

Topics:

- Gradient Descent
- Neural Networks

CS 4644-DL / 7643-A
ZSOLT KIRA

- **Assignment 1 out!**
 - **Due Jun 5th (with grace period June 7th)**
 - Start now, start now, start now!
 - Start now, start now, start now!
 - Start now, start now, start now!
- **Piazza**
 - Be active!!!
- **Office hours**
 - Schedules coming out today
- Note: Course will start to get math heavy!
- [Matrix calculus for deep learning](#)

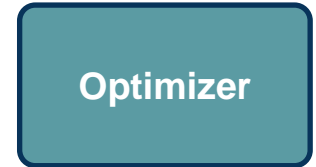
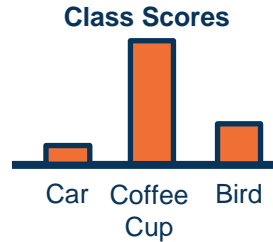
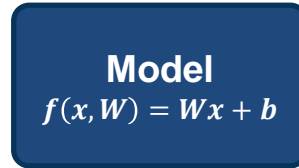
- **Input** (and representation)
- **Functional form** of the model
 - Including parameters
- **Performance measure** to improve
 - Loss or objective function
- **Algorithm** for finding best parameters
 - Optimization algorithm



Data: Image



Features: Histogram



Components of a Parametric Model

Several issues with scores:

- Not very interpretable (no bounded value)

We often want **probabilities**

- More interpretable
- Can relate to probabilistic view of machine learning

We use the **softmax** function to convert scores to probabilities

$$s = f(x, W) \quad \text{Scores}$$

$$= Wx$$

$$P(Y = k | X = x) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{Softmax Function}$$

- If we use the softmax function to convert scores to probabilities, the right loss function to use is **cross-entropy**
- Can be derived by looking at the distance between two probability distributions (output of model and ground truth)
- Can also be derived from a maximum likelihood estimation perspective

$$s = f(x, W) \quad \text{Scores}$$

$$P(Y = k|X = x) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{Softmax Function}$$

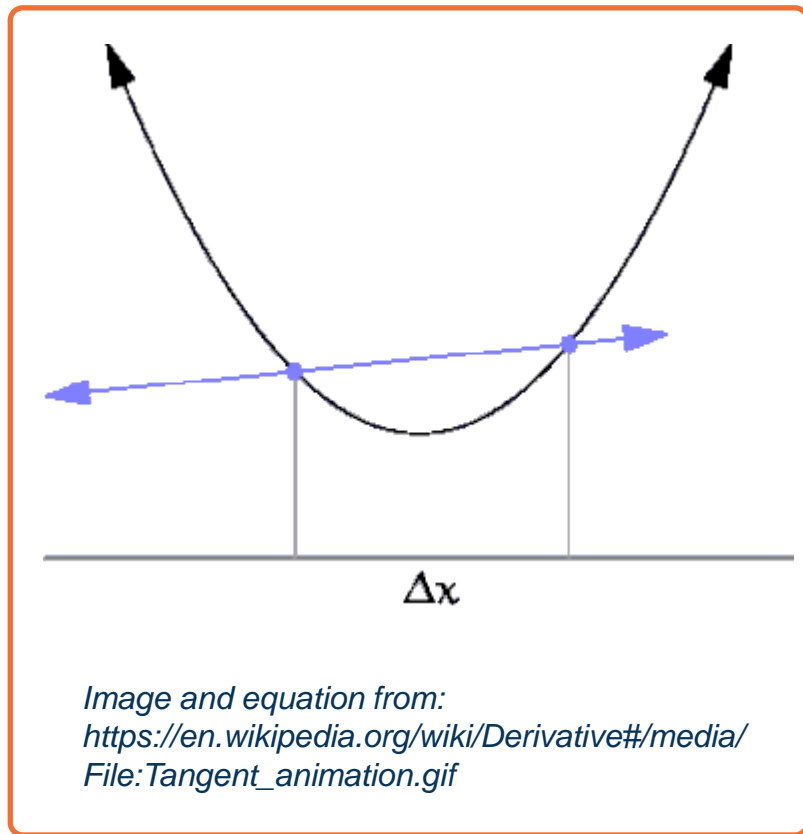
$$L_i = -\log P(Y = y_i|X = x_i)$$

Maximize log-prob of correct class =
Maximize the log likelihood
= Minimize the negative log likelihood

- 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

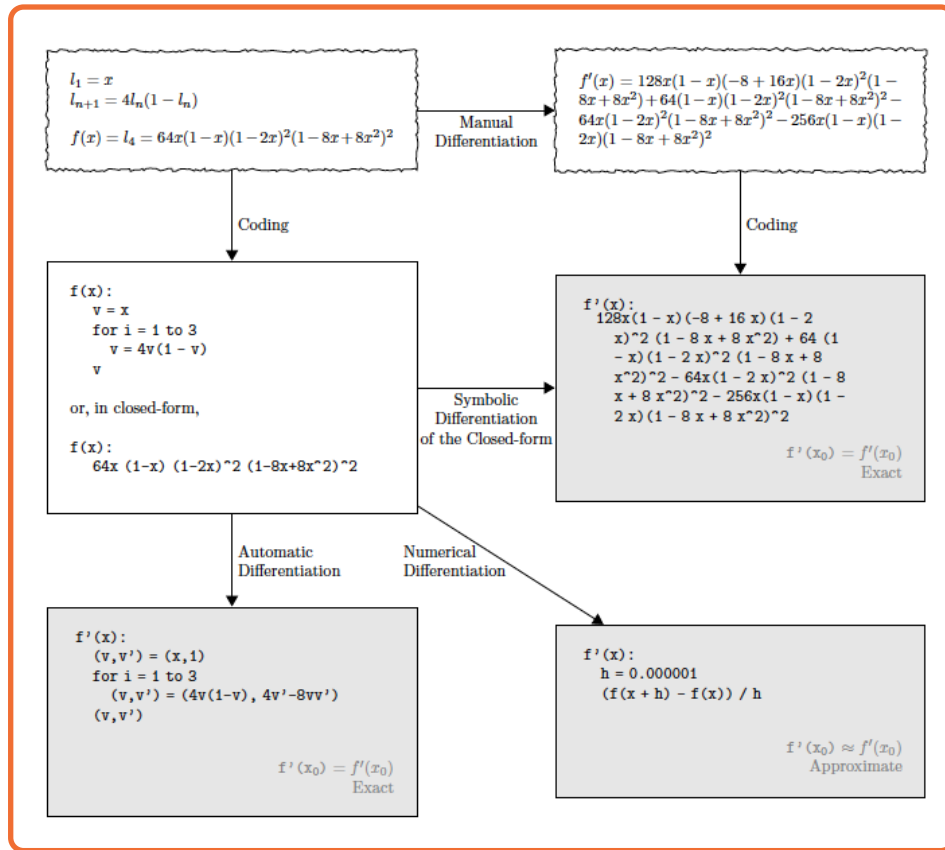


- Input: **Vector**
- Functional form of the model: **Softmax(Wx)**
- Performance measure to improve: **Cross-Entropy**
- Algorithm for finding best parameters: **Gradient Descent**
 - Compute $\frac{\partial L}{\partial w_i}$
 - Update Weights $w_i = w_i - \alpha \frac{\partial L}{\partial w_i}$

We know how to compute the **model output and loss function**

Several ways to compute $\frac{\partial L}{\partial w_i}$

- Manual differentiation
- Symbolic differentiation
- Numerical differentiation
- Automatic differentiation



current W:

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]

loss 1.25347

gradient dW:

[?,
?,
?,
?,
?,
?,
?,
?,
?,...]

current W:

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]

loss 1.25347

W + h (first dim):

[0.34 + **0.0001**,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]

loss 1.25322

gradient dW:

[?,
?,
?,
?,
?,
?,
?,
?,
?,...]

current W:

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]

loss 1.25347

W + h (first dim):

[0.34 + **0.0001**,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]

loss 1.25322

gradient dW:

[-2.5,
?,
?,


$$(1.25322 - 1.25347)/0.0001 = -2.5$$

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

?,
?,...]

current W:

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]

loss 1.25347

W + h (second dim):

[0.34,
-1.11 + **0.0001**,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]

loss 1.25353

gradient dW:

[-2.5,
?,
?,
?,
?,
?,
?,
?,
?,...]

current W:

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]

loss 1.25347

W + h (second dim):

[0.34,
-1.11 + **0.0001**,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]

loss 1.25353

gradient dW:

[-2.5,
0.6,
?,
?,

$$(1.25353 - 1.25347)/0.0001 = 0.6$$

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

?,...]

current W:

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]

loss 1.25347

W + h (third dim):

[0.34,
-1.11,
0.78 + **0.0001**,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]

loss 1.25347

gradient dW:

[-2.5,
0.6,
?,
?,
?,
?,
?,
?,
?,...]

current W:

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]

loss 1.25347

W + h (third dim):

[0.34,
-1.11,
0.78 + **0.0001**,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]

loss 1.25347

gradient dW:

[-2.5,
0.6,
0,
?,
?,...]


$$(1.25347 - 1.25347)/0.0001 = 0$$

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

?,...]

