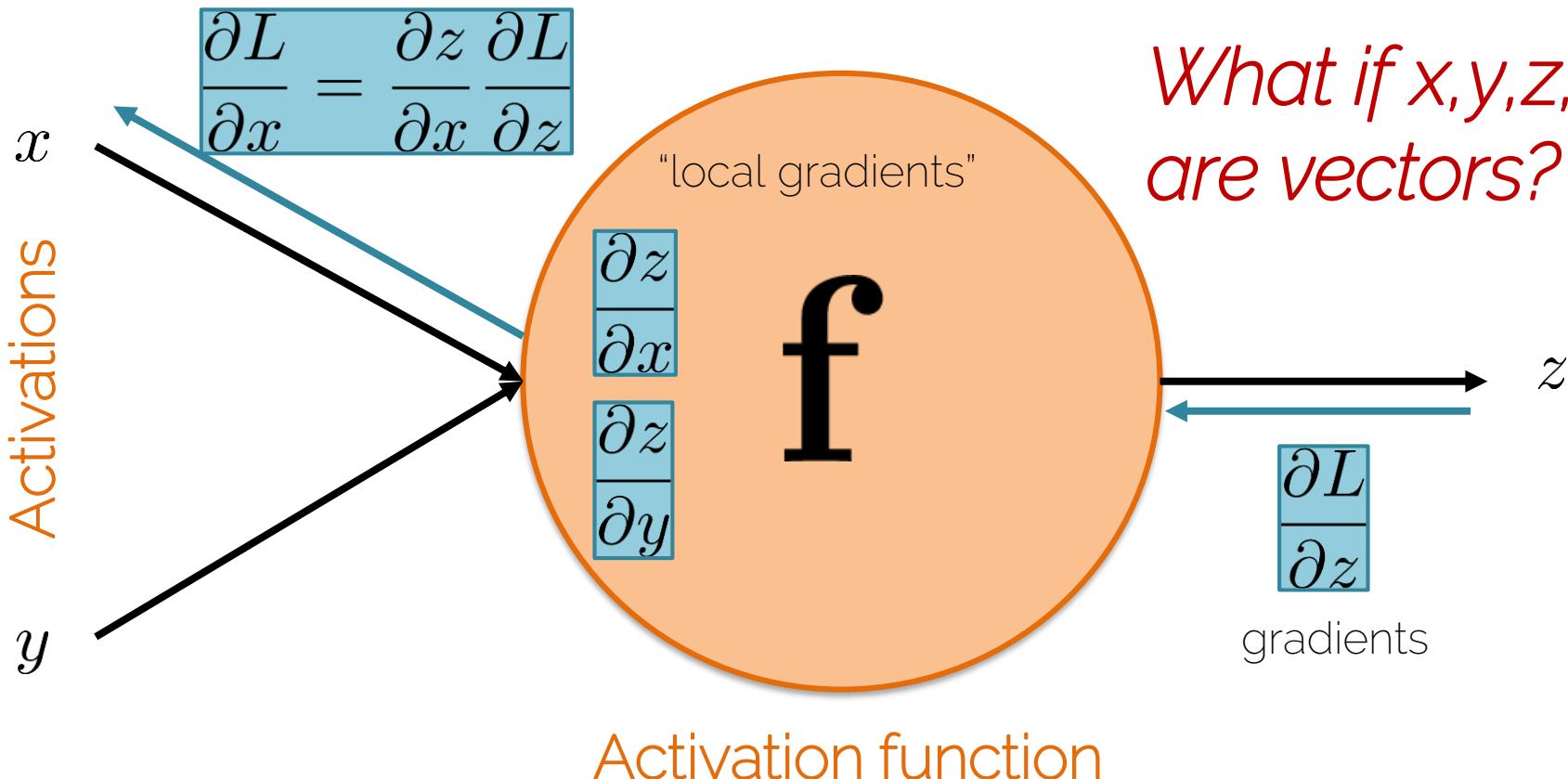
Forward()

```
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    }
```

