

M2177.0043 Introduction to Deep Learning

Lecture 7: Model fitting / Computation graph and Backpropagation¹

Hyun Oh Song¹

¹Dept. of Computer Science and Engineering, Seoul National University

April 7, 2020

¹Many slides and figures adapted Justin Johnson

Last time

- ▶ Score functions
- ▶ Loss functions
- ▶ Regularization

Outline

Model fitting

Computational graphs

Neural networks

Bias-Variance tradeoff²

- ▶ Suppose we have a training set of data and labels
 $D = \{(x_i, y_i)\}_{i=1}^n$
- ▶ Assume there is a function with noise $y = f(x) + \epsilon$, where the ϵ has zero mean and variance σ^2
- ▶ Want to find a function $\hat{f}(x; D)$ that approximates the true function $f(x)$ as well as possible
- ▶ Concretely, we want MSE $(y - \hat{f}(x; D))^2$ to be minimal for the training data x_1, \dots, x_n and for unseen sample x outside the training set.

²https://en.wikipedia.org/wiki/Bias-variance_tradeoff

- The expected error on an unseen sample x is,

$$\begin{aligned}
 \mathbb{E}_D[(y - \hat{f}(x; D))^2] &= \mathbb{E}_D[(f(x) + \epsilon - \hat{f}(x; D))^2] \\
 &= \mathbb{E}_D[(f(x) + \epsilon - \hat{f}(x; D) + \mathbb{E}_{D'}[\hat{f}(x; D')] - \mathbb{E}_{D'}[\hat{f}(x; D')])^2] \\
 &= \mathbb{E}_D[(f(x) - \mathbb{E}_{D'}[\hat{f}(x; D')])^2] + \mathbb{E}_D[\epsilon^2] + \mathbb{E}_D[(\mathbb{E}_{D'}[\hat{f}(x; D')] - \hat{f}(x; D))^2] \\
 &\quad + 2\mathbb{E}_D[(f(x) - \mathbb{E}_{D'}[\hat{f}(x; D')])\epsilon] + 2\mathbb{E}_D[\epsilon(\mathbb{E}_{D'}[\hat{f}(x; D')] - \hat{f}(x; D))] \\
 &\quad + 2\mathbb{E}_D[(\mathbb{E}_{D'}[\hat{f}(x; D')] - \hat{f}(x; D))(f(x) - \mathbb{E}_{D'}[\hat{f}(x; D')])] \\
 &= (f(x) - \mathbb{E}_{D'}[\hat{f}(x; D')])^2 + \text{Var}[\epsilon] + \text{Var}[\hat{f}(x; D)] \\
 &= (\text{Bias}[\hat{f}(x; D)])^2 + \sigma^2 + \text{Var}[\hat{f}(x; D)]
 \end{aligned}$$

- Finally, MSE loss is obtained by taking the expectation over x

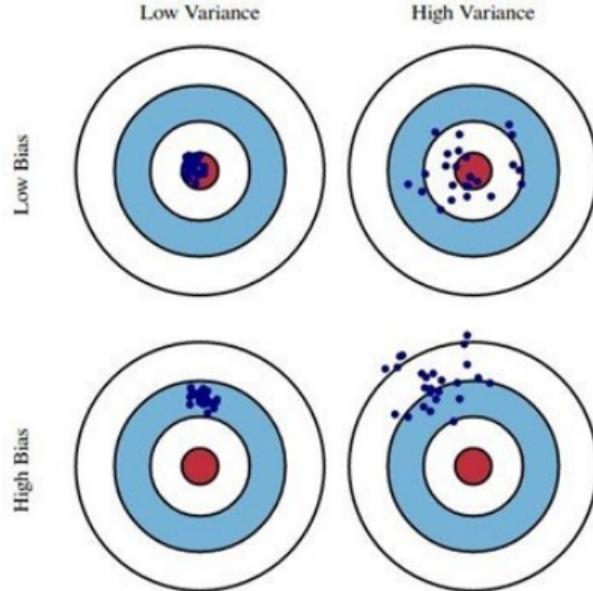
$$\text{MSE} = \mathbb{E}_x[(\text{Bias}[\hat{f}(x; D)])^2 + \text{Var}[\hat{f}(x; D)]] + \sigma^2$$

- ▶ Summarizing, the expected error on an unseen sample x ,

$$\begin{aligned}\mathbb{E}_D[(y - \hat{f}(x; D))^2] &= (f(x) - \mathbb{E}_{D'}[\hat{f}(x; D')])^2 + \text{Var}[\epsilon] + \text{Var}[\hat{f}(x; D)] \\ &= (\text{Bias}[\hat{f}(x; D)])^2 + \sigma^2 + \text{Var}[\hat{f}(x; D)]\end{aligned}$$

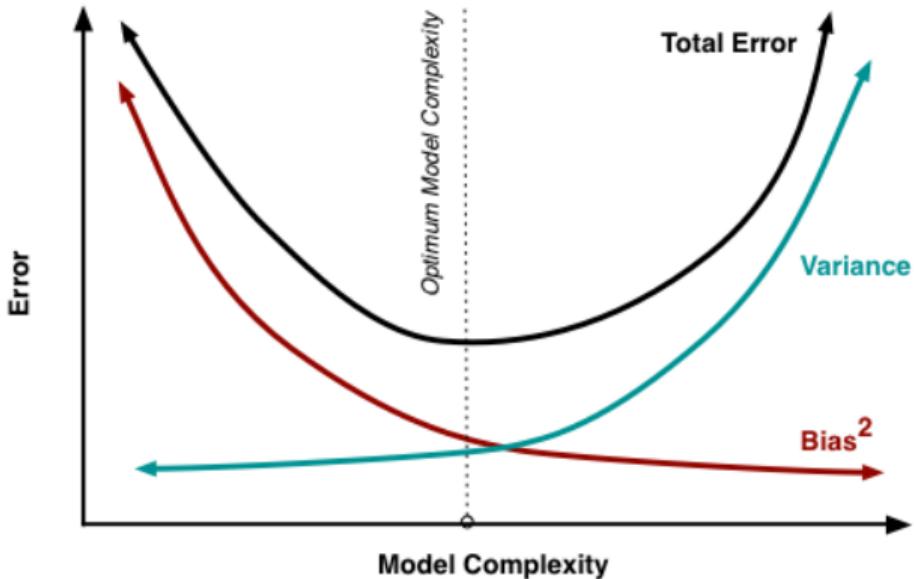
- ▶ The error can be decomposed as the sum of the following terms
 - The square of the *bias*: error caused by the simplifying assumptions built into the method
 - *variance*: how much the learning method $\hat{f}(x; D)$ will move around its mean
 - *irreducible error*: σ^2
- ▶ Bias-variance decomposition also applies to classification and other ML tasks. Look at the wiki page for more details.