Numerical vs Analytic Gradients

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

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

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

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

This is called a **gradient check**.

For some functions, we can analytically derive the partial derivative

Example:

Function

$$f(\mathbf{w}, \mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i$$

(Assume \mathbf{w} and \mathbf{x}_i are column vectors, so same as $\mathbf{w} \cdot \mathbf{x}_i$)

Loss

$$\sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

Dataset: N examples (indexed by i)

Update Rule

$$\mathbf{w}_j \leftarrow \mathbf{w}_j + 2\alpha \sum_{i=1}^N \delta_i \mathbf{x}_{ij}$$

Derivation of Update Rule

$$L = \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

Gradient descent tells us we should update \mathbf{w} as follows to minimize L :

$$\mathbf{w}_j \leftarrow \mathbf{w}_j - \alpha \frac{\partial L}{\partial \mathbf{w}_j}$$

So what's $\frac{\partial L}{\partial \mathbf{w}_j}$?

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{w}_j} &= \sum_{i=1}^N \frac{\partial}{\partial \mathbf{w}_j} (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \\ &= \sum_{i=1}^N 2(y_i - \mathbf{w}^T \mathbf{x}_i) \frac{\partial}{\partial \mathbf{w}_j} (y_i - \mathbf{w}^T \mathbf{x}_i) \\ &= -2 \sum_{i=1}^N \delta_i \frac{\partial}{\partial \mathbf{w}_j} \mathbf{w}^T \mathbf{x}_i \\ &\quad \dots \text{where} \dots \\ &\quad \delta_i = y_i - \mathbf{w}^T \mathbf{x}_i \\ &= -2 \sum_{i=1}^N \delta_i \frac{\partial}{\partial \mathbf{w}_j} \sum_{k=1}^N \mathbf{w}_k \mathbf{x}_{ik} \\ &= -2 \sum_{i=1}^N \delta_i \mathbf{x}_{ij} \end{aligned}$$

If we add a **non-linearity (sigmoid)**, derivation is more complex

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

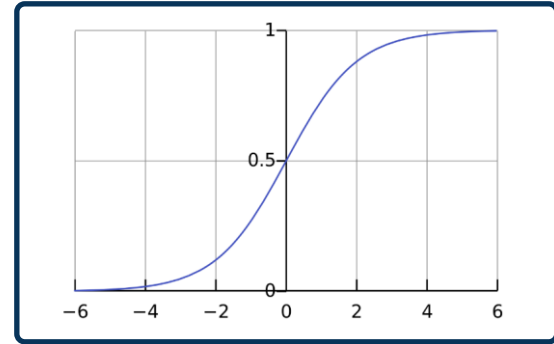
First, one can derive that: $\sigma'(x) = \sigma(x)(1 - \sigma(x))$

$$\mathbf{f}(\mathbf{x}) = \sigma\left(\sum_k w_k x_k\right)$$

$$L = \sum_i \left(y_i - \sigma\left(\sum_k w_k x_{ik}\right) \right)^2$$

$$\begin{aligned} \frac{\partial L}{\partial w_j} &= \sum_i 2 \left(y_i - \sigma\left(\sum_k w_k x_{ik}\right) \right) \left(-\frac{\partial}{\partial w_j} \sigma\left(\sum_k w_k x_{ik}\right) \right) \\ &= \sum_i -2 \left(y_i - \sigma\left(\sum_k w_k x_{ik}\right) \right) \sigma'\left(\sum_k w_k x_{ik}\right) \frac{\partial}{\partial w_j} \sum_k w_k x_{ik} \\ &= \sum_i -2 \delta_i \sigma(d_i) (1 - \sigma(d_i)) x_{ij} \end{aligned}$$

$$\text{where } \delta_i = y_i - \mathbf{f}(x_i) \quad d_i = \sum_k w_k x_{ik}$$



The sigmoid perception update rule:

$$w_j \leftarrow w_j + 2\alpha \sum_{k=1}^N \delta_i \sigma_i (1 - \sigma_i) x_{ij}$$

where $\sigma_i = \sigma\left(\sum_{j=1}^d w_j x_{ij}\right)$

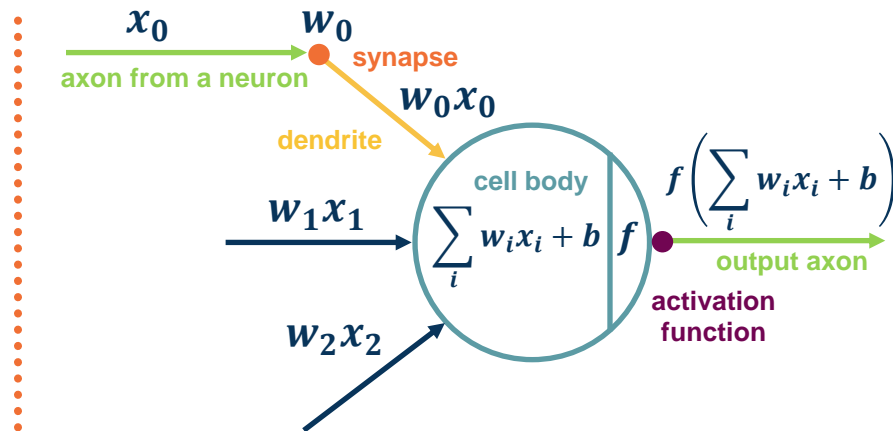
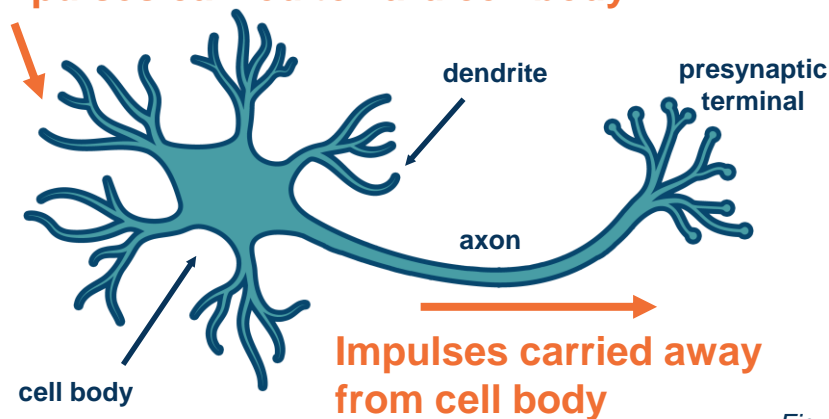
$$\delta_i = y_i - \sigma_i$$

Neural Network View of a Linear Classifier

A simple **neural network** has similar structure as our linear classifier:

- A neuron takes input (firings) from other neurons (-> **input to linear classifier**)
- The inputs are summed in a weighted manner (-> **weighted sum**)
 - Learning is through a modification of the weights
- If it receives enough input, it “fires” (threshold or if weighted sum plus bias is high enough)