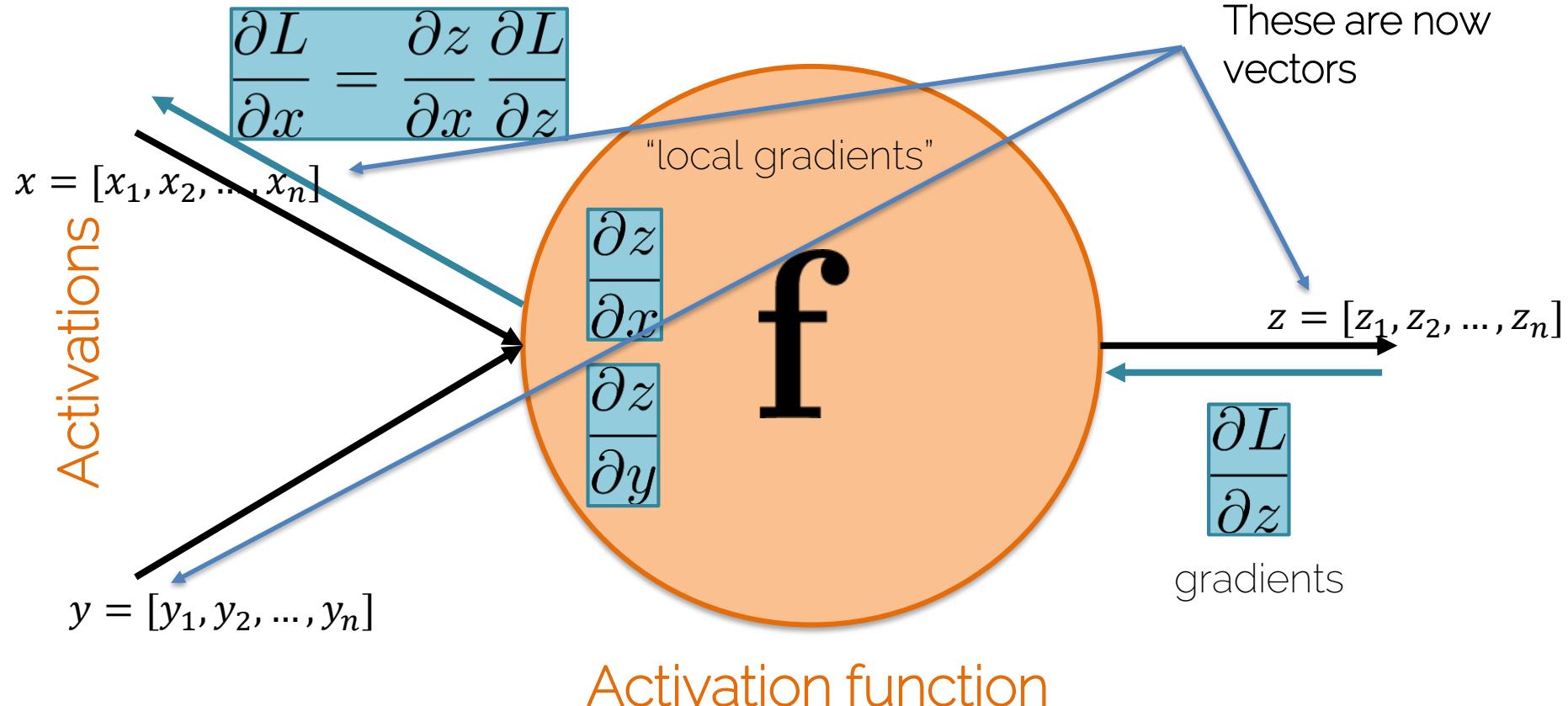
$$\sigma'(x) = (1 - \sigma(x))\sigma(x)$$

Backward()