Graphical illustration³



³<http://scott.fortmann-roes.com/docs/BiasVariance.html>

Graphical illustration⁴



⁴<http://scott.fortmann-roe.com/docs/BiasVariance.html>

Outline

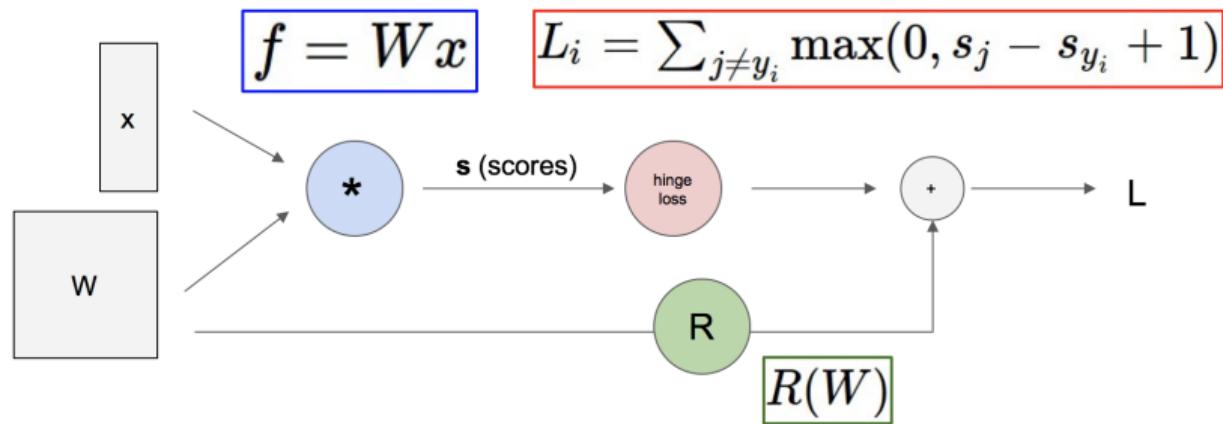
Model fitting

Computational graphs

Neural networks

Computational graphs

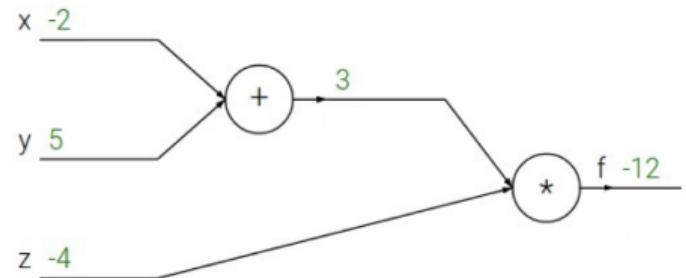
Computational graph



Backpropagation example

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

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

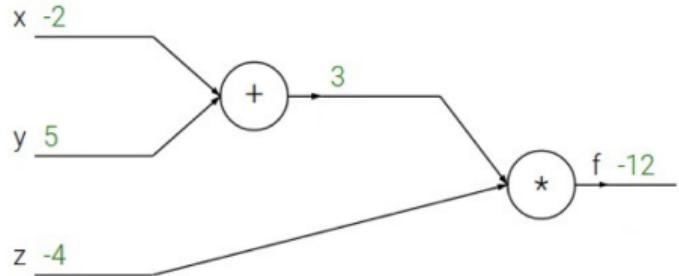


Backpropagation example

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

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

$$\begin{aligned} q &= x + y, \quad \frac{\partial q}{\partial x} = 1, \quad \frac{\partial q}{\partial y} = 1 \\ f &= qz, \quad \frac{\partial f}{\partial q} = z, \quad \frac{\partial f}{\partial z} = q \end{aligned}$$



- Want to compute $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

Backpropagation example

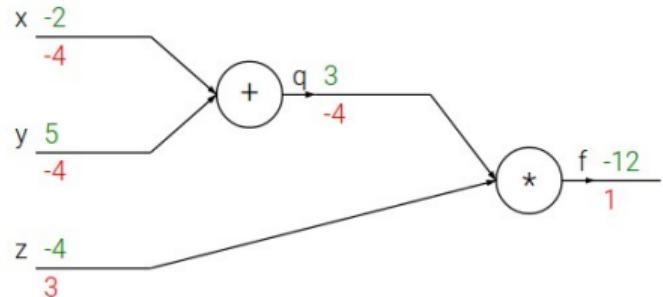
$$f(x, y, z) = (x + y)z$$

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

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

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

- Want to compute $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

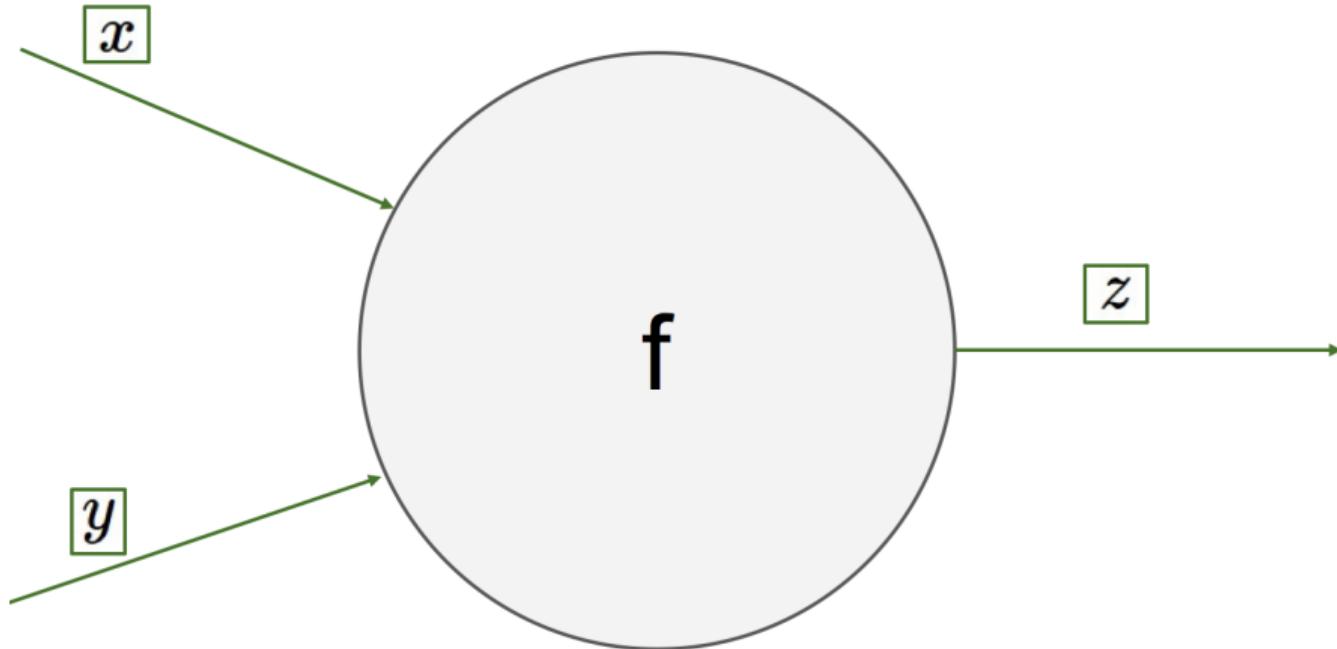


$$\frac{\partial f}{\partial f} = 1, \frac{\partial f}{\partial z} = 3, \frac{\partial f}{\partial q} = -4$$

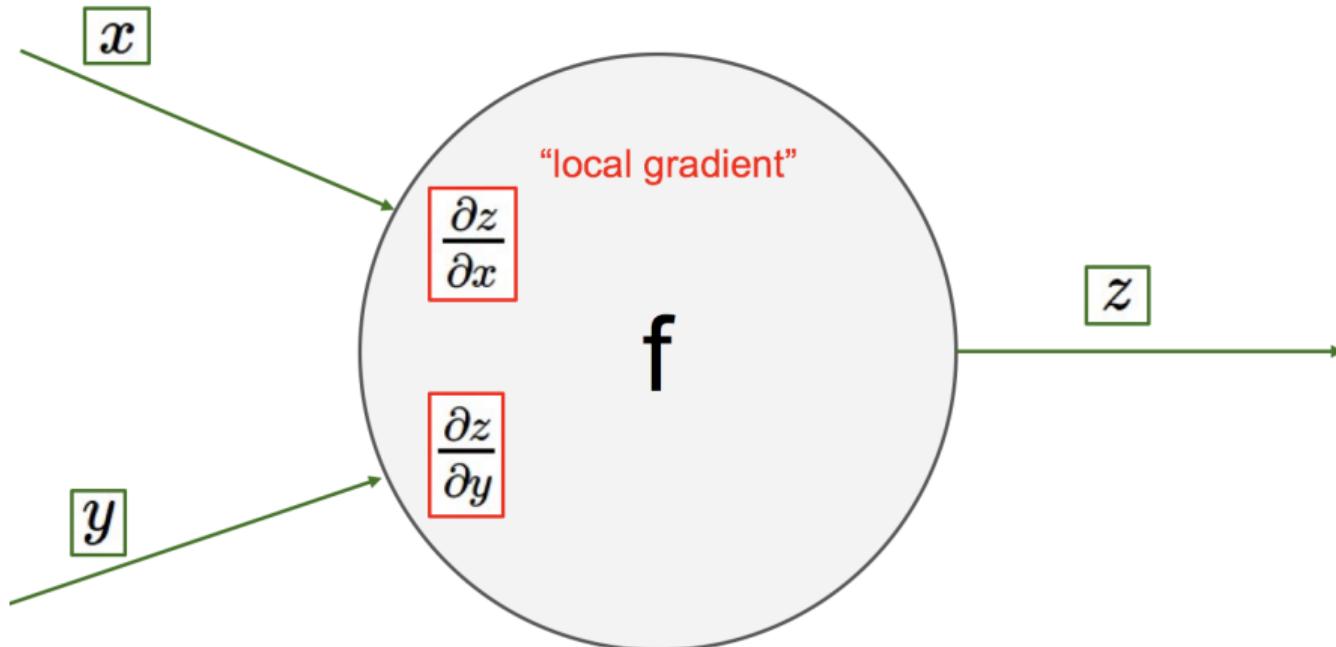
$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y} = -4(1) = -4$$

$$\frac{\partial f}{\partial x} = \underbrace{\frac{\partial f}{\partial q}}_{\text{Upstream gradient}} \underbrace{\frac{\partial q}{\partial x}}_{\text{Downstream gradient}} = -4(1) = -4$$