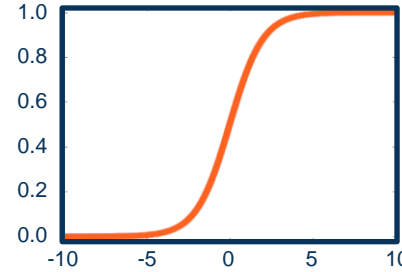
Impulses carried toward cell body



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

Origins of the Term Neural Network

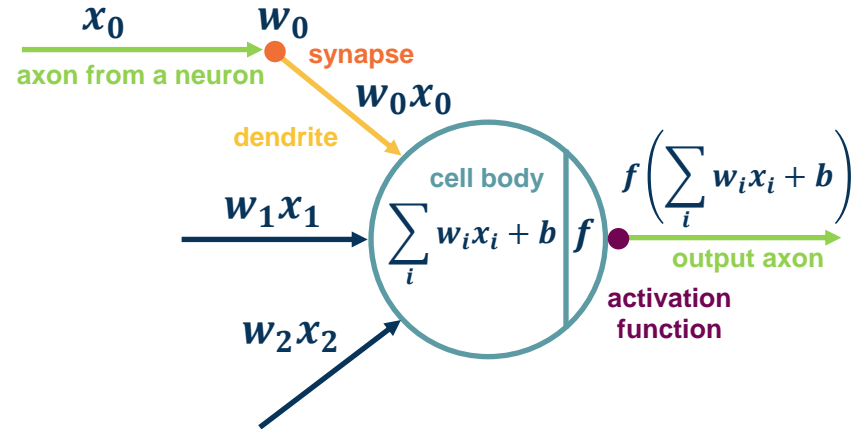
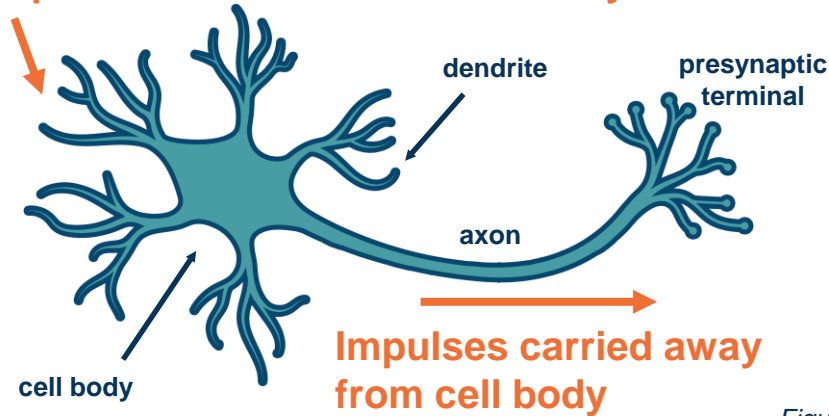
As we did before, the output of a neuron can be modulated by a non-linear function (e.g. sigmoid)



**Sigmoid
Activation
Function**

$$\frac{1}{1 + e^{-x}}$$

Impulses carried toward cell body



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

Adding Non-Linearities

We can have **multiple** neurons connected to the same input

Corresponds to a multi-class classifier

- Each output node outputs the score for a class

$$f(x, W) = \sigma(Wx + b) \quad \begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1m} & b1 \\ w_{21} & w_{22} & \cdots & w_{2m} & b2 \\ w_{31} & w_{32} & \cdots & w_{3m} & b3 \end{bmatrix}$$

- Often called fully connected layers
 - Also called a linear *projection layer*

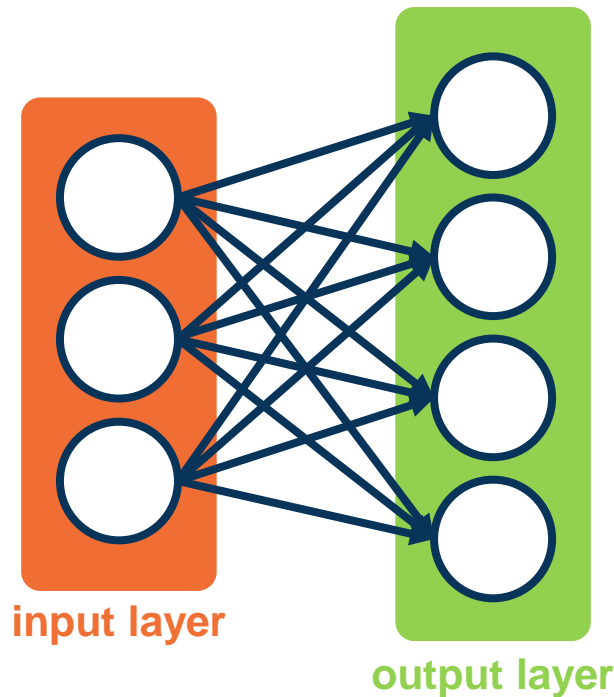


Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

- Each input/output is a **neuron (node)**
- A linear classifier (+ optional non-linearity) is called a **fully connected** layer
- Connections are represented as **edges**
- Output of a particular neuron is referred to as **activation**
- This will be expanded as we view computation in a neural network as a **graph**

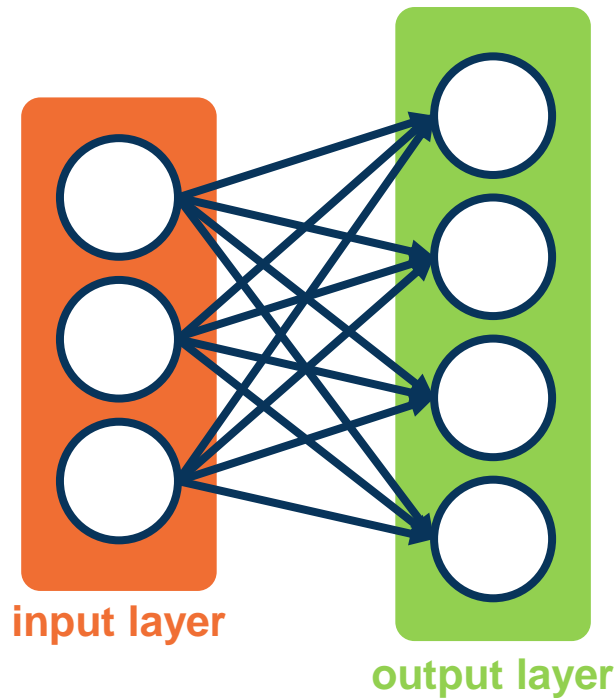


Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

We can **stack** multiple layers together

- Input to second layer is output of first layer

Called a **2-layered neural network** (input is not counted)

Because the middle layer is neither input or output, and we don't know what their values represent, we call them **hidden** layers

- We will see that they end up learning effective features

This **increases** the representational power of the function!

- Two layered networks can represent any continuous function

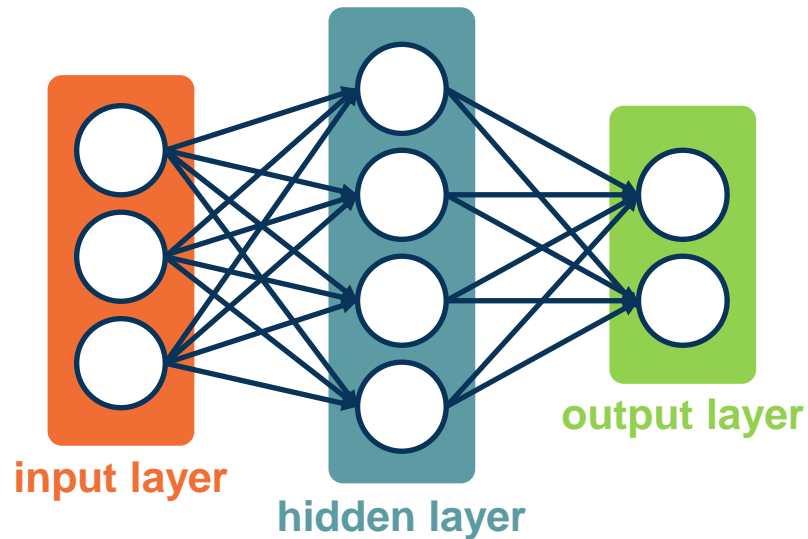


Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

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)

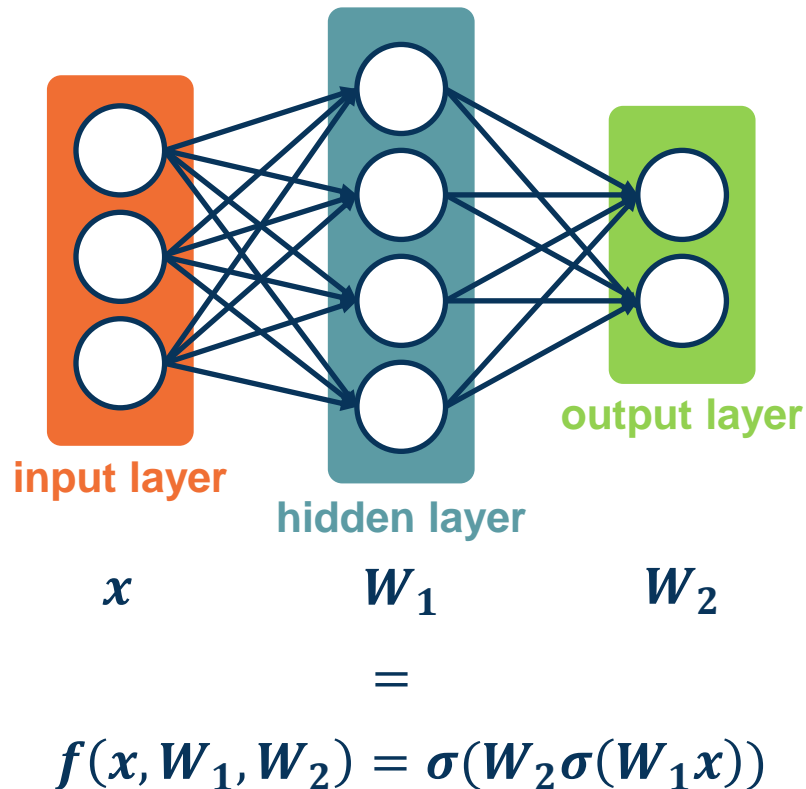
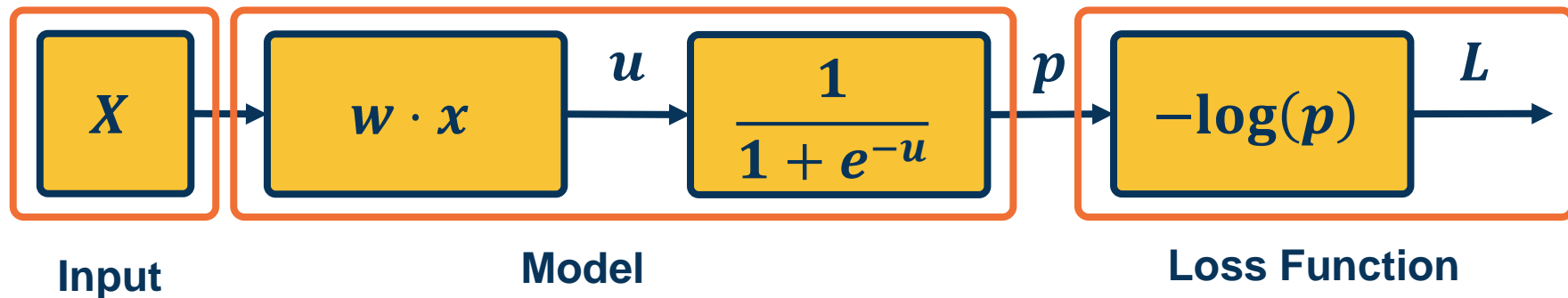


Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

A **classifier** can be broken down into:

- Input
- A function of the input
- A loss function

It's all just one function that can be **decomposed** into building blocks



What Does a Linear Classifier Consist of?

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**:

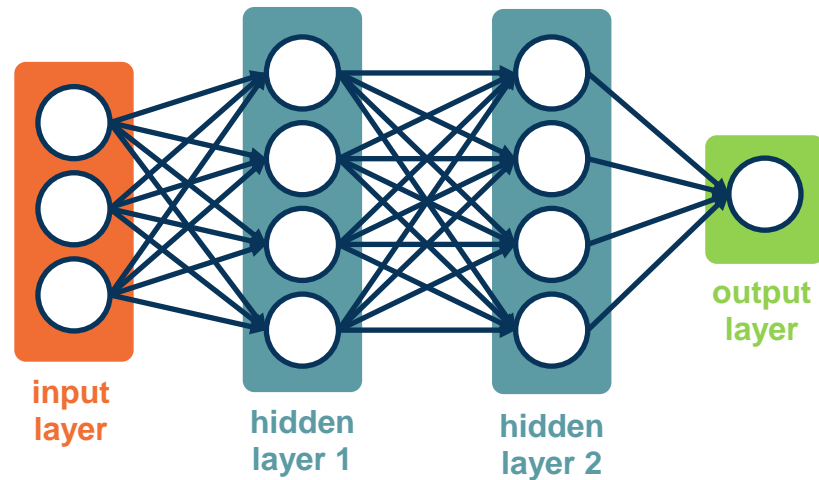
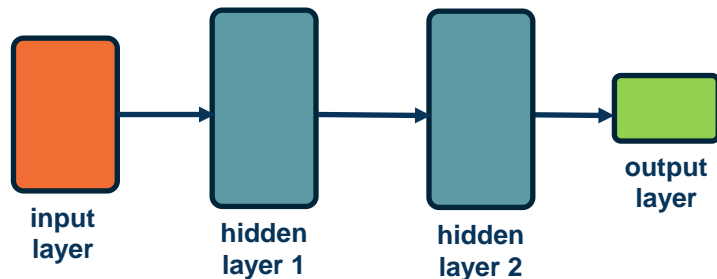


Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Adding More Layers!

Computation Graphs

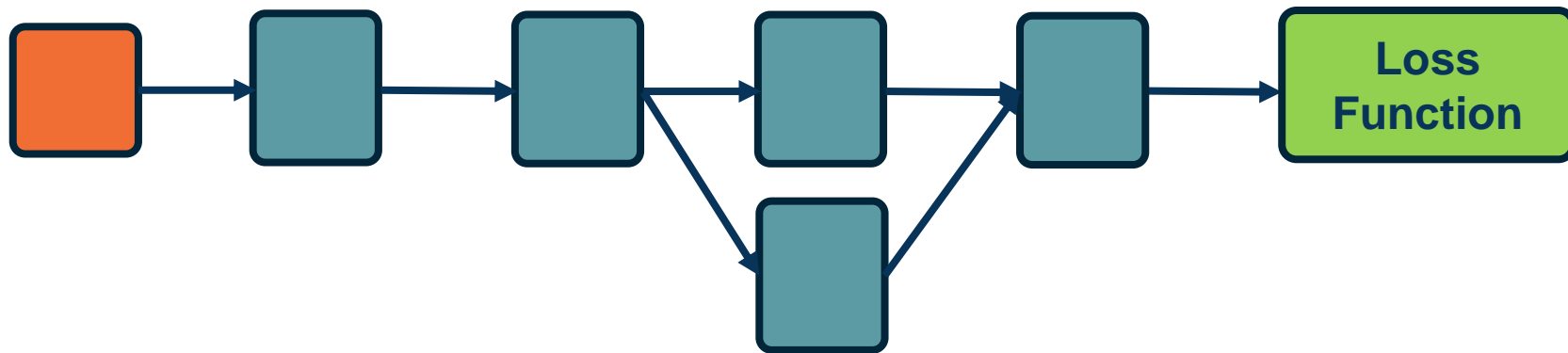
Functions can be made **arbitrarily complex** (subject to memory and computational limits), e.g.:

$$f(x, W) = \sigma(W_5 \sigma(W_4 \sigma(W_3 \sigma(W_2 \sigma(W_1 x))))$$

We can use **any type of differentiable function (layer)** we want!

✦ At the end, **add the loss function**

Composition can have **some structure**



Adding Even More Layers

The world is **compositional**!

We want our **model** to reflect this

Empirical and theoretical evidence that it makes **learning complex functions easier**

Note that **prior state of art engineered features** often had this compositionality as well