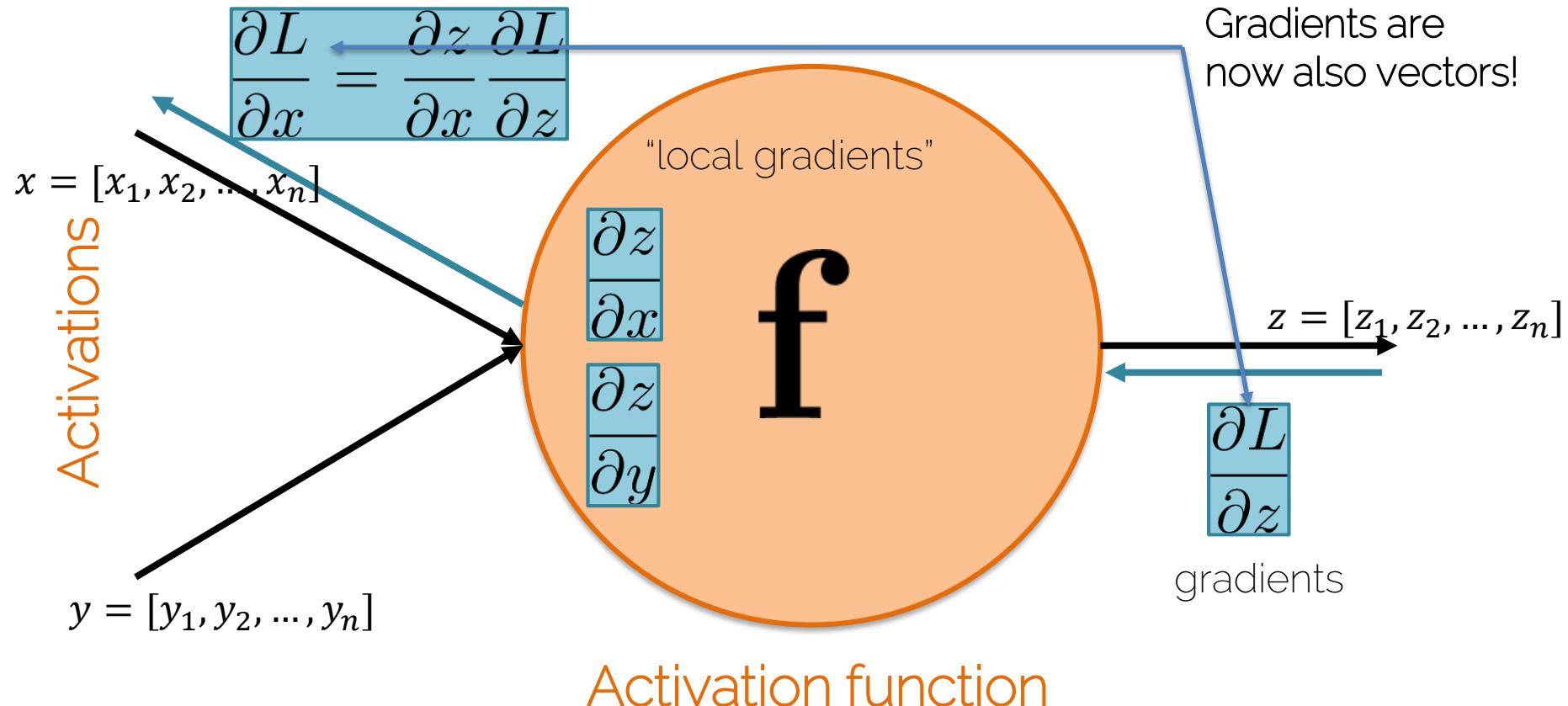
Vectorized Operations



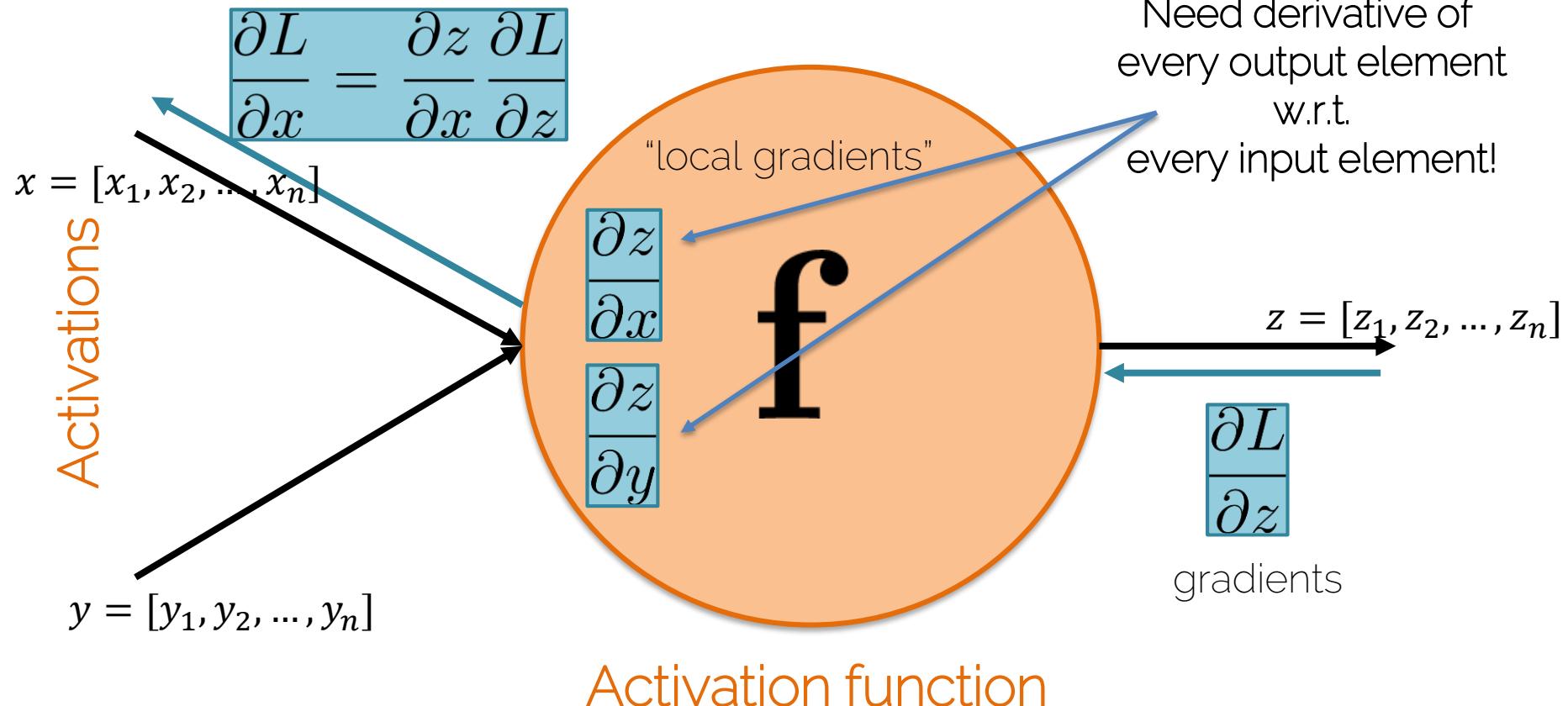
Vectorized Operations



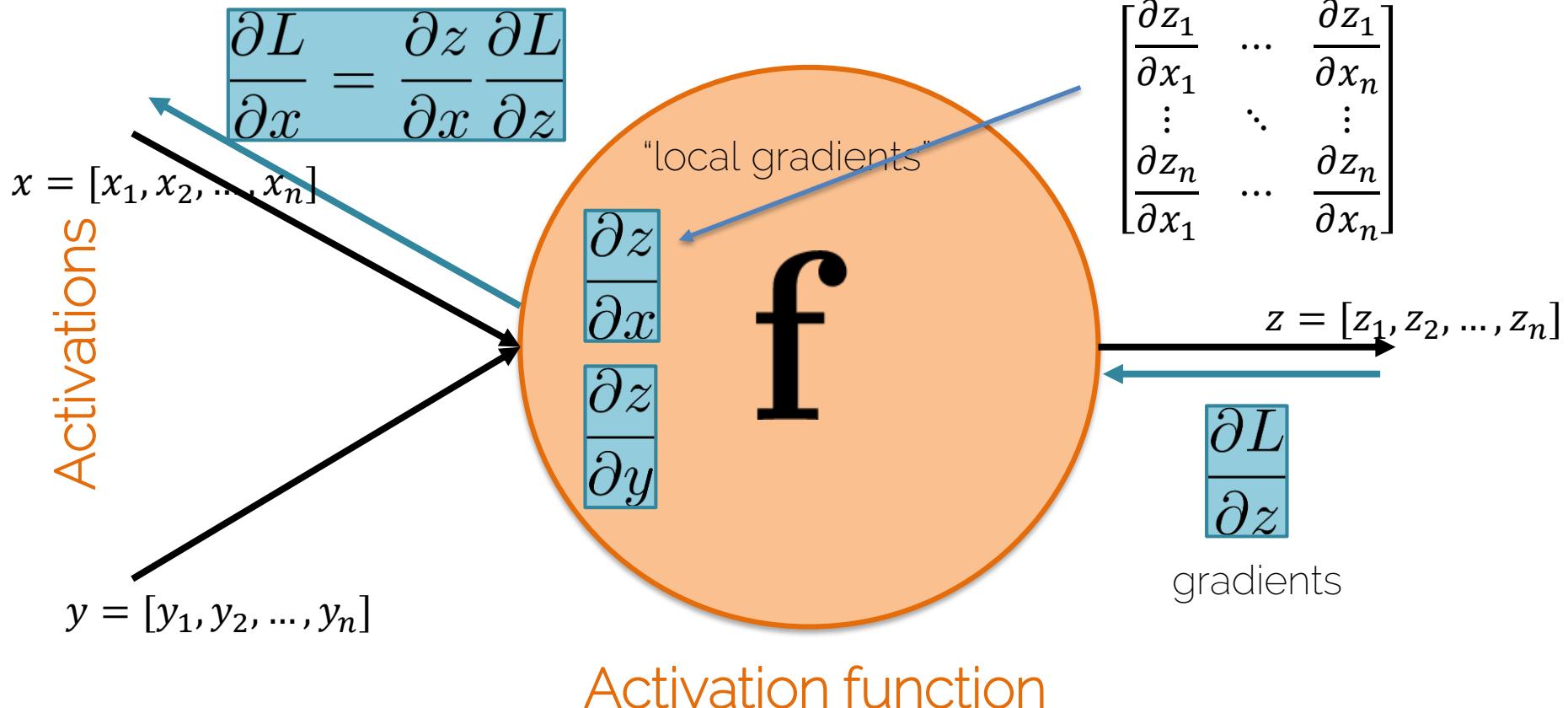