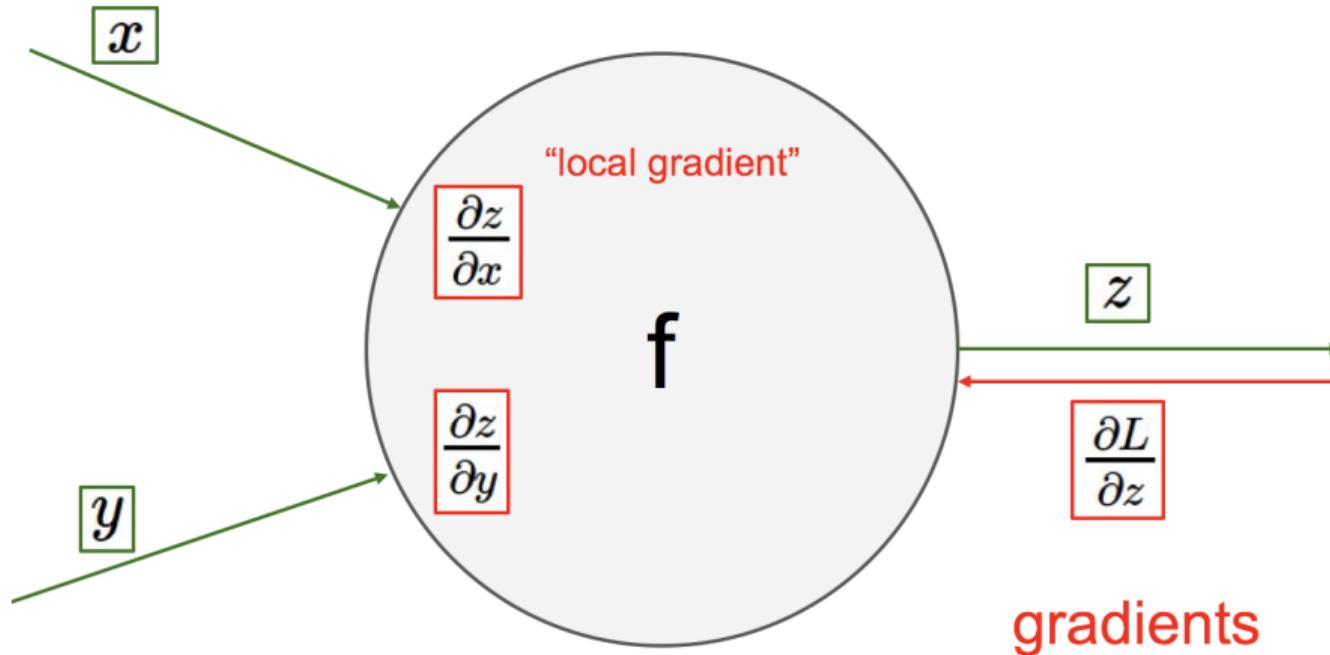
Nodes in computational graph



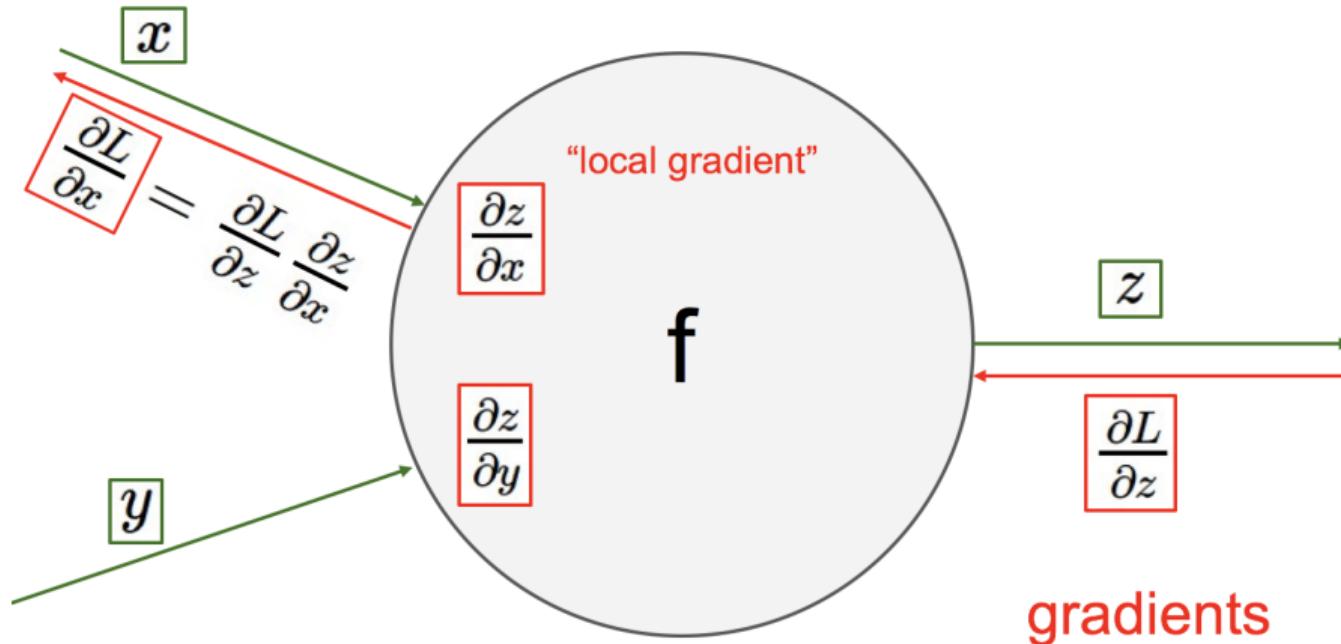
Nodes in computational graph



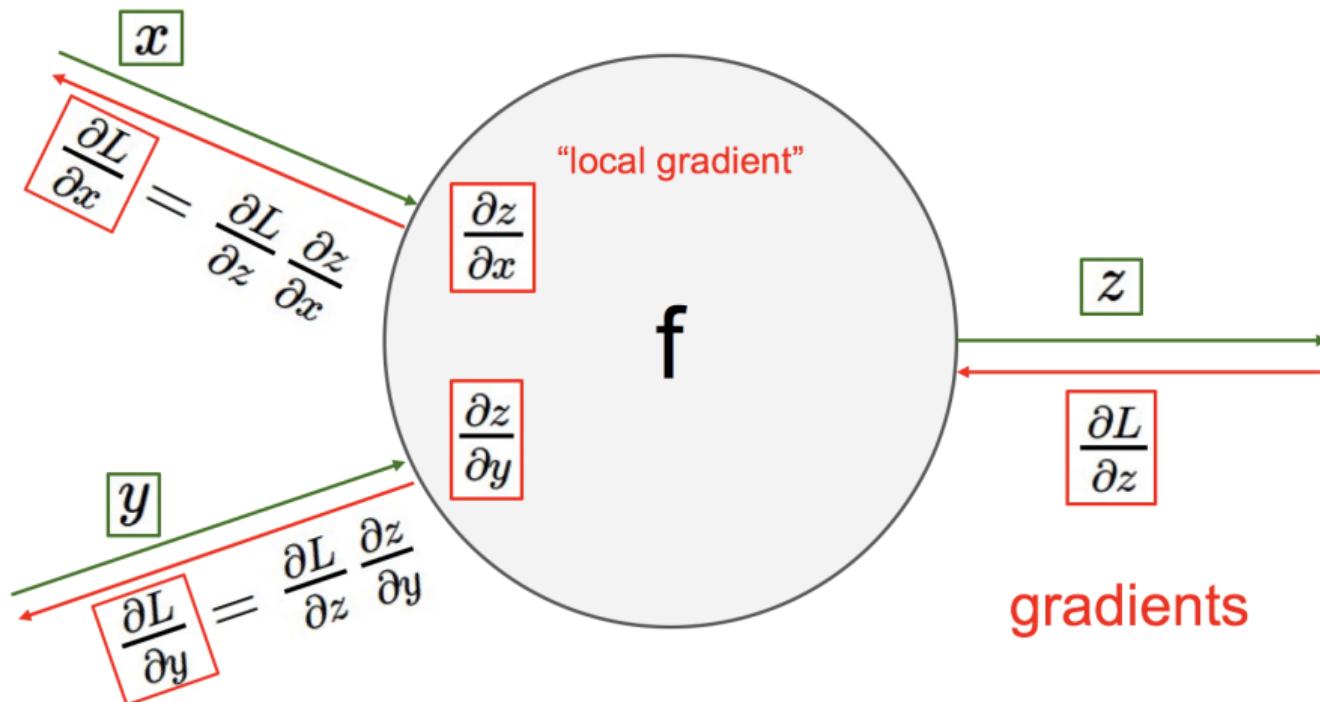
Nodes in computational graph



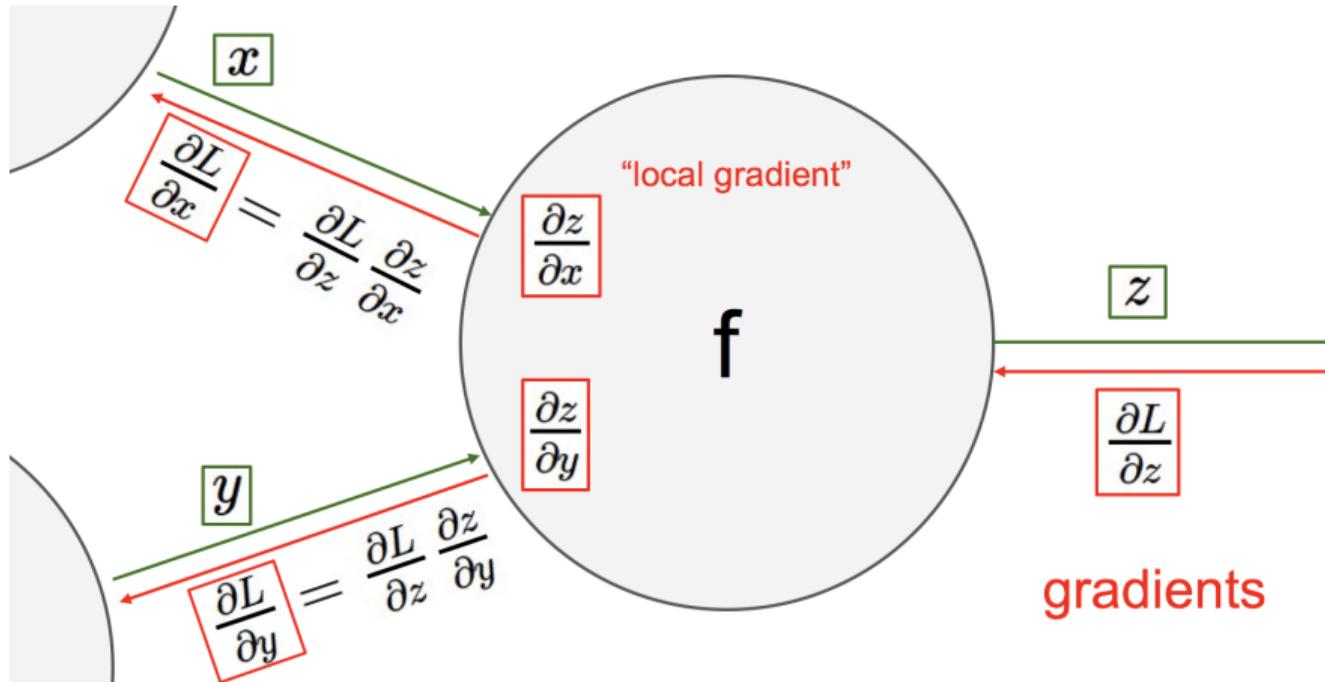
Nodes in computational graph



Nodes in computational graph



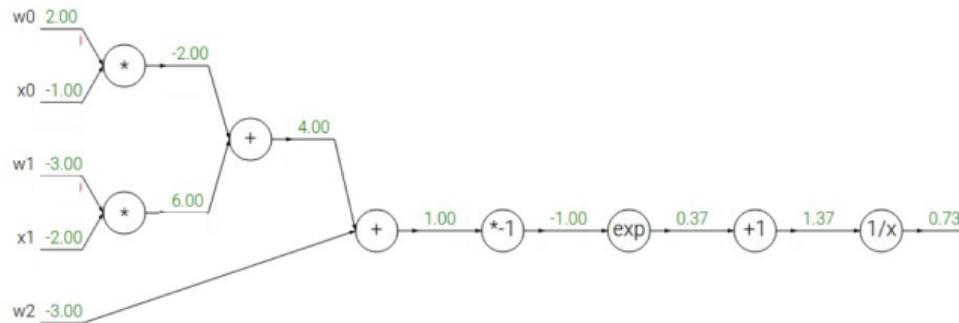
Nodes in computational graph



Sigmoid example

Another example:

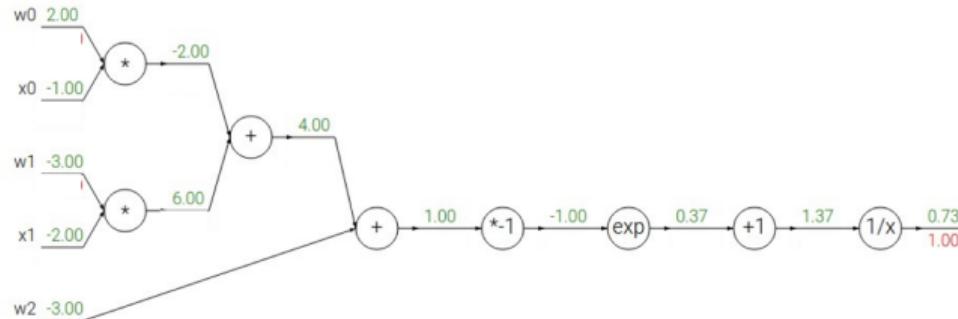
$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$



Sigmoid example

Another example:

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



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

→

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

$$f_a(x) = ax$$

→

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

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

→

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

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

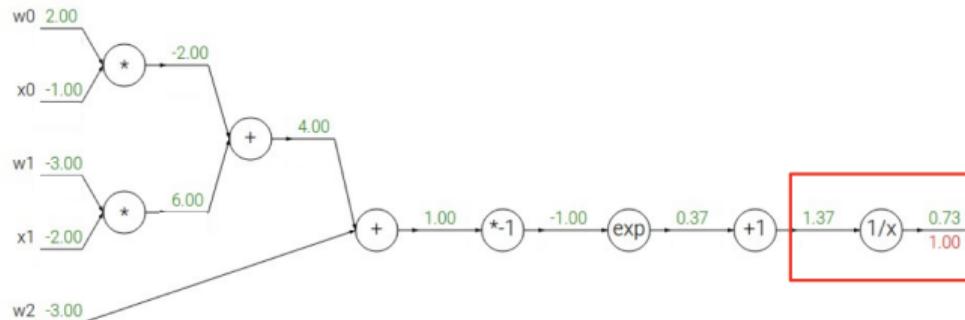
→

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

Sigmoid example

Another example:

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



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

→

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

$$f_a(x) = ax$$

→

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

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

→

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

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

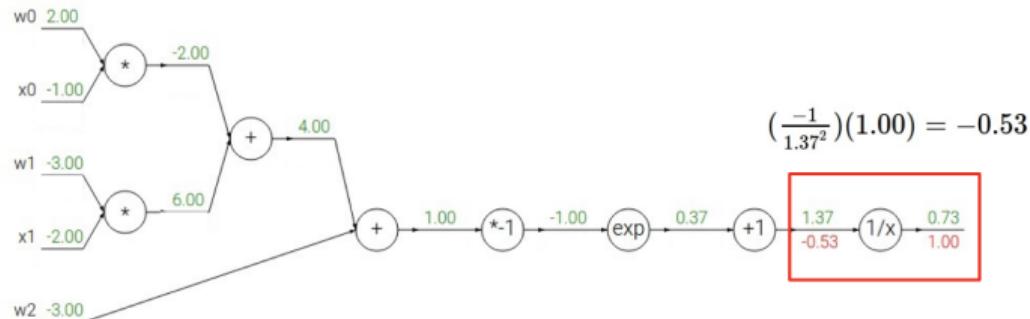
→

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

Sigmoid example

Another example:

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



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

→

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

$$f_a(x) = ax$$

→

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

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

→

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

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

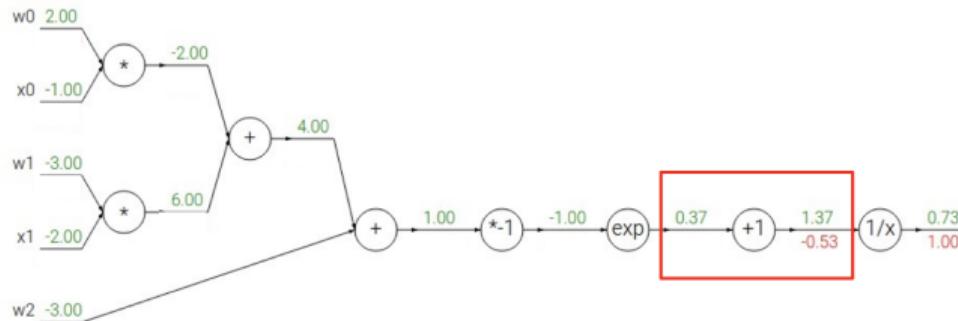
→

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

Sigmoid example

Another example:

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



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

→

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

$$f_a(x) = ax$$

→

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

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

→

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

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

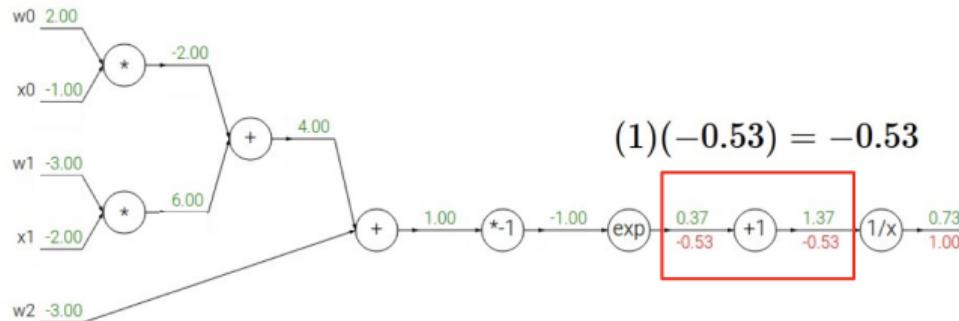
→

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

Sigmoid example

Another example:

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



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

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

→

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

$$f_a(x) = ax$$

→

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

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

→

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

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

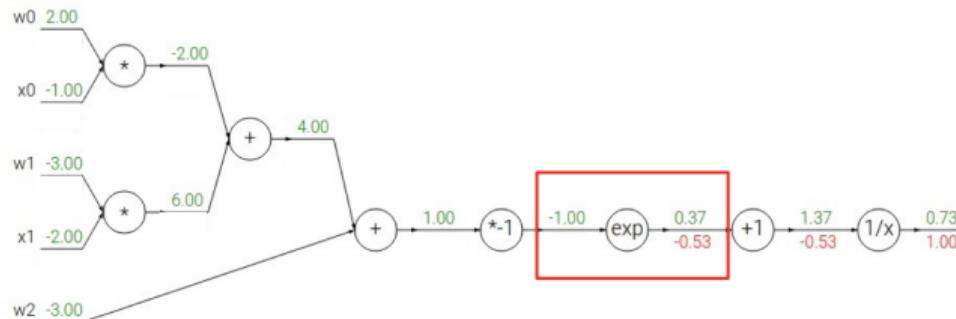
→

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

Sigmoid example

Another example:

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



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

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

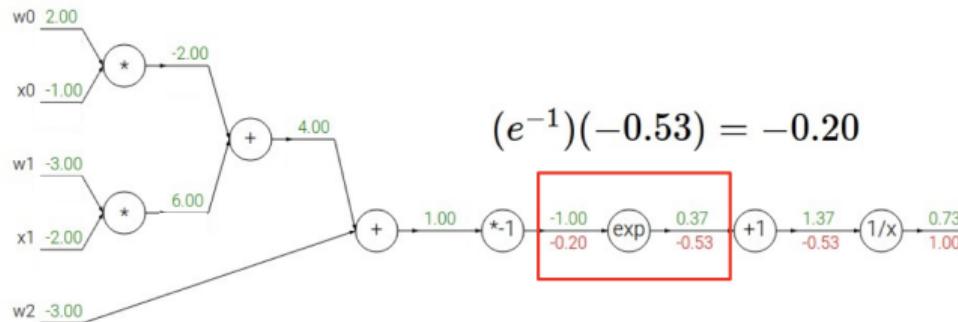
$$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -1/x^2$$

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

Sigmoid example

Another example:

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



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

→

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

$$f_a(x) = ax$$

→

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

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

→

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

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

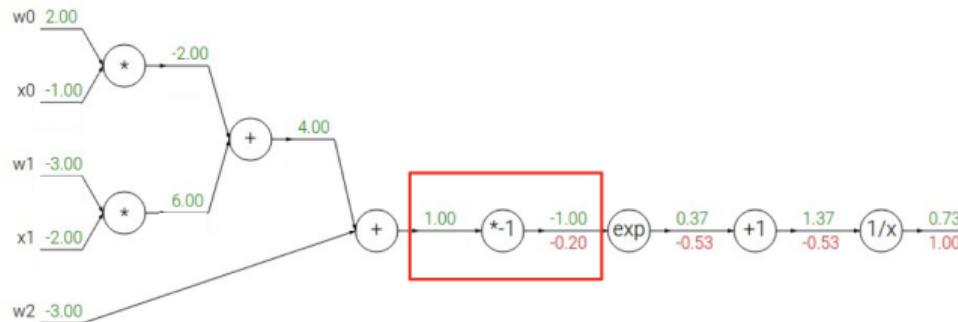
→

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

Sigmoid example

Another example:

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



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

→

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

$$f_a(x) = ax$$

→

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

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

→

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

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

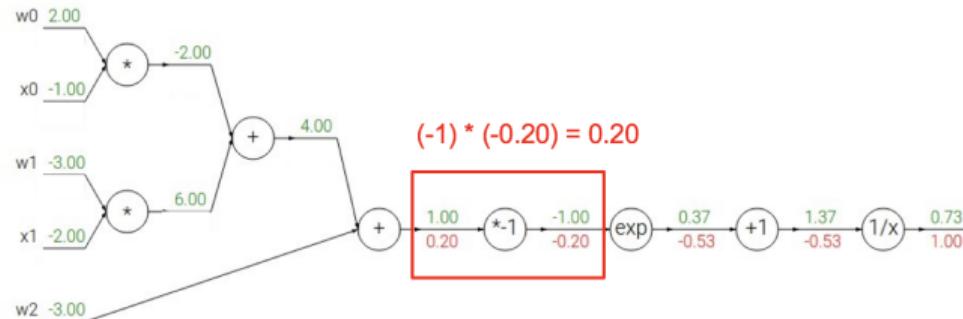
→

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

Sigmoid example

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

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

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

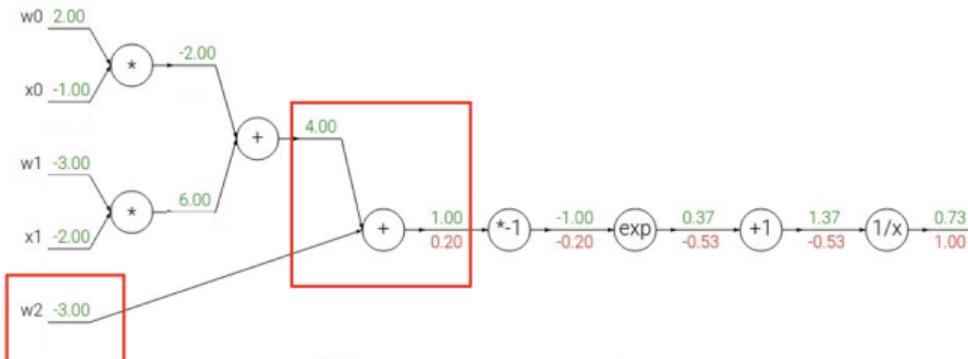
\rightarrow

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

Sigmoid example

Another example:

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



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

→

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

$$f_a(x) = ax$$

→

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

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

→

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

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

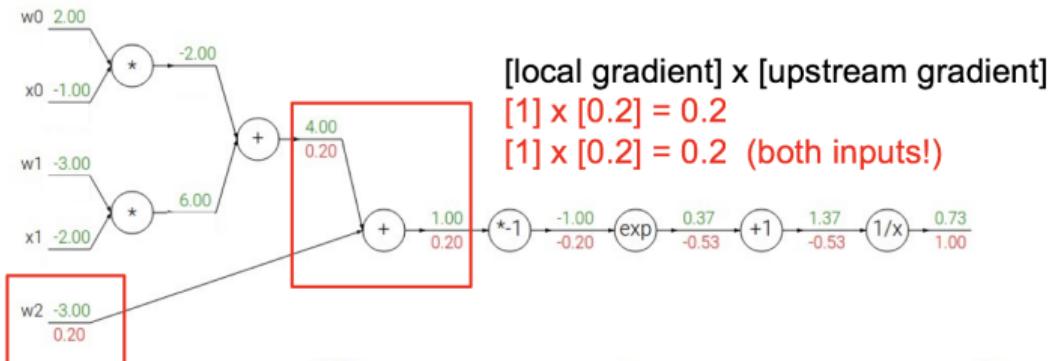
→

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

Sigmoid example

Another example:

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



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

→

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

$$f_a(x) = ax$$

→

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

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

→

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

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

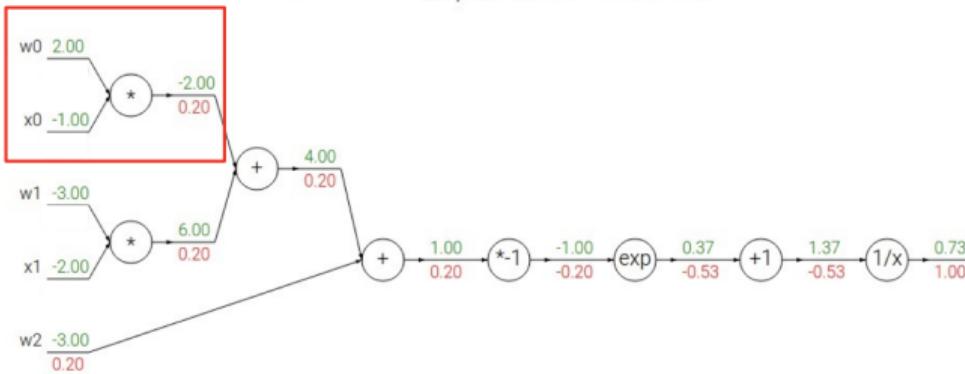
→

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

Sigmoid example

Another example:

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



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

→

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

$$f_a(x) = ax$$

→

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

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

→

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

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

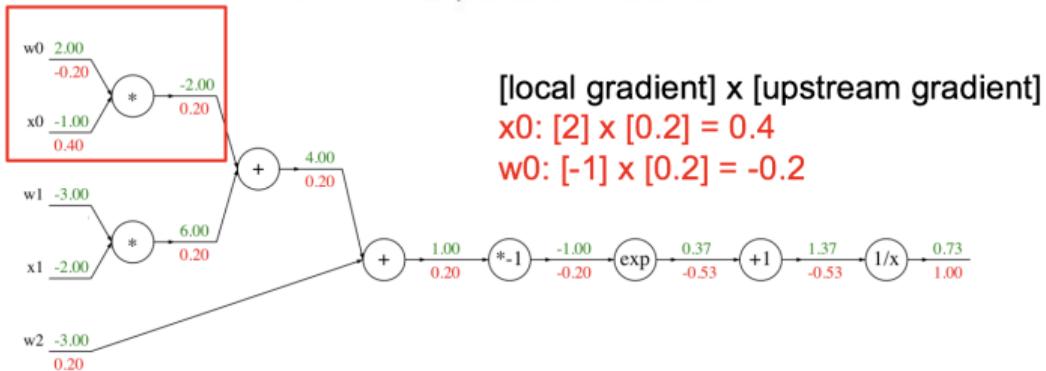
→

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

Sigmoid example

Another example:

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



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

→

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

$$f_a(x) = ax$$

→

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

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

→

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

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

→

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

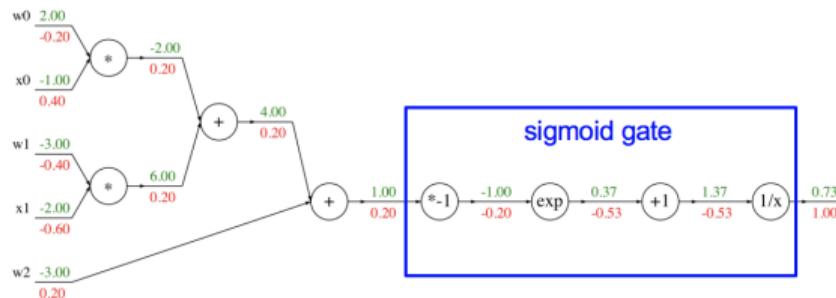
Sigmoid example

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

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

sigmoid function

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



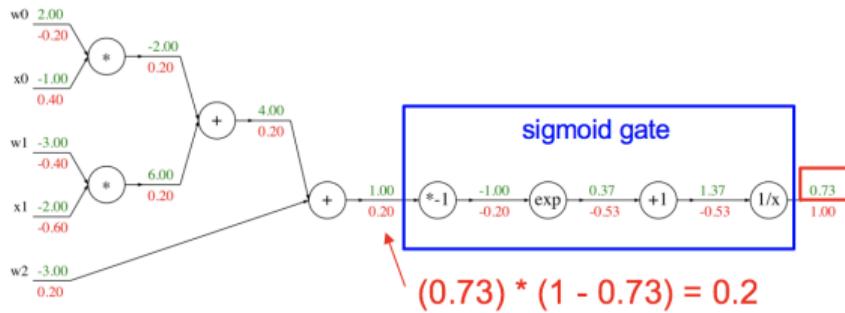
Sigmoid example

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

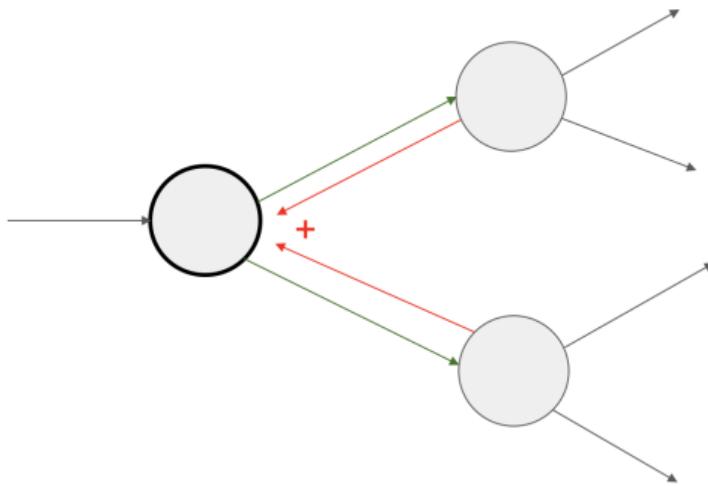
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