VISION

pixels → edge → texton → motif → part → object

SPEECH

sample → spectral band → formant → motif → phone → word

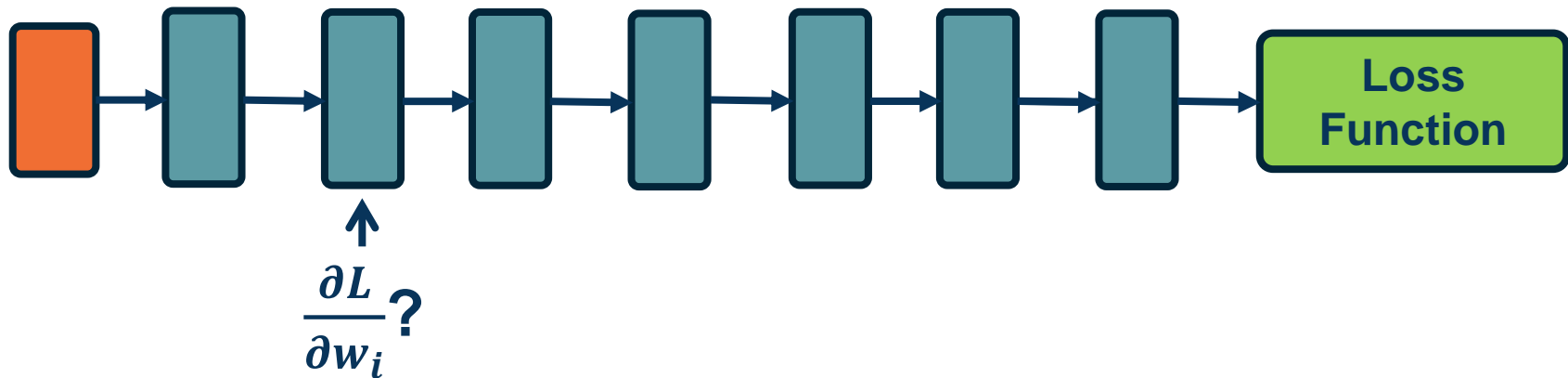
NLP

character → word → NP/VP/.. → clause → sentence → story

Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

🟡 **Pixels -> edges -> object parts -> objects**

- 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

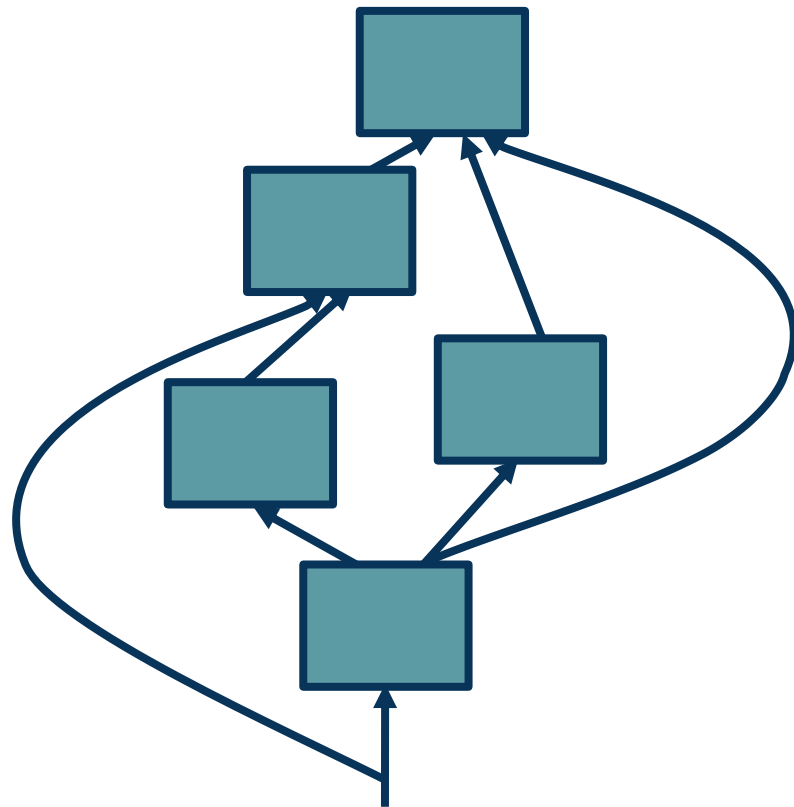


To develop a general algorithm for this, we will view the function as a **computation graph**

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

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

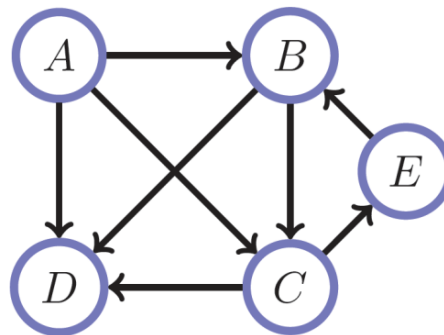
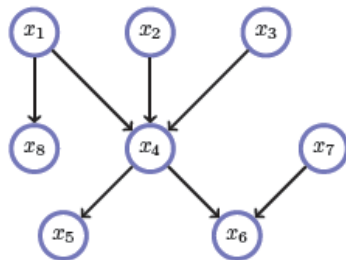
A **training algorithm** will then process this graph, **one module at a time**



Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

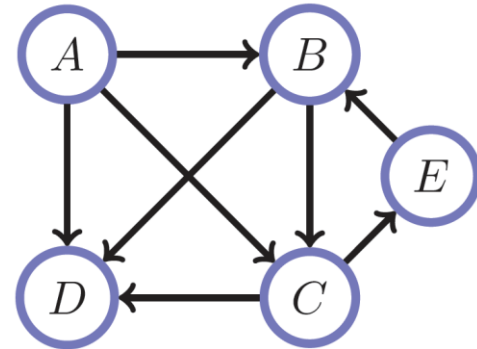
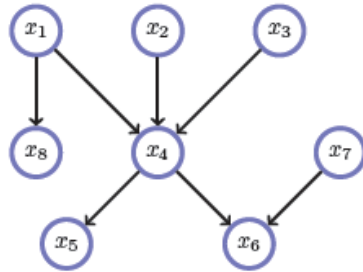
Directed Acyclic Graphs (DAGs)

- Exactly what the name suggests
 - Directed edges
 - No (directed) cycles
 - Underlying undirected cycles okay

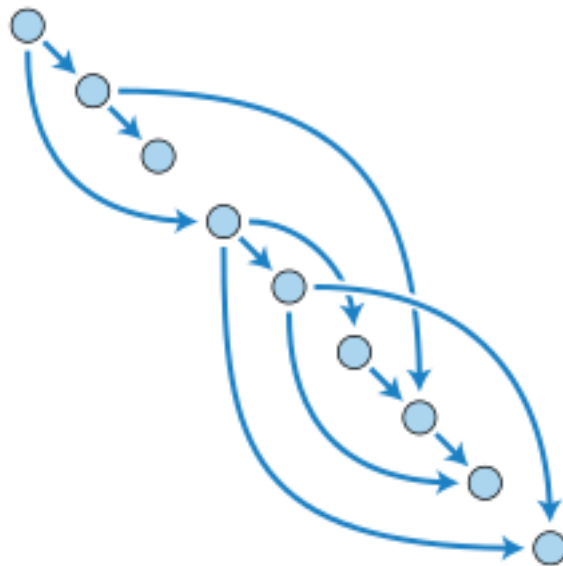


Directed Acyclic Graphs (DAGs)

- Concept
 - Topological Ordering



Directed Acyclic Graphs (DAGs)



Backpropagation

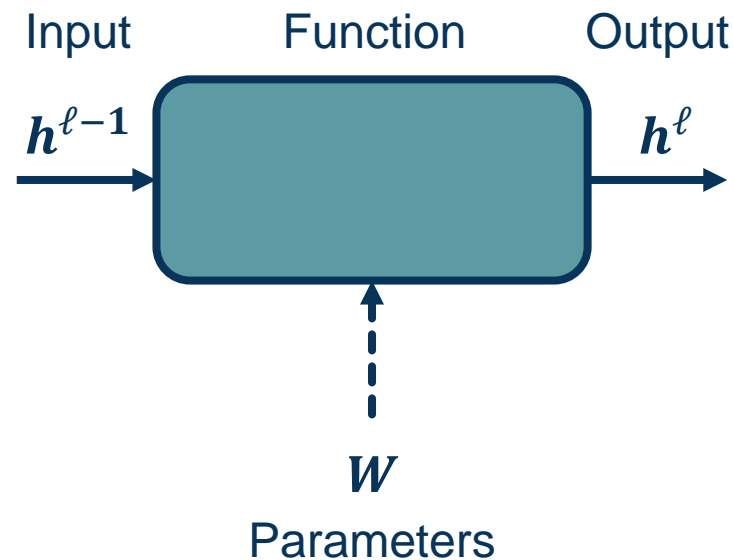
Given this computation graph, the training algorithm will:

- Calculate the current model's outputs (called the **forward pass**)
- Calculate the gradients for each module (called the **backward pass**)

Backward pass is a recursive algorithm that:

- Starts at **loss function** where we know how to calculate the gradients
- Progresses back through the modules
- Ends in the **input layer** where we do not need gradients (no parameters)

This algorithm is called **backpropagation**



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**



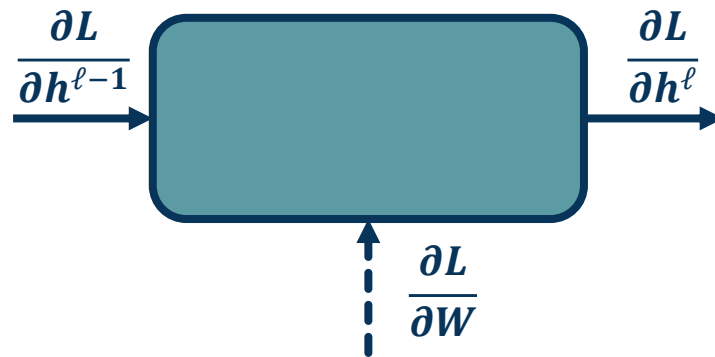
Note that we must store the **intermediate outputs of all layers!**

- ⬢ 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

In the **backward pass**, we seek to calculate the gradients of the loss with respect to the module's parameters

- Assume that we have the gradient of the loss with respect to the **module's outputs** (given to us by upstream module)
- We will also pass the gradient of the loss with respect to the **module's inputs**
 - This is not required for update the module's weights, but passes the gradients back to the previous module

