

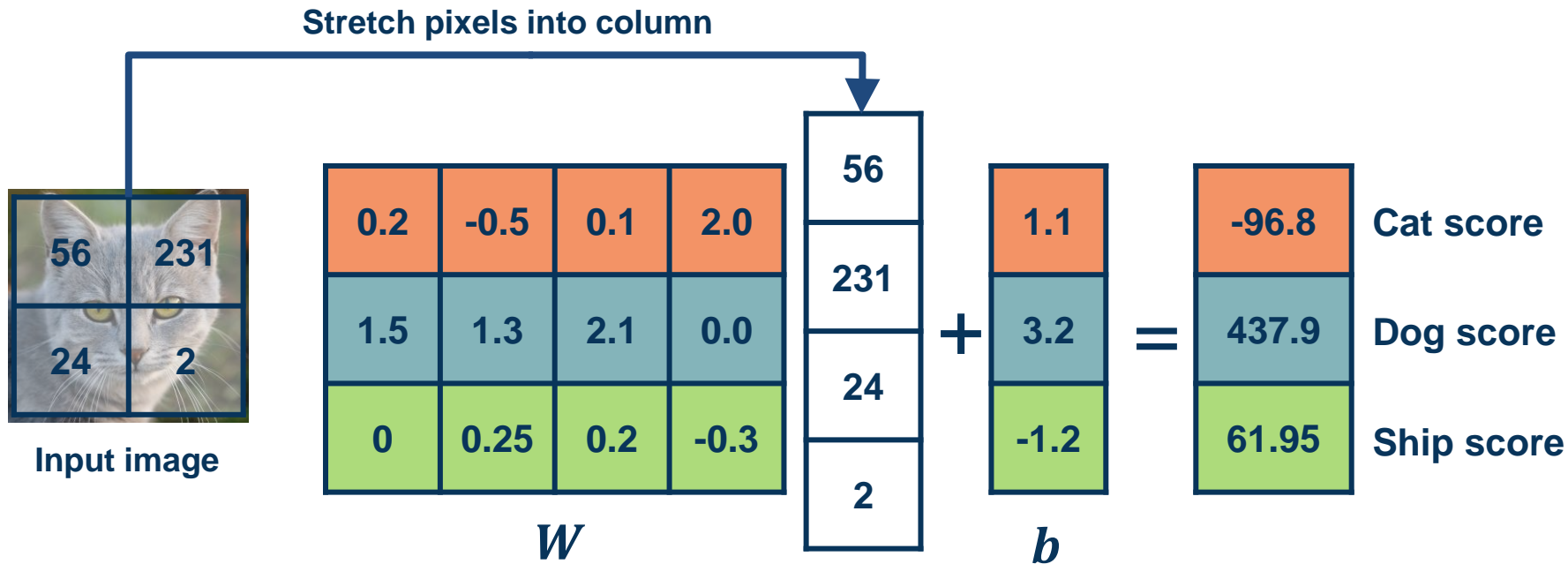
Topics:

- Backpropagation
- Matrix/Linear Algebra view

CS 4644-DL / 7643-A
ZSOLT KIRA

- **Assignment 1 out!**
 - **Due Jan 31st (with grace period Feb 2nd)**
 - Start now, start now, start now!
 - Start now, start now, start now!
 - Start now, start now, start now!
- Resources:
 - These lectures
 - [Matrix calculus for deep learning](#)
 - [Gradients notes](#) and [MLP/ReLU Jacobian notes](#).
 - **Topic OH:** Assignment 1 and matrix calculus (@93)
- Piazza: Project teaming thread
 - **Project Proposal: Feb. 14th, Project Check-in: Mar. 14th.**
 - Project proposal overview during my OH (Thursday 3:30pm ET, recorded)

Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

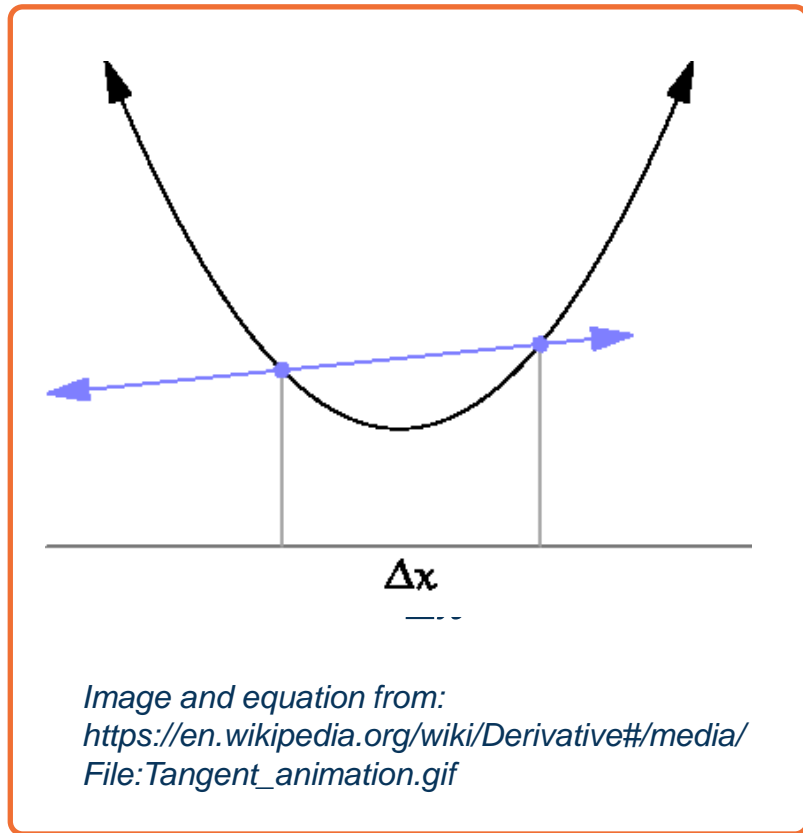


Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n

- We can find the steepest descent direction by computing the **derivative (gradient)**:

$$f'(a) = \lim_{h \rightarrow 0} \frac{f(a + h) - f(a)}{h}$$

- Steepest descent direction is the **negative gradient**
- **Intuitively:** Measures how the function changes as the argument a changes by a small step size
 - As step size goes to zero
- **In Machine Learning:** Want to know how the **loss function** changes **as weights** are varied
 - Can consider each parameter separately by taking **partial derivative** of loss function with respect to that parameter



The same two-layered neural network corresponds to adding another weight matrix

- ◆ We will prefer the linear algebra view, but use some terminology from neural networks (& biology)