sigmoid function

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



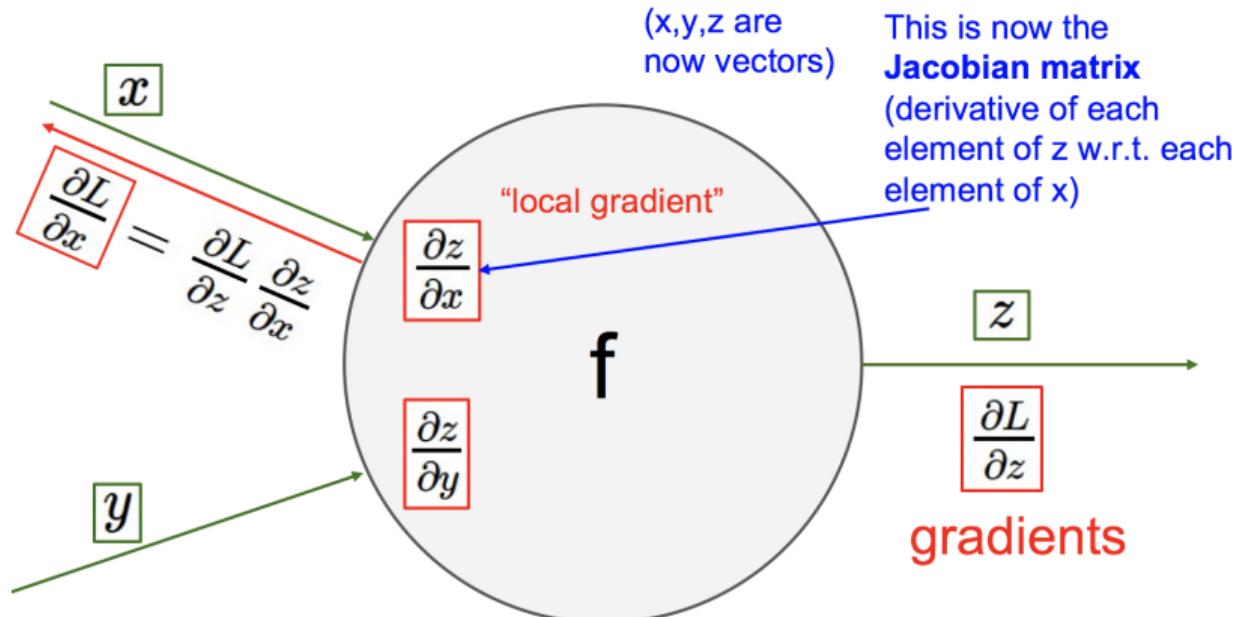
Gradients add at branches



Chain rule

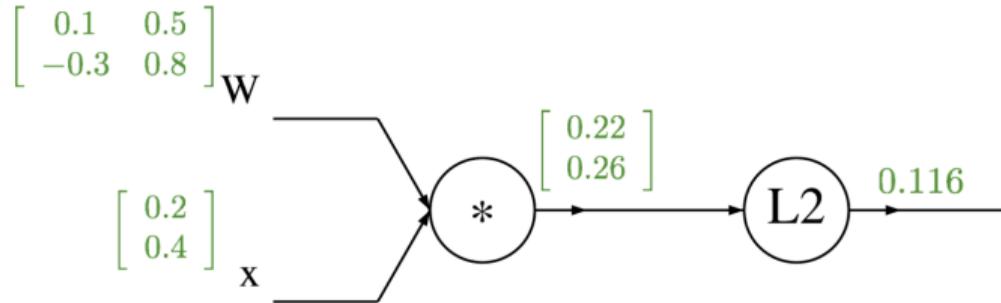
$$\frac{\partial f}{\partial x} = \sum_i \frac{\partial f}{\partial q_i} \frac{\partial q_i}{\partial x}$$

Gradients for vectorized code



Vectorized example

$$f(x, W) = \|Wx\|^2 = \sum_{i=1}^n (Wx)_i^2$$



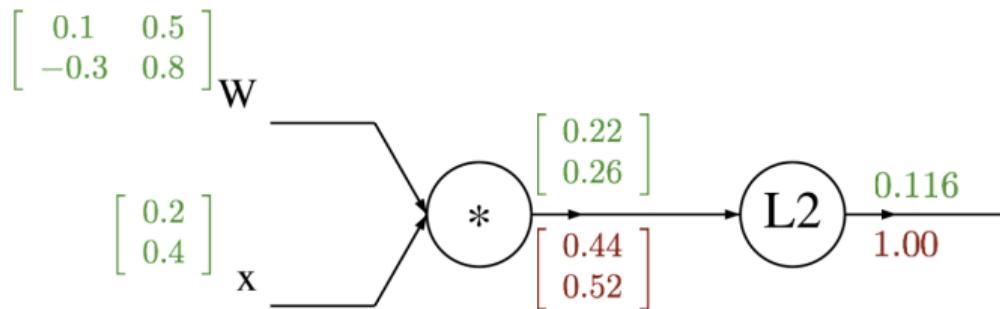
$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \cdots + q_n^2$$

Computational graphs

Vectorized example

$$f(x, W) = \|Wx\|^2 = \sum_{i=1}^n (Wx)_i^2$$



$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \cdots + q_n^2$$

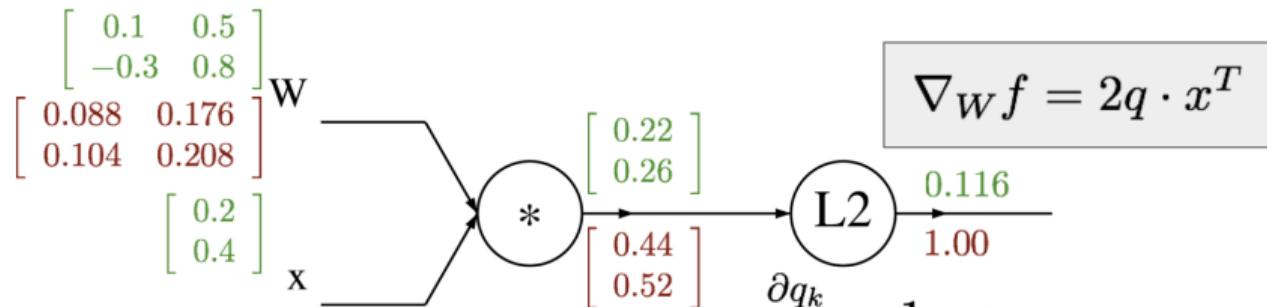
$$\frac{\partial f}{\partial q_i} = 2q_i$$

$$\boxed{\nabla_q f = 2q}$$

Computational graphs

Vectorized example

$$f(x, W) = \|Wx\|^2 = \sum_{i=1}^n (Wx)_i^2$$



$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

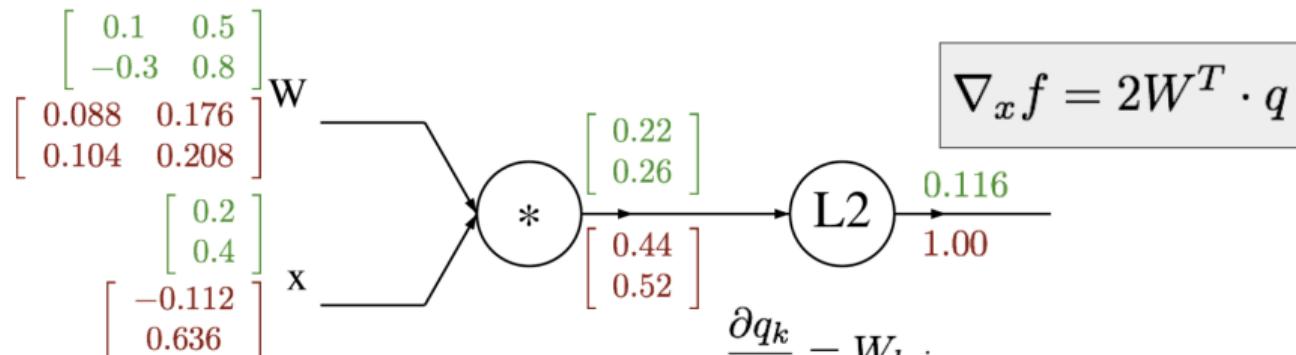
$$f(q) = \|q\|^2 = q_1^2 + \cdots + q_n^2$$

$$\frac{\partial q_k}{\partial W_{i,j}} = \mathbf{1}_{k=i} x_j$$

$$\begin{aligned} \frac{\partial f}{\partial W_{i,j}} &= \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial W_{i,j}} \\ &= \sum_k (2q_k)(\mathbf{1}_{k=i} x_j) \\ &= 2q_i x_j \end{aligned}$$