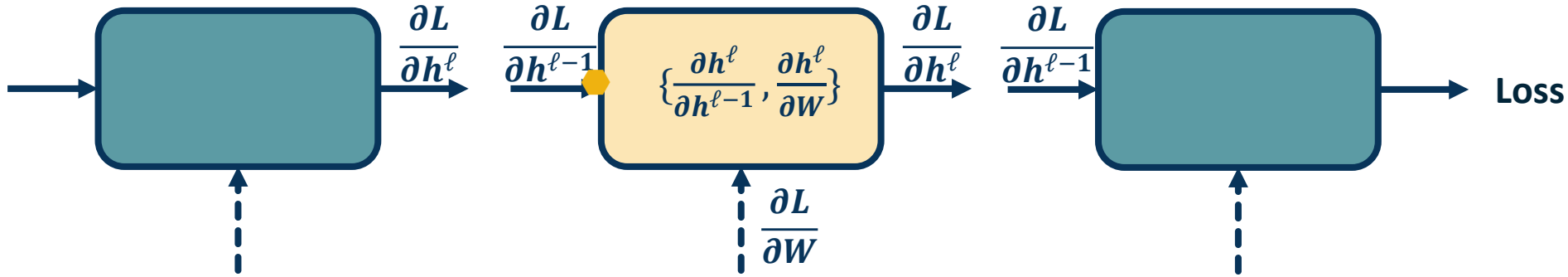


Problem:

- We are given: $\frac{\partial L}{\partial h^{\ell}}$
- We can compute local gradients: $\left\{ \frac{\partial h^{\ell}}{\partial h^{\ell-1}}, \frac{\partial h^{\ell}}{\partial W} \right\}$
- Compute: $\left\{ \frac{\partial L}{\partial h^{\ell-1}}, \frac{\partial L}{\partial W} \right\}$

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

- We can compute **local gradients**: $\{\frac{\partial h^\ell}{\partial h^{\ell-1}}, \frac{\partial h^\ell}{\partial W}\}$
- This is just the **derivative of our function** with respect to its parameters and inputs!

Example: If $h^\ell = Wh^{\ell-1}$

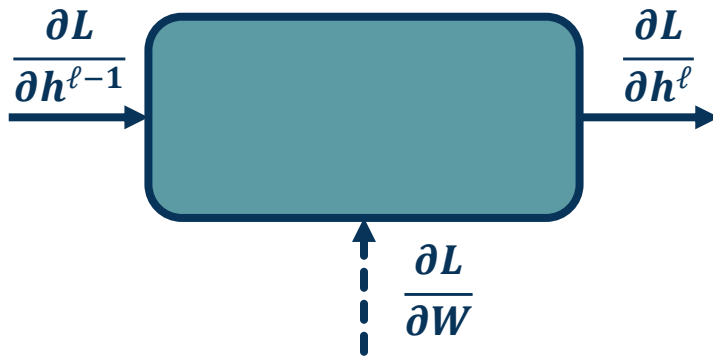
$$\text{then } \frac{\partial h^\ell}{\partial h^{\ell-1}} = W$$

$$\text{and } \frac{\partial h_i^\ell}{\partial w_i} = h^{\ell-1,T}$$

- We will use the **chain rule** to compute: $\{\frac{\partial L}{\partial h^{\ell-1}}, \frac{\partial L}{\partial W}\}$

- **Gradient of loss w.r.t. inputs:** $\frac{\partial L}{\partial h^{\ell-1}} = \frac{\partial L}{\partial h^{\ell}} \frac{\partial h^{\ell}}{\partial h^{\ell-1}}$ └─ Given by upstream module (**upstream gradient**)

- **Gradient of loss w.r.t. weights:** $\frac{\partial L}{\partial W} = \frac{\partial L}{\partial h^{\ell}} \frac{\partial h^{\ell}}{\partial W}$



Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

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

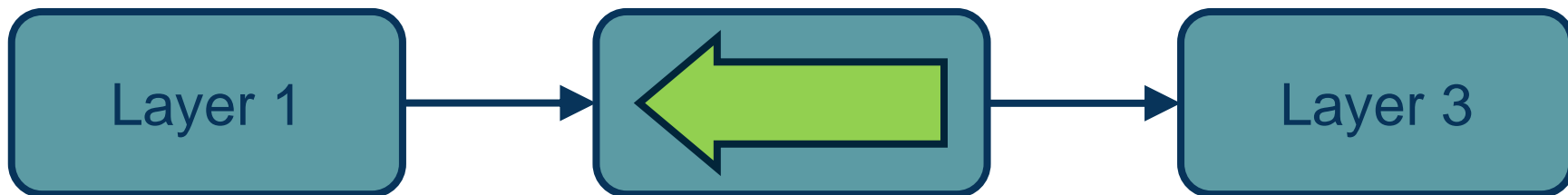
Step 2: Compute Gradients wrt parameters: **Backward Pass**



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

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!



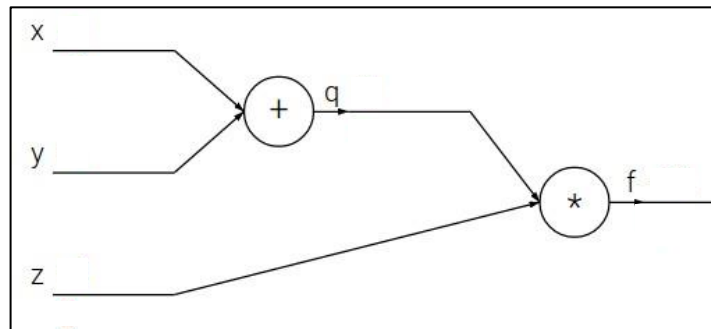
Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

Backpropagation: a simple example

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

Backpropagation: a simple example

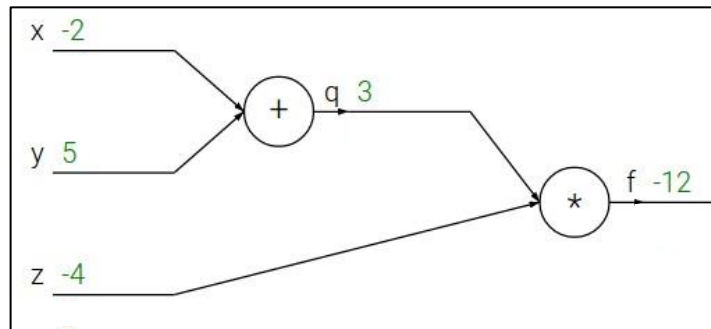
$$f(x, y, z) = (x + y)z$$



Backpropagation: a simple example

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

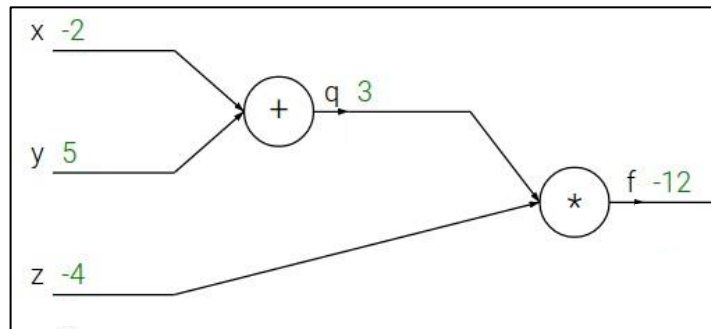
e.g. $x = -2$, $y = 5$, $z = -4$



Backpropagation: a simple example

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

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



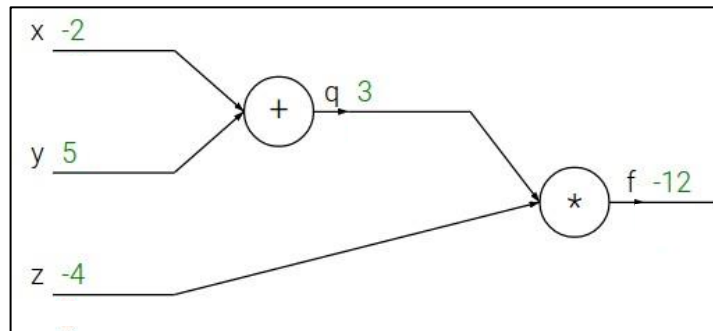
Want: $\frac{\partial f}{\partial x}$, $\frac{\partial f}{\partial y}$, $\frac{\partial f}{\partial z}$