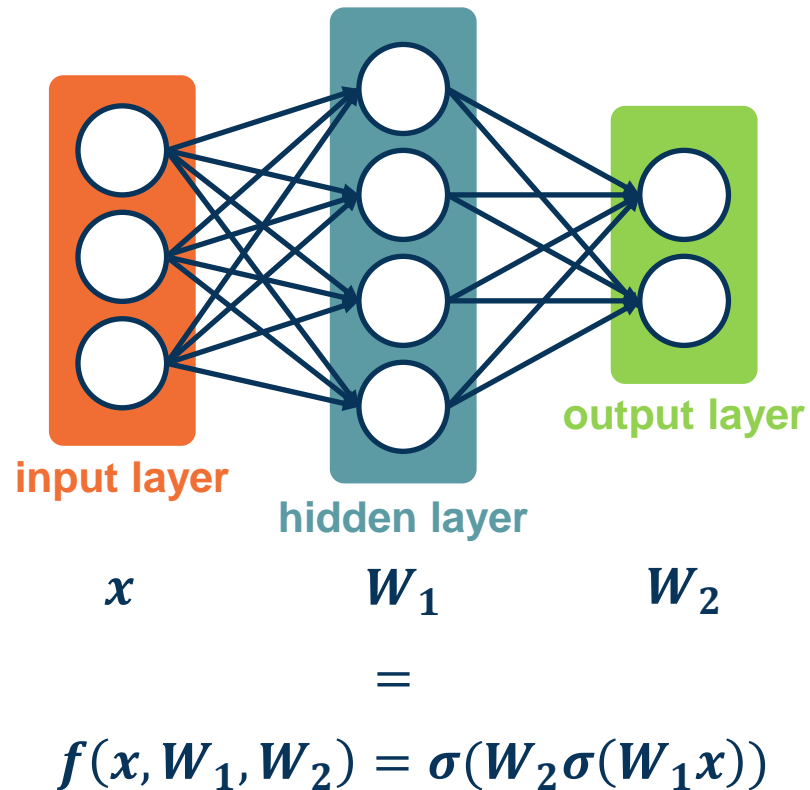


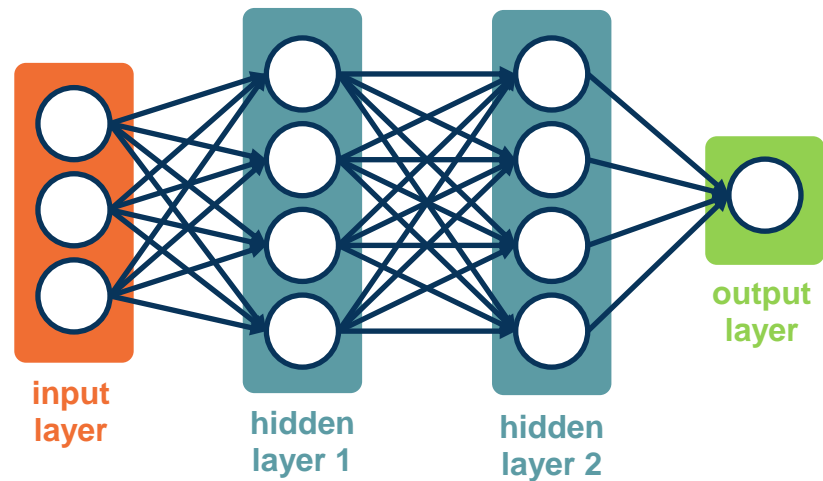
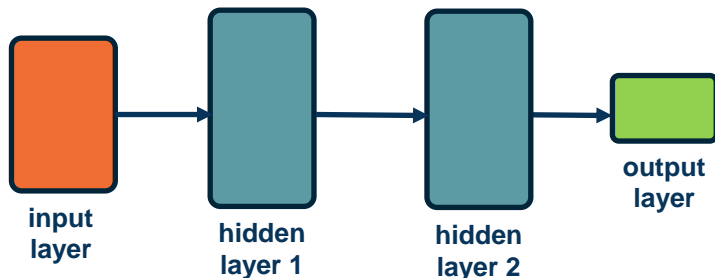
Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Large (deep) networks can be built by adding more and more layers

Three-layered neural networks can represent **any function**

- The number of nodes could grow unreasonably (exponential or worse) with respect to the complexity of the function

We will show them **without edges**:

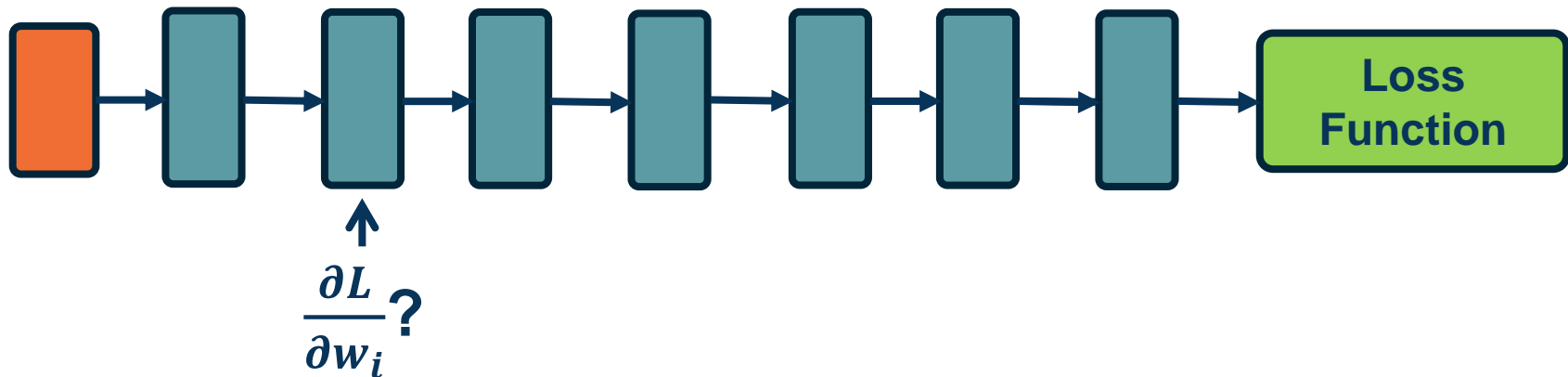


$$f(x, W_1, W_2, W_3) = \sigma(W_2 \sigma(W_1 x))$$

Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Adding More Layers!

- We are learning **complex models** with significant amount of parameters (millions or billions)
- How do we compute the gradients of the **loss** (at the end) with respect to **internal** parameters?
- Intuitively, want to understand how **small changes** in weight deep inside **are propagated** to affect the **loss function** at the end

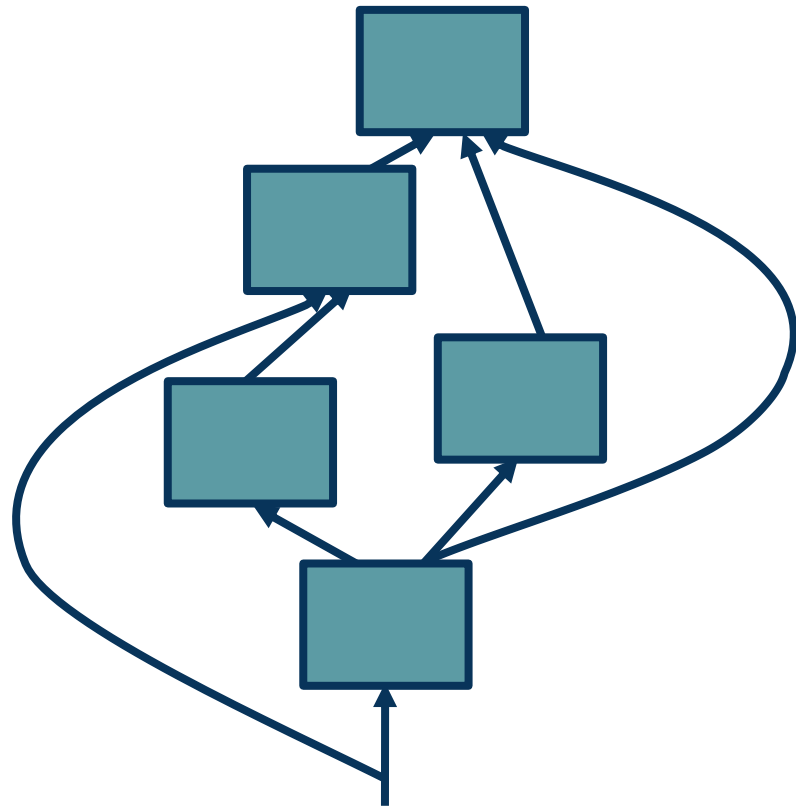


To develop a general algorithm for this, we will view the function as a **computation graph**