Vectorized Operations



Vectorized Operations



Vectorized Operations

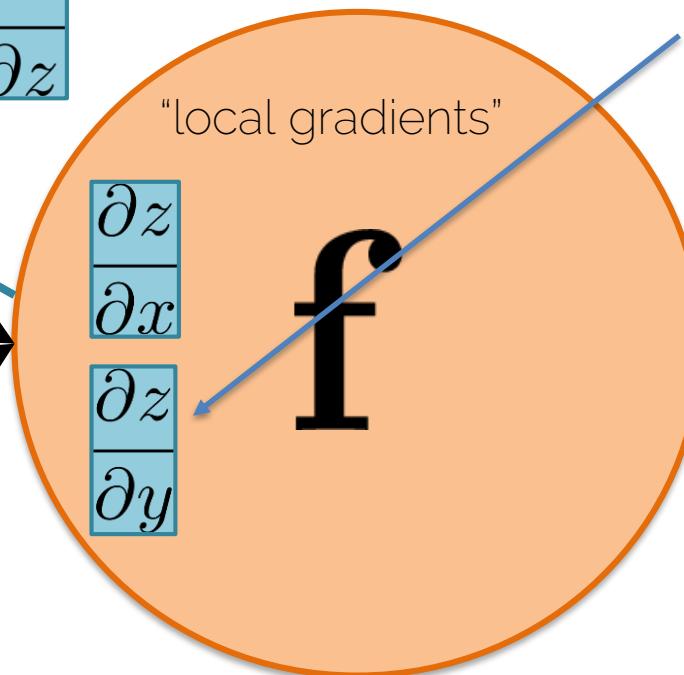


Vectorized Operations

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

Activations

$$x = [x_1, x_2, \dots, x_n]$$
$$y = [y_1, y_2, \dots, y_n]$$



Activation function

Jacobian Matrix:

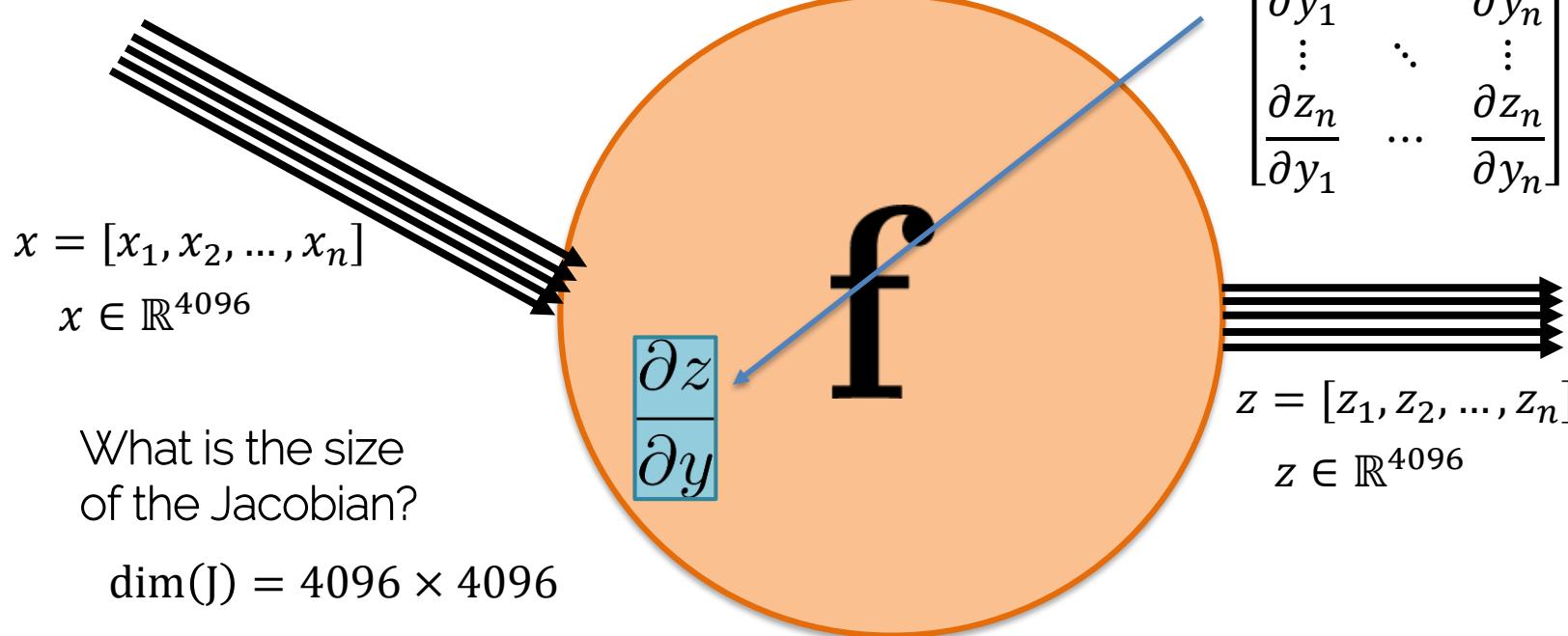
$$\begin{bmatrix} \frac{\partial z_1}{\partial y_1} & \dots & \frac{\partial z_1}{\partial y_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_n}{\partial y_1} & \dots & \frac{\partial z_n}{\partial y_n} \end{bmatrix}$$

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

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

gradients

Vectorized Operations



Vectorized Operations

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 efficient is that:

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

Typically, networks are run in batches:

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

How to handle this?

Neural Networks

Neural Networks

- Are compute graphs
- Goal: for given train set, find optimal weights
- Optimization using gradient-based solvers
 - Many options (more in the next lectures)
- Gradients are computed via backpropagation
 - Nice because can easily modularize complex functions

Administrative Things

- This week: tutorial on Nov 13th
 - No tutorial
- Next Lecture on Nov 15th
 - Optimization and Regularization
 - More on neural networks ☺
- Next week: first exercise
 - Nov 20th --> start of exercise 1

Administrative Things

Bestehender Termin

Zeit: Do, 22.11.2018 18:00-20:00

Ort: 00.02.001, MI HS 1, Friedrich L. Bauer Hörsaal (5602.EG.001)

wurde verschoben auf

Zeit: Do, 22.11.2018 18:00-20:00

Ort: MW 2001, Rudolf-Diesel-Hörsaal (5510.02.001)

Termintyp: geplant

durchgeführt von Dr.sc.nat. Monika Hanesch

Es hat sich Folgendes geändert: Ort

See you next week!