Backpropagation: a simple example

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

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

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



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

Backpropagation: a simple example

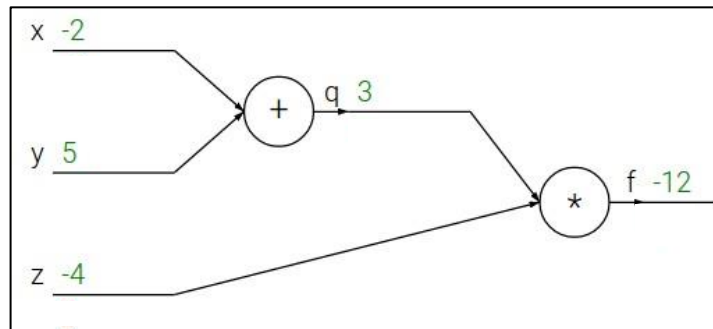
$$f(x, y, z) = (x + y)z$$

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



Backpropagation: a simple example

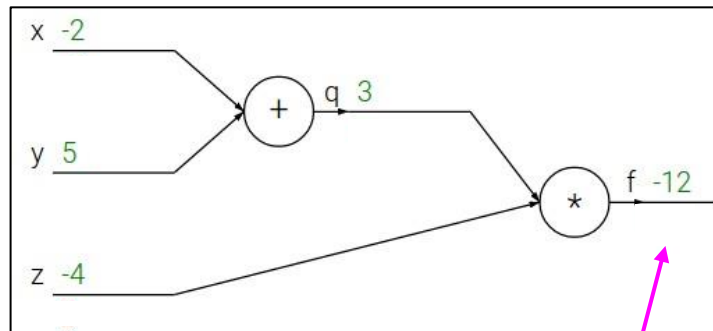
$$f(x, y, z) = (x + y)z$$

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



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

Backpropagation: a simple example

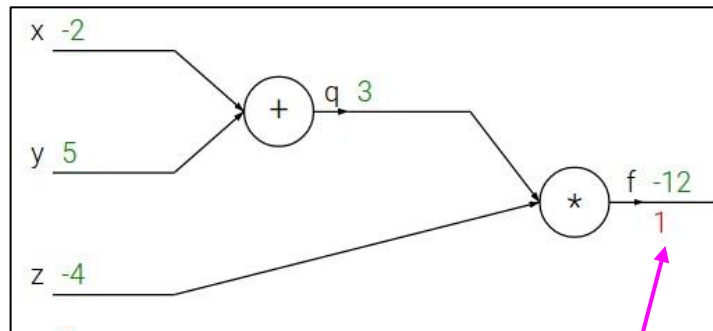
$$f(x, y, z) = (x + y)z$$

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



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

Backpropagation: a simple example

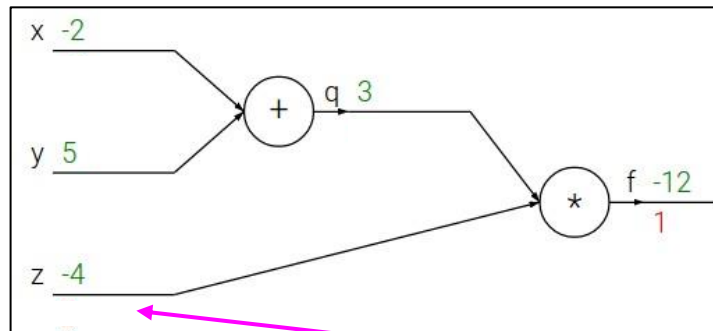
$$f(x, y, z) = (x + y)z$$

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



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

Backpropagation: a simple example

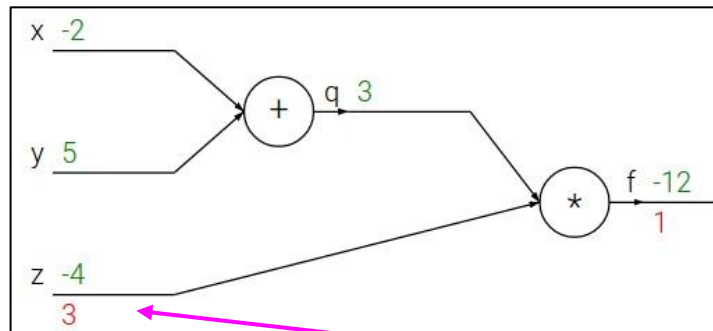
$$f(x, y, z) = (x + y)z$$

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



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

Backpropagation: a simple example

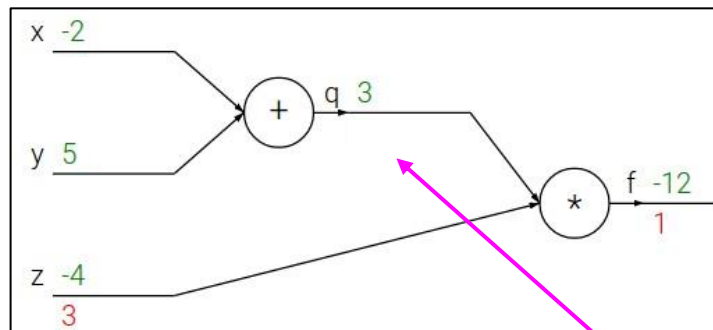
$$f(x, y, z) = (x + y)z$$

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



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

Backpropagation: a simple example

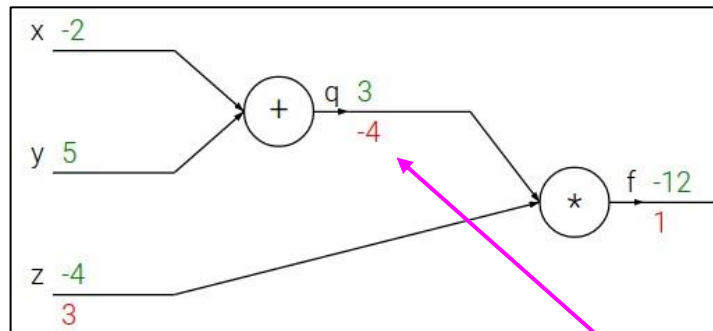
$$f(x, y, z) = (x + y)z$$

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



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

Backpropagation: a simple example

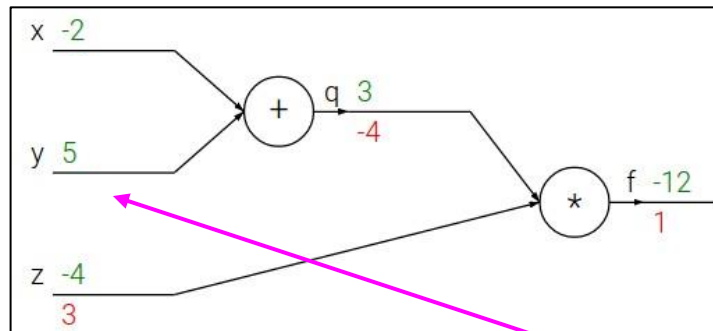
$$f(x, y, z) = (x + y)z$$

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



Chain rule:

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

Upstream
gradient

Local
gradient

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

Backpropagation: a simple example

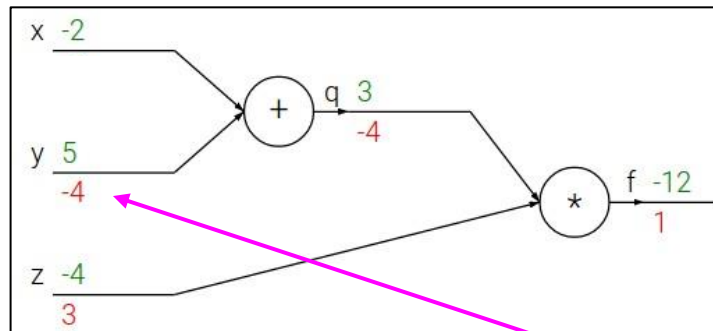
$$f(x, y, z) = (x + y)z$$

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



Chain rule:

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

Upstream
gradient

Local
gradient

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

Backpropagation: a simple example

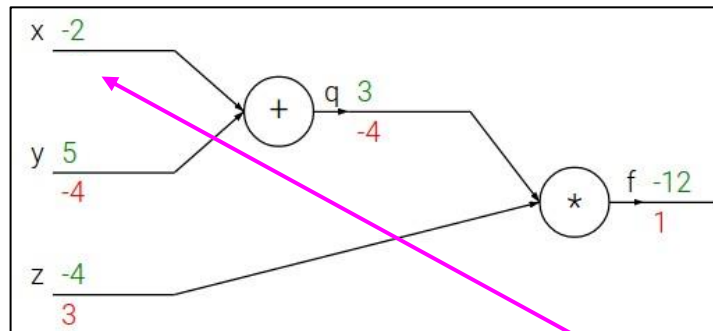
$$f(x, y, z) = (x + y)z$$

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



Chain rule:

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

Upstream
gradient

Local
gradient

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

Backpropagation: a simple example

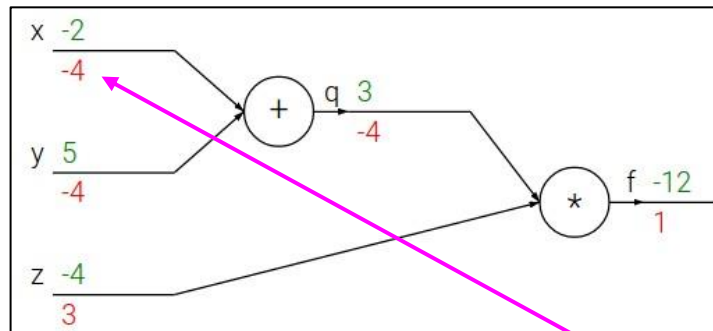
$$f(x, y, z) = (x + y)z$$

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



Chain rule:

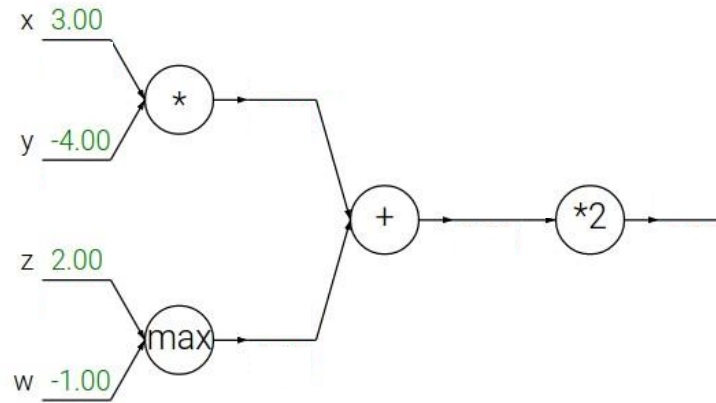
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

Upstream
gradient

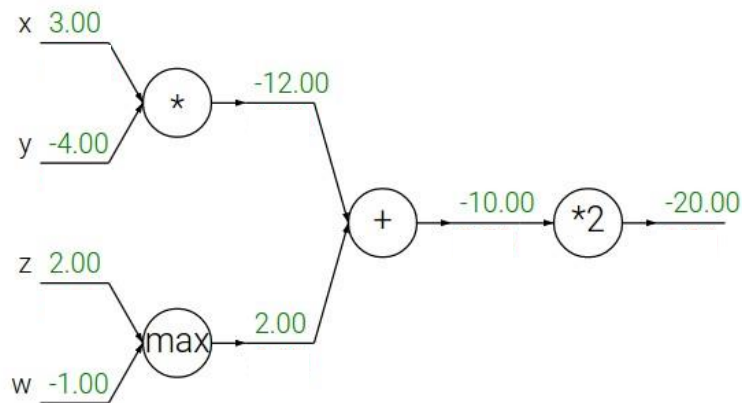
Local
gradient

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

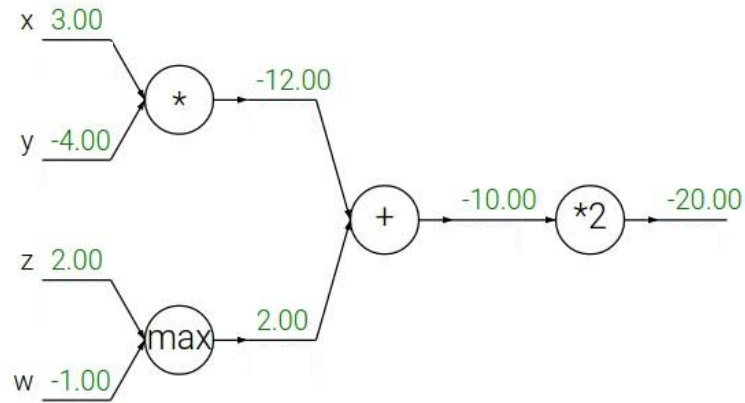
Backpropagation: a simple example



Backpropagation: a simple example

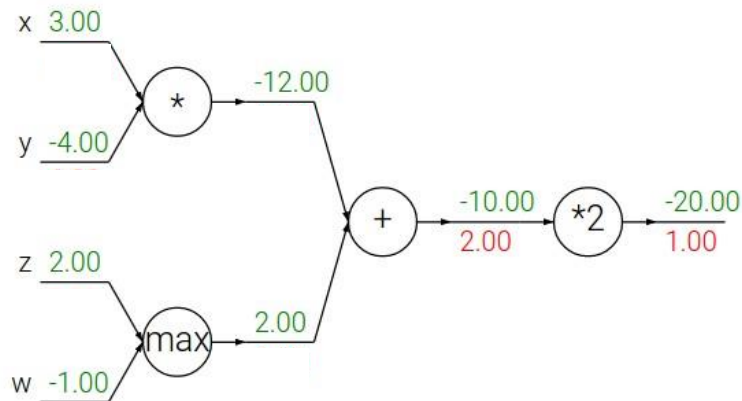


Patterns in backward flow



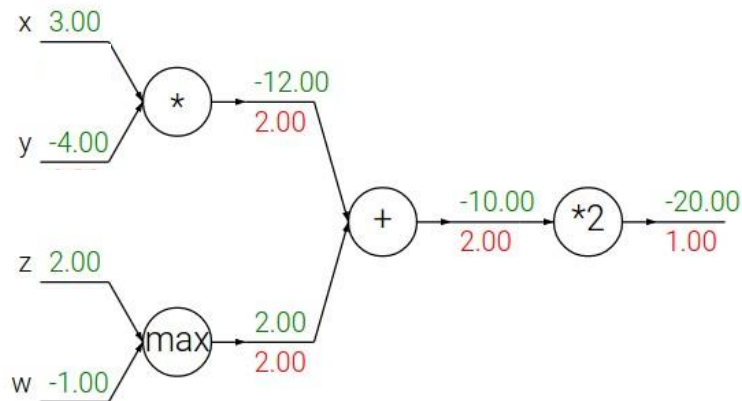
Patterns in backward flow

Q: What is an **add** gate?



Patterns in backward flow

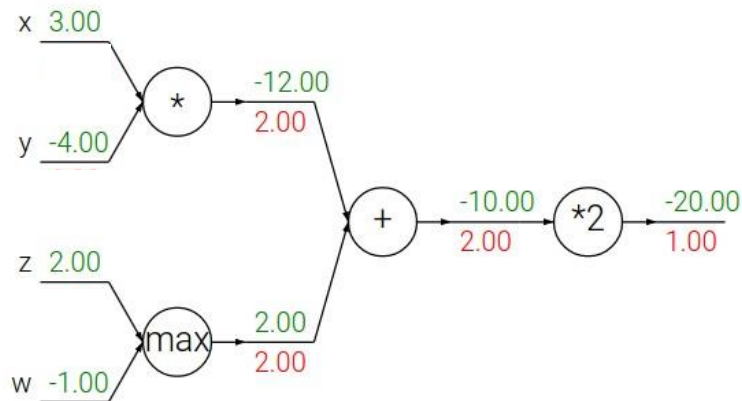
add gate: gradient distributor



Patterns in backward flow

add gate: gradient distributor

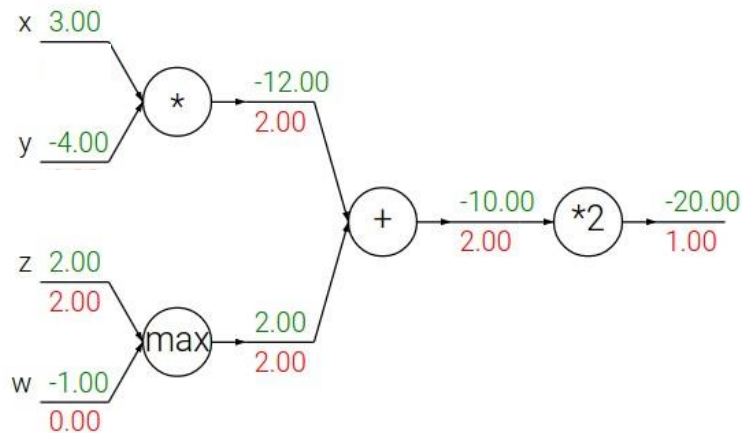
Q: What is a **max** gate?



Patterns in backward flow

add gate: gradient distributor

max gate: gradient router

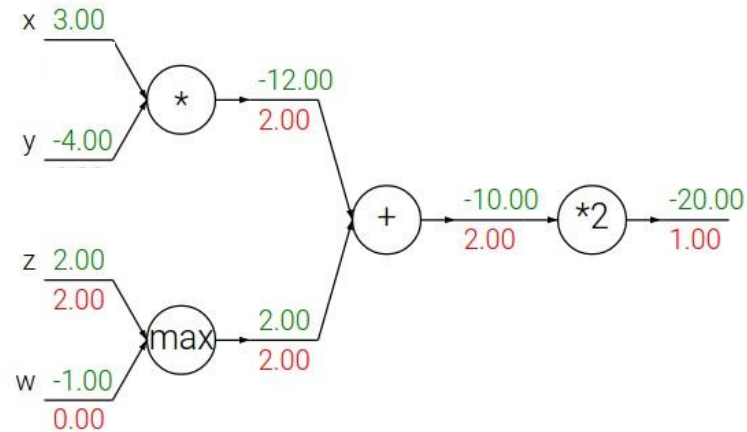


Patterns in backward flow

add gate: gradient distributor

max gate: gradient router

Q: What is a **mul** gate?