Graph can be any **directed acyclic graph (DAG)**

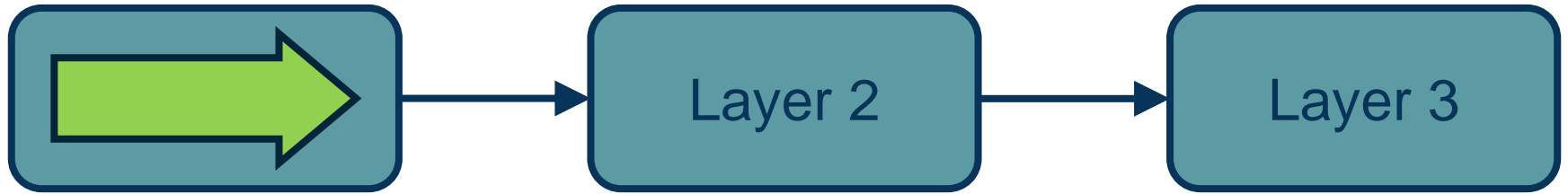
- ◆ Modules must be differentiable to support gradient computations for gradient descent

A **training algorithm** will then process this graph, **one module at a time**



Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

Step 1: Compute Loss on Mini-Batch: Forward Pass



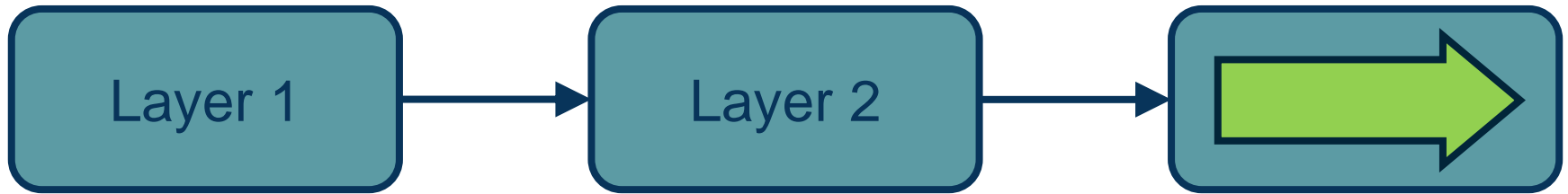
Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

Step 1: Compute Loss on Mini-Batch: Forward Pass



Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

Step 1: Compute Loss on Mini-Batch: Forward Pass



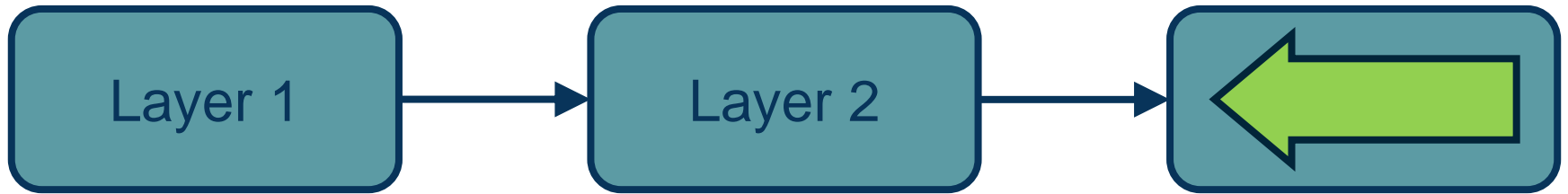
Note that we must store the **intermediate outputs of all layers!**

- ◆ This is because we will need them to **compute the gradients** (the gradient equations will have terms with the output values in them)

Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

Step 1: Compute Loss on Mini-Batch: Forward Pass

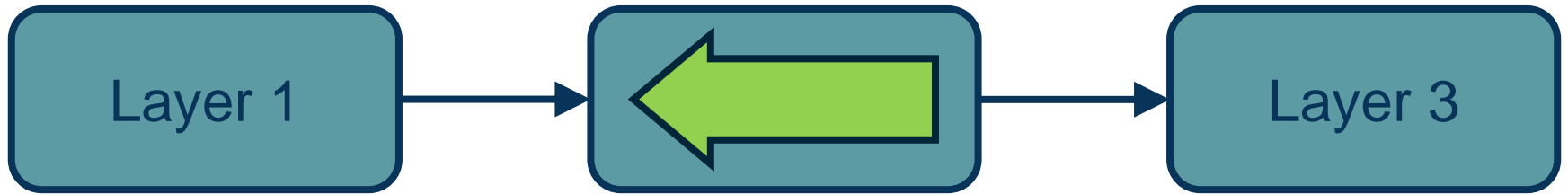
Step 2: Compute Gradients wrt parameters: Backward Pass



Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

Step 1: Compute Loss on Mini-Batch: Forward Pass

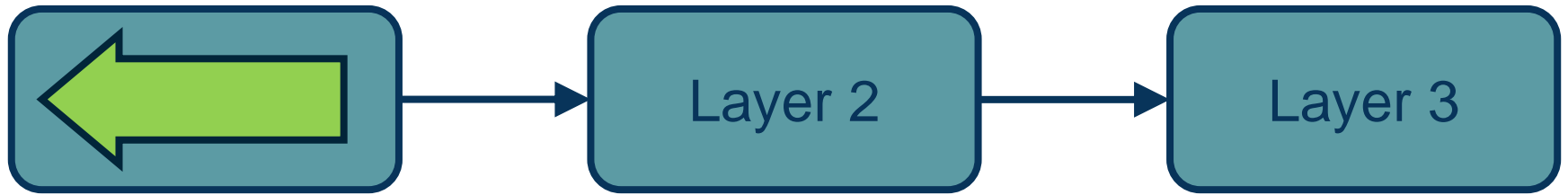
Step 2: Compute Gradients wrt parameters: Backward Pass



Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

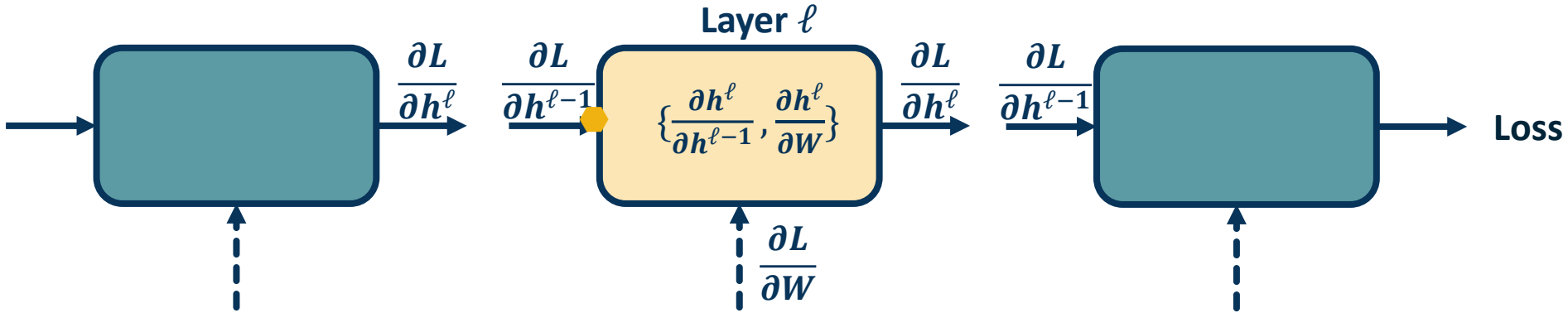
Step 1: Compute Loss on Mini-Batch: Forward Pass

Step 2: Compute Gradients wrt parameters: Backward Pass



Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

- ◆ We want to compute: $\left\{ \frac{\partial L}{\partial h^{\ell-1}}, \frac{\partial L}{\partial W} \right\}$



- ◆ We will use the *chain rule* to do this:

Chain Rule: $\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial x}$

$$\frac{\partial L}{\partial h^{\ell-1}} = \frac{\partial L}{\partial h^{\ell}} \frac{\partial h^{\ell}}{\partial h^{\ell-1}}$$

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial h^{\ell}} \frac{\partial h^{\ell}}{\partial W}$$

Step 1: Compute Loss on Mini-Batch: **Forward Pass**

Step 2: Compute Gradients wrt parameters: **Backward Pass**

Step 3: Use **gradient** to update **all parameters** at the end



$$w_i = w_i - \alpha \frac{\partial L}{\partial w_i}$$

Backpropagation is the application of gradient descent to a computation graph via the chain rule!



Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

Backpropagation: a simple example

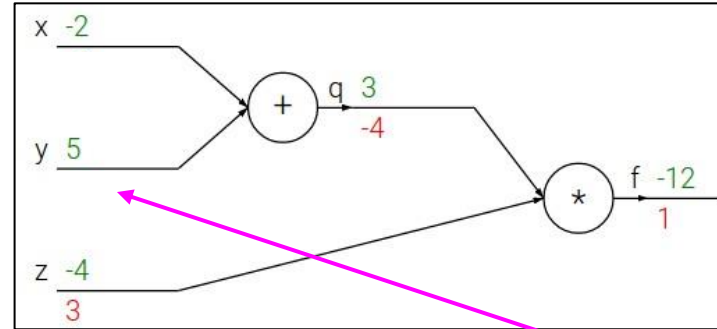
$$f(x, y, z) = (x + y)z$$

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

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

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

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



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

Chain rule:

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

Upstream
gradient

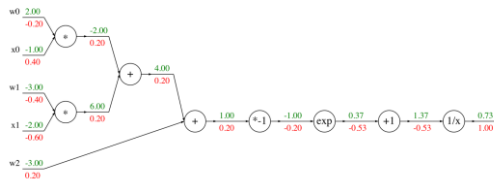
Local
gradient

Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Deep Learning = Differentiable Programming

- Computation = Graph
 - Input = Data + Parameters
 - Output = Loss
 - Scheduling = Topological ordering
- What do we need to do?
 - Generic code for representing the graph of modules
 - Specify modules (both forward and backward function)
 - Backpropagation implementation on the graph

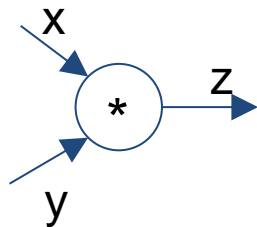
Modularized implementation: forward / backward API



Graph (or Net) object (*rough psuedo code*)

```
class ComputationalGraph(object):  
    #...  
    def forward(inputs):  
        # 1. [pass inputs to input gates...]  
        # 2. forward the computational graph:  
        for gate in self.graph.nodes_topologically_sorted():  
            gate.forward()  
        return loss # the final gate in the graph outputs the loss  
    def backward():  
        for gate in reversed(self.graph.nodes_topologically_sorted()):  
            gate.backward() # little piece of backprop (chain rule applied)  
        return inputs_gradients
```

Modularized implementation: forward / backward API



(x,y,z are scalars)

```
class MultiplyGate(object):
```

```
    def forward(x,y):
```

```
        z = x*y
```

```
        return z
```

```
    def backward(dz):
```

```
        # dx = ... #todo
```

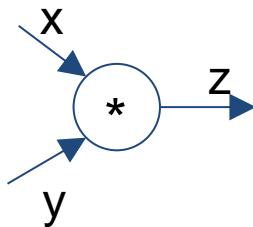
```
        # dy = ... #todo
```

```
        return [dx, dy]
```