Vectorized example

$$f(x, W) = \|Wx\|^2 = \sum_{i=1}^n (Wx)_i^2$$



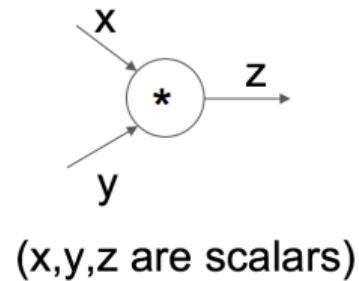
$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \cdots + q_n^2$$

$$\begin{aligned} \frac{\partial q_k}{\partial x_i} &= W_{k,i} \\ \frac{\partial f}{\partial x_i} &= \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial x_i} \\ &= \sum_k 2q_k W_{k,i} \end{aligned}$$

Computational graphs

Modularized 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]
```

Example: Caffe layers

branch: master	cafe / src / caffe / layers /	Create new file	Upload files	Find file	History
sheilhamer committed on GitHub Merge pull request #4630 from 80Dev/load_hdf5_fix	Latest commit: e87ea71 21 days ago				
abval_layer.cpp	dismantle layer headers	a year ago	oudrn_lcn_layer.cpp	dismantle layer headers	a year ago
abval_layer.cu	dismantle layer headers	a year ago	oudrn_lcn_layer.cu	dismantle layer headers	a year ago
accuracy_layer.cpp	dismantle layer headers	a year ago	oudrn_im_layer.cpp	dismantle layer headers	a year ago
argmax_layer.cpp	dismantle layer headers	a year ago	oudrn_im_layer.cu	dismantle layer headers	a year ago
base_conv_layer.cpp	enable dilated deconvolution	a year ago	oudrn_pooling_layer.cpp	dismantle layer headers	a year ago
base_data_layer.cpp	Using default from proto for prefetch	3 months ago	oudrn_pooling_layer.cu	dismantle layer headers	a year ago
base_data_layer.cu	Switched multi-GPU to NCCL	3 months ago	oudrn_relu_layer.cpp	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
batch_norm_layer.cpp	Add missing spaces besides equal signs in batch_norm_layer.cpp	4 months ago	oudrn_relu_layer.cu	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
batch_norm_layer.cu	dismantle layer headers	a year ago	oudrn_sigmoid_layer.cpp	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
batch_reindex_layer.cpp	dismantle layer headers	a year ago	oudrn_sigmoid_layer.cu	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
batch_reindex_layer.cu	dismantle layer headers	a year ago	oudrn_softmax_layer.cpp	dismantle layer headers	a year ago
bias_layer.cpp	Remove incorrect cast of gemm Int arg to Dtype in BiasLayer	a year ago	oudrn_softmax_layer.cu	dismantle layer headers	a year ago
bias_layer.cu	Separation and generalization of ChannelwiseAffineLayer into BiasLayer	a year ago	oudrn_tanh_layer.cpp	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
bnl_layer.cpp	dismantle layer headers	a year ago	oudrn_tanh_layer.cu	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
bnl_layer.cu	dismantle layer headers	a year ago	data_layer.cpp	Switched multi-GPU to NCCL	3 months ago
concat_layer.cpp	dismantle layer headers	a year ago	deconv_layer.cpp	enable dilated deconvolution	a year ago
concat_layer.cu	dismantle layer headers	a year ago	deconv_layer.cu	dismantle layer headers	a year ago
contrastive_loss_layer.cpp	dismantle layer headers	a year ago	dropout_layer.cpp	supporting N-D Blobs in Dropout layer Reshape	a year ago
contrastive_loss_layer.cu	dismantle layer headers	a year ago	dropout_layer.cu	dismantle layer headers	a year ago
conv_layer.cpp	add support for 2D dilated convolution	a year ago	dummy_data_layer.cpp	dismantle layer headers	a year ago
conv_layer.cu	dismantle layer headers	a year ago	eltwise_layer.cpp	dismantle layer headers	a year ago
crop_layer.cpp	remove redundant operations in Crop layer (#5138)	2 months ago	eltwise_layer.cu	dismantle layer headers	a year ago
crop_layer.cu	remove redundant operations in Crop layer (#5138)	2 months ago	elu_layer.cpp	ELU layer with basic tests	a year ago
euclidean_conv_layer.cpp	dismantle layer headers	a year ago	elu_layer.cu	ELU layer with basic tests	a year ago
euclidean_conv_layer.cu	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago	embed_layer.cpp	dismantle layer headers	a year ago
			embed_layer.cu	dismantle layer headers	a year ago
			euclidean_loss_layer.cpp	dismantle layer headers	a year ago
			euclidean_loss_layer.cu	dismantle layer headers	a year ago
			exp_layer.cpp	Solving issue with exp layer with base e	a year ago
			exp_layer.cu	dismantle layer headers	a year ago

Caffe Sigmoid layers

```
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<blob>& propagate_down,
27                                              const vector<blob>& 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_CPU(SigmoidLayer);
42 #endif
43
44 INSTANTIATE_CLASS(SigmoidLayer);
45
46
47 } // namespace caffe
```

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

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

Dynamic programming interpretation

Outline

Model fitting

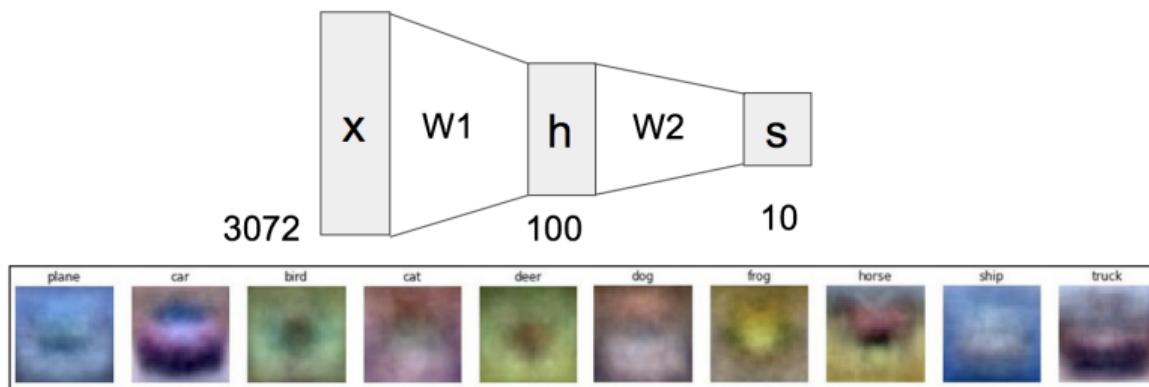
Computational graphs

Neural networks

Neural networks

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

(Now) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$



- ▶ Linear combination of intermediate templates

Multi-layer perceptron (MLP)

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

(Now) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$
or 3-layer Neural Network

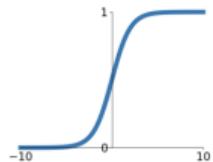
$$f = W_3 \max(0, W_2 \max(0, W_1 x))$$

- ▶ Can go many layers deep

Non-linearities

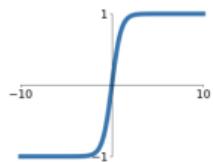
Sigmoid

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



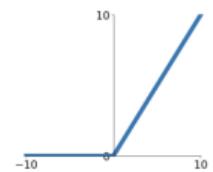
tanh

$$\tanh(x)$$



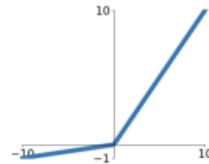
ReLU

$$\max(0, x)$$



Leaky ReLU

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

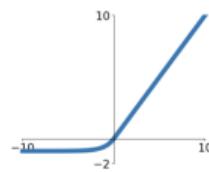


Maxout

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

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



- ▶ What happens to MLP without non-linearities?