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

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

Modularized implementation: forward / backward API



(x,y,z are scalars)

```
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 [caffe](#) / [src](#) / [caffe](#) / [layers](#) /

Create new file Upload files Find file History

shelhamer committed on GitHub Merge pull request #4630 from BIGene/load_hdf5_fix Latest commit e687a71 21 days ago

..

absval_layer.cpp	dismantle layer headers	a year ago
absval_layer.cu	dismantle layer headers	a year ago
accuracy_layer.cpp	dismantle layer headers	a year ago
argmax_layer.cpp	dismantle layer headers	a year ago
base_conv_layer.cpp	enable dilated deconvolution	a year ago
base_data_layer.cpp	Using default from proto for prefetch	3 months ago
base_data_layer.cu	Switched multi-GPU to NCCL	3 months ago
batch_norm_layer.cpp	Add missing spaces besides equal signs in batch_norm_layer.cpp	4 months ago
batch_norm_layer.cu	dismantle layer headers	a year ago
batch_reindex_layer.cpp	dismantle layer headers	a year ago
batch_reindex_layer.cu	dismantle layer headers	a year ago
bias_layer.cpp	Remove incorrect cast of gemm int arg to Dtype in BiasLayer	a year ago
bias_layer.cu	Separation and generalization of ChannelwiseAffineLayer into BiasLayer	a year ago
bn1_layer.cpp	dismantle layer headers	a year ago
bn1_layer.cu	dismantle layer headers	a year ago
concat_layer.cpp	dismantle layer headers	a year ago
concat_layer.cu	dismantle layer headers	a year ago
contrastive_loss_layer.cpp	dismantle layer headers	a year ago
contrastive_loss_layer.cu	dismantle layer headers	a year ago
conv_layer.cpp	add support for 2D dilated convolution	a year ago
conv_layer.cu	dismantle layer headers <small>Caffe is licensed under BSD 2-Clause</small>	a year ago
crop_layer.cpp	remove redundant operations in Crop layer (#5138)	2 months ago
crop_layer.cu	remove redundant operations in Crop layer (#5138)	2 months ago
cudnn_conv_layer.cpp	dismantle layer headers	a year ago
cudnn_conv_layer.cu	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago

cudnn_lcn_layer.cpp	dismantle layer headers	a year ago
cudnn_lcn_layer.cu	dismantle layer headers	a year ago
cudnn_lrn_layer.cpp	dismantle layer headers	a year ago
cudnn_lrn_layer.cu	dismantle layer headers	a year ago
cudnn_pooling_layer.cpp	dismantle layer headers	a year ago
cudnn_pooling_layer.cu	dismantle layer headers	a year ago
cudnn_relu_layer.cpp	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
cudnn_relu_layer.cu	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
cudnn_sigmoid_layer.cpp	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
cudnn_sigmoid_layer.cu	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
cudnn_softmax_layer.cpp	dismantle layer headers	a year ago
cudnn_softmax_layer.cu	dismantle layer headers	a year ago
cudnn_tanh_layer.cpp	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
cudnn_tanh_layer.cu	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
data_layer.cpp	Switched multi-GPU to NCCL	3 months ago
deconv_layer.cpp	enable dilated deconvolution	a year ago
deconv_layer.cu	dismantle layer headers	a year ago
dropout_layer.cpp	supporting N-D Blobs in Dropout layer Reshape	a year ago
dropout_layer.cu	dismantle layer headers	a year ago
dummy_data_layer.cpp	dismantle layer headers	a year ago
eltwise_layer.cpp	dismantle layer headers	a year ago
eltwise_layer.cu	dismantle layer headers	a year ago
elu_layer.cpp	ELU layer with basic tests	a year ago
elu_layer.cu	ELU layer with basic tests	a year 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 Layer

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

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

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

[Caffe is licensed under BSD 2-Clause](#)

- Neural networks involves composing simple functions into a **computation graph**
- Optimization (updating weights) of this graph is through backpropagation
 - Recursive algorithm: Gradient descent (partial derivatives) plus chain rule
- Remaining questions:
 - How does this work with vectors, matrices, tensors?
 - Across a composed function?
 - How can we implement this algorithmically to make these calculations automatic? **Automatic Differentiation**

**Linear
Algebra
View:
Vector and
Matrix Sizes**

$$\begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1m} & b_1 \\ w_{21} & w_{22} & \cdots & w_{2m} & b_2 \\ w_{31} & w_{32} & \cdots & w_{3m} & b_3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \\ 1 \end{bmatrix}$$

W

x

Sizes: $[c \times (m + 1)]$ $[(m + 1) \times 1]$

Where c is number of classes

m is dimensionality of input

Conventions:

- Size of derivatives for scalars, vectors, and matrices:

Assume we have scalar $s \in \mathbb{R}^1$, vector $v \in \mathbb{R}^m$, i.e. $v = [v_1, v_2, \dots, v_m]^T$ and matrix $M \in \mathbb{R}^{m_1 \times m_2}$

	S []	V []	M []
S	$\frac{\partial s_1}{\partial s_2}$ []	$\frac{\partial s}{\partial v}$ []	$\frac{\partial s}{\partial M}$ []
V	$\frac{\partial v}{\partial s}$ []	$\frac{\partial v_1}{\partial v_2}$ []	Tensors
M	$\frac{\partial M}{\partial s}$ []		

Conventions:

- Size of derivatives for scalars, vectors, and matrices:

Assume we have scalar $s \in \mathbb{R}^1$, vector $\mathbf{v} \in \mathbb{R}^m$, i.e. $\mathbf{v} = [v_1, v_2, \dots, v_m]^T$ and matrix $\mathbf{M} \in \mathbb{R}^{m_1 \times m_2}$

- What is the size of $\frac{\partial \mathbf{v}}{\partial s}$? $\mathbb{R}^{m \times 1}$ (column vector of size m)

- What is the size of $\frac{\partial s}{\partial \mathbf{v}}$? $\mathbb{R}^{1 \times m}$ (row vector of size m)

$$\begin{bmatrix} \frac{\partial v_1}{\partial s} \\ \frac{\partial v_2}{\partial s} \\ \vdots \\ \frac{\partial v_m}{\partial s} \end{bmatrix}$$

$$\left[\frac{\partial s}{\partial v_1} \quad \frac{\partial s}{\partial v_2} \quad \dots \quad \frac{\partial s}{\partial v_m} \right]$$

Conventions:

- What is the size of $\frac{\partial v^1}{\partial v^2}$? A matrix:

$$\begin{array}{c} \text{Row } i \\ \left[\begin{array}{cccc} \frac{\partial v^1_1}{\partial v^1_1} & \dots & \dots & \dots \\ \frac{\partial v^2_1}{\partial v^1_1} & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ \frac{\partial v^1_i}{\partial v^1_1} & \dots & \frac{\partial v^1_i}{\partial v^2_j} & \dots \\ \frac{\partial v^2_1}{\partial v^1_1} & \dots & \frac{\partial v^2_1}{\partial v^2_j} & \dots \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \end{array} \right] \\ \text{Col } j \\ m_1 \times m_2 \end{array}$$

- This matrix of partial derivatives is called a **Jacobian**

(Note this is slightly different convention than on [Wikipedia](https://en.wikipedia.org/wiki/Jacobian_matrix)). Also, computationally other conventions are used.

Conventions:

- What is the size of $\frac{\partial s}{\partial M}$? A matrix:

$$\begin{bmatrix} \frac{\partial s}{\partial m_{[1,1]}} & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \frac{\partial s}{\partial m_{[i,j]}} & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \end{bmatrix}$$

(Note this is slightly different convention than on [Wikipedia](#)). Also, computationally other conventions are used.

Example 1:

$$y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} x \\ x^2 \end{bmatrix} \quad \frac{\partial y}{\partial x} = \begin{bmatrix} 1 \\ 2x \end{bmatrix}$$

Example 2:

$$y = w^T x = \sum_k w_k x_k$$

$$\frac{\partial y}{\partial x} = \left[\frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_m} \right]$$

$$= [w_1, \dots, w_m]$$

$$= w^T$$

because

$$\frac{\partial (\sum_k w_k x_k)}{\partial x_i} = w_i$$

Example 3:

$$y = Wx$$

$$\frac{\partial y}{\partial x} = W$$

$$\text{Row } i \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \dots & \dots & \dots \\ \dots & \dots & \frac{\partial y_i}{\partial x_j} & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \end{bmatrix} = \begin{bmatrix} \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & w_{ij} & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \end{bmatrix}$$

$$y_i = \sum_j w_{ij} x_j$$

Example 4:

$$\frac{\partial (wAw)}{\partial w} = 2w^T A \text{ (assuming } A \text{ is symmetric)}$$

What is the size of $\frac{\partial L}{\partial W}$?

Remember that loss is a **scalar** and W is a matrix:

$$\begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1m} & b_1 \\ w_{21} & w_{22} & \cdots & w_{2m} & b_2 \\ w_{31} & w_{32} & \cdots & w_{3m} & b_3 \end{bmatrix}$$

Jacobian is also a matrix:

$$\begin{matrix} & & & & W \\ \begin{bmatrix} \frac{\partial L}{\partial w_{11}} & \frac{\partial L}{\partial w_{12}} & \cdots & \frac{\partial L}{\partial w_{1m}} & \frac{\partial L}{\partial b_1} \\ \frac{\partial L}{\partial w_{21}} & \cdots & \cdots & \frac{\partial L}{\partial w_{2m}} & \frac{\partial L}{\partial b_2} \\ \cdots & \cdots & \cdots & \frac{\partial L}{\partial w_{3m}} & \frac{\partial L}{\partial b_3} \end{bmatrix} & & & & \end{matrix}$$

Batches of data are **matrices** or **tensors** (multi-dimensional matrices)

Examples:

- Each instance is a vector of size m , our batch is of size $[B \times m]$
- Each instance is a matrix (e.g. grayscale image) of size $W \times H$, our batch is $[B \times W \times H]$
- Each instance is a multi-channel matrix (e.g. color image with R,B,G channels) of size $C \times W \times H$, our batch is $[B \times C \times W \times H]$

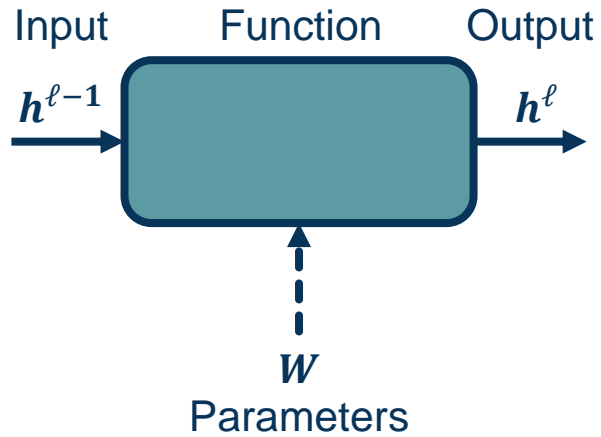
Jacobians become tensors which is complicated

- Instead, flatten input to a vector and get a vector of derivatives!
- This can also be done for partial derivatives between two vectors, two matrices, or two tensors

$$\begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nn} \end{bmatrix}$$

Flatten 

$$\begin{bmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{21} \\ x_{22} \\ \vdots \\ x_{n1} \\ \vdots \\ x_{nn} \end{bmatrix}$$



Define:

$$h_i^l = w_i^T h^{l-1}$$

$$h^l = W h^{l-1}$$

$$\begin{array}{ccc}
 \left[\begin{array}{c} | \\ | \\ | \end{array} \right] & \left[\begin{array}{c} \leftarrow w_i^T \rightarrow \\ | \\ | \\ | \end{array} \right] & \left[\begin{array}{c} | \\ | \\ | \end{array} \right] \\
 |h^l| \times 1 & |h^l| \times |h^{l-1}| & |h^{l-1}| \times 1
 \end{array}$$

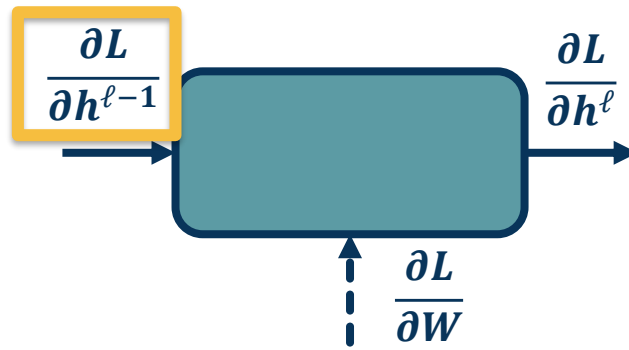
Fully Connected (FC) Layer: Forward Function

$$\mathbf{h}^\ell = \mathbf{W}\mathbf{h}^{\ell-1}$$

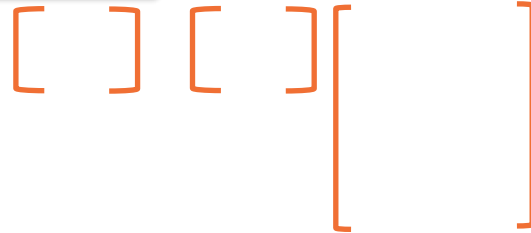
$$\frac{\partial \mathbf{h}^\ell}{\partial \mathbf{h}^{\ell-1}} = \mathbf{W}$$

Define:

$$\mathbf{h}_i^\ell = \mathbf{w}_i^T \mathbf{h}^{\ell-1}$$



$$\frac{\partial L}{\partial \mathbf{h}^{\ell-1}} = \frac{\partial L}{\partial \mathbf{h}^\ell} \frac{\partial \mathbf{h}^\ell}{\partial \mathbf{h}^{\ell-1}}$$



$$1 \times |\mathbf{h}^{\ell-1}| \quad 1 \times |\mathbf{h}^\ell| \quad |\mathbf{h}^\ell| \times |\mathbf{h}^{\ell-1}|$$

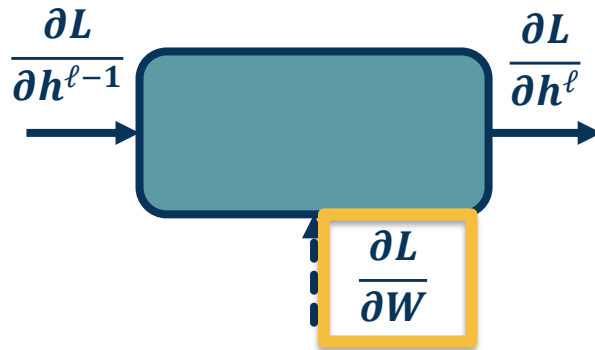
Fully Connected (FC) Layer

$$\mathbf{h}^\ell = \mathbf{W}\mathbf{h}^{\ell-1}$$

$$\frac{\partial \mathbf{h}^\ell}{\partial \mathbf{h}^{\ell-1}} = \mathbf{W}$$

Define:

$$h_i^\ell = w_i^T \mathbf{h}^{\ell-1}$$



Note doing this on full W matrix would result in Jacobian tensor!

But it is *sparse* – each output only affected by corresponding weight row

~~$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial \mathbf{h}^\ell} \frac{\partial \mathbf{h}^\ell}{\partial W}$$~~

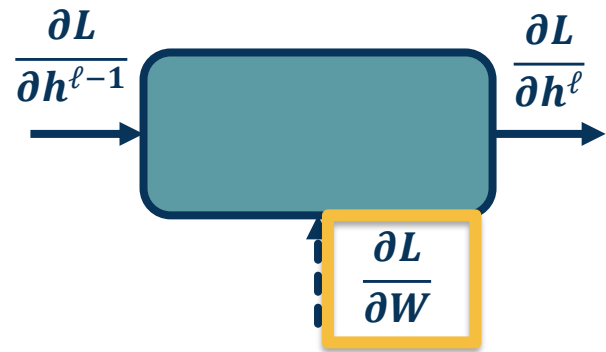
$$\mathbf{h}^\ell = \mathbf{W}\mathbf{h}^{\ell-1}$$

$$\frac{\partial \mathbf{h}^\ell}{\partial \mathbf{h}^{\ell-1}} = \mathbf{W}$$

Define:

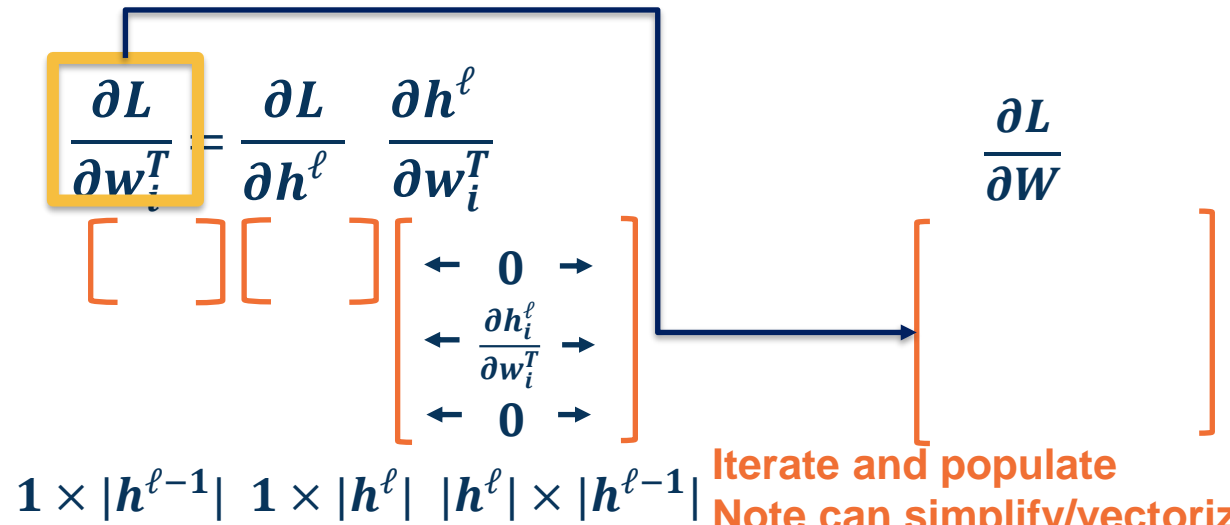
$$\mathbf{h}_i^\ell = \mathbf{w}_i^T \mathbf{h}^{\ell-1}$$

$$\frac{\partial \mathbf{h}_i^\ell}{\partial \mathbf{w}_i^T} = \mathbf{h}^{(\ell-1),T}$$



Note doing this on full W matrix would result in Jacobian tensor!

But it is *sparse* – each output only affected by corresponding weight row



Iterate and populate
Note can simplify/vectorize!

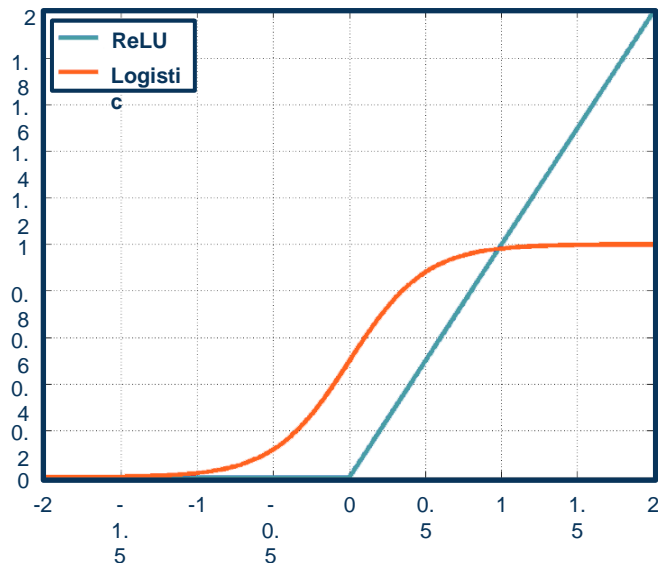
Fully Connected (FC) Layer

We can employ **any differentiable (or piecewise differentiable) function**

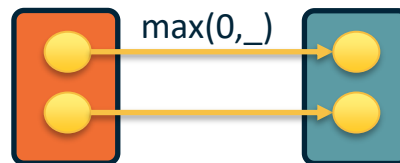
A common choice is the **Rectified Linear Unit**

- Provides non-linearity but better gradient flow than sigmoid
- Performed **element-wise**

How many parameters for this layer?



$$h^\ell = \max(0, h^{\ell-1})$$



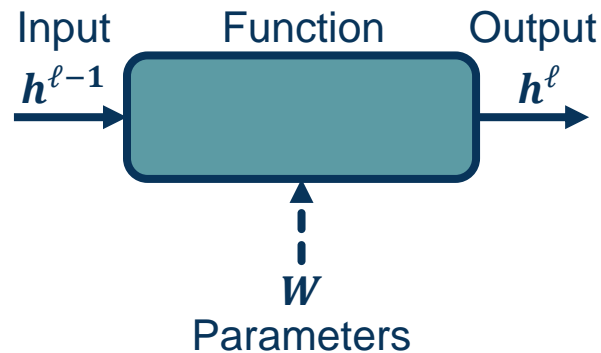
Rectified Linear Unit (ReLU)

Full Jacobian of ReLU layer is **large**
(output dim x input dim)

- But again it is **sparse**
- Only **diagonal values non-zero** because it is element-wise
- An output value affected only by **corresponding input value**

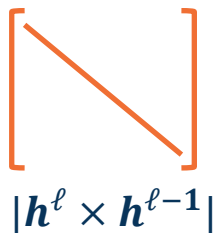
Max function **funnels gradients through selected max**

- Gradient will be **zero** if input ≤ 0



Forward: $h^l = \max(0, h^{l-1})$

Backward: $\frac{\partial L}{\partial h^{l-1}} = \frac{\partial L}{\partial h^l} \frac{\partial h^l}{\partial h^{l-1}}$



For diagonal

$$\frac{\partial h^l}{\partial h^{l-1}} = \begin{cases} 1 & \text{if } h^{l-1} > 0 \\ 0 & \text{otherwise} \end{cases}$$

4D input x :
$$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix}$$

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

(elementwise)

4D output z :
$$\begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$$

What does $\frac{\partial z}{\partial x}$ look like?

4D dL/dz :
$$\begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$$

Upstream
gradient

4D input x:

$$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix}$$

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

(elementwise)

4D output z:

$$\begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$$

4D dL/dx:

$$\begin{bmatrix} 4 \\ 0 \\ 5 \\ 0 \end{bmatrix}$$

[dz/dx] [dL/dz]

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$$

4D dL/dz:

$$\begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$$
Upstream
gradient

For element-wise ops, jacobian is **sparse**: off-diagonal entries always zero!
Never **explicitly** form Jacobian -- instead use elementwise multiplication

- Neural networks involves composing simple functions into a **computation graph**
- Optimization (updating weights) of this graph is through backpropagation
 - Recursive algorithm: Gradient descent (partial derivatives) plus chain rule
- Remaining questions:
 - How does this work with vectors, matrices, tensors?
 - Across a composed function? **Next!**
 - How can we implement this algorithmically to make these calculations automatic? **Automatic Differentiation**

Composition of Functions: $f(g(x)) = (f \circ g)(x)$

A complex function (e.g. defined by a neural network):

$$f(x) = g_\ell (g_{\ell-1} (\dots g_1(x)))$$

$$f(x) = g_\ell \circ g_{\ell-1} \dots \circ g_1(x)$$

(Many of these will be parameterized)

(Note you might find the opposite notation as well!)



Scalar Case

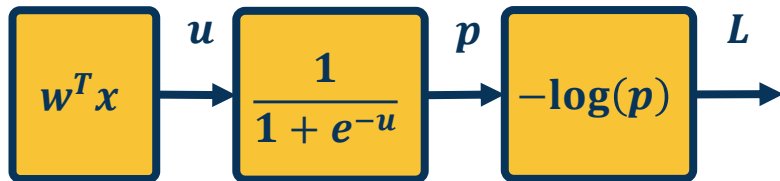
Vector Case

Jacobian View of Chain Rule



Graphical View of Chain Rule

Chain Rule: Cascaded



$$\bar{L} = 1$$

$$\bar{p} = \frac{\partial L}{\partial p} = -\frac{1}{p}$$

where $p = \sigma(w^T x)$ and $\sigma(x) = \frac{1}{1+e^{-x}}$

$$\bar{u} = \frac{\partial L}{\partial u} = \frac{\partial L}{\partial p} \frac{\partial p}{\partial u} = \bar{p} \sigma(1 - \sigma)$$

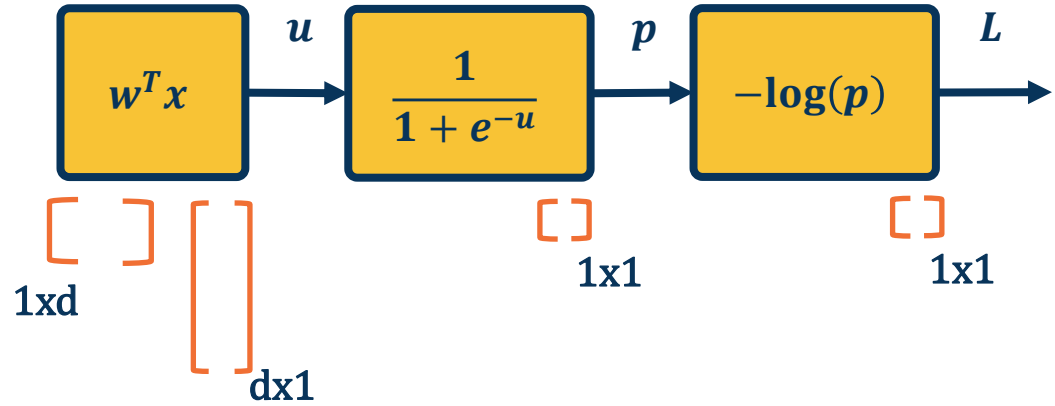
$$\bar{w} = \frac{\partial L}{\partial w} = \frac{\partial L}{\partial u} \frac{\partial u}{\partial w} = \bar{u} x^T$$

We can do this in a combined way to see all terms together:

$$\begin{aligned} \bar{w} &= \frac{\partial L}{\partial p} \frac{\partial p}{\partial u} \frac{\partial u}{\partial w} = -\frac{1}{\sigma(w^T x)} \sigma(w^T x) (1 - \sigma(w^T x)) x^T \\ &= -\left(1 - \sigma(w^T x)\right) x^T \end{aligned}$$

This effectively shows gradient flow along path from L to w

The chain rule can be computed as a **series of scalar, vector, and matrix linear algebra operations**

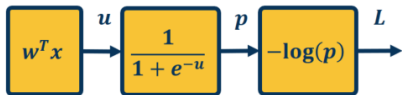


Extremely efficient in graphics processing units (GPUs)

$$\bar{w} = - \frac{1}{\sigma(w^T x)} \sigma(w^T x) (1 - \sigma(w^T x)) x^T$$

Dimensions are indicated by brackets below the terms:

- $\frac{1}{\sigma(w^T x)}$: 1×1
- $\sigma(w^T x)$: 1×1
- $(1 - \sigma(w^T x))$: 1×1
- x^T : $1 \times d$



$$L = -\log(p)$$

$$\bar{p} = \frac{\partial L}{\partial p} = -\frac{1}{p}$$

where $p = \sigma(w^T x)$ and $\sigma(x) = \frac{1}{1+e^{-x}}$

$$\bar{u} = \frac{\partial L}{\partial u} = \frac{\partial L}{\partial p} \frac{\partial p}{\partial u} = \bar{p} \sigma(1-\sigma)$$

$$\bar{w} = \frac{\partial L}{\partial w} = \frac{\partial L}{\partial u} \frac{\partial u}{\partial w} = \bar{u} x^T$$

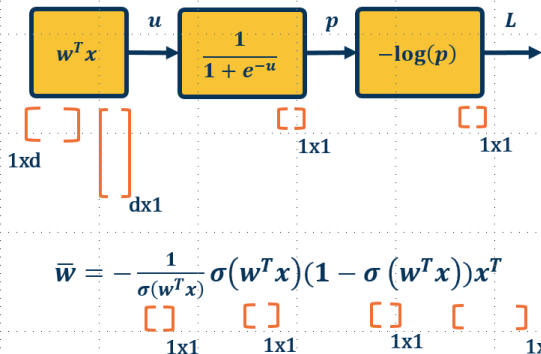
We can do this in a combined way to see all terms together:

$$\bar{w} = \frac{\partial L}{\partial p} \frac{\partial p}{\partial u} \frac{\partial u}{\partial w} = -\frac{1}{\sigma(w^T x)} \sigma(w^T x) (1 - \sigma(w^T x)) x^T$$

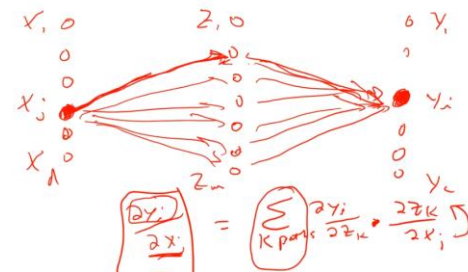
$$= -(1 - \sigma(w^T x)) x^T$$

This effectively shows gradient flow along path from L to w

Computation Graph / Global View of Chain Rule

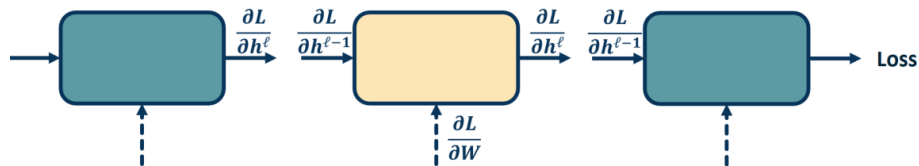


Computational / Tensor View



Graph View

- We want to compute: $\left\{ \frac{\partial L}{\partial h^{\ell-1}}, \frac{\partial L}{\partial W} \right\}$



Backpropagation View (Recursive Algorithm)

Different Views of Equivalent Ideas

- **Backpropagation:** Recursive, modular algorithm for chain rule + gradient descent
- **When we move to vectors and matrices:**
 - Composition of functions (scalar)
 - Composition of functions (vectors/matrices)
 - Jacobian view of chain rule
 - Can view entire set of calculations as linear algebra operations (matrix-vector or matrix-matrix multiplication)
- **Automatic differentiation:**
 - Reduction of modules to simple operations we know (simple multiplication, etc.)
 - Automatically build computation graph in background as write code
 - Automatically compute gradients via backward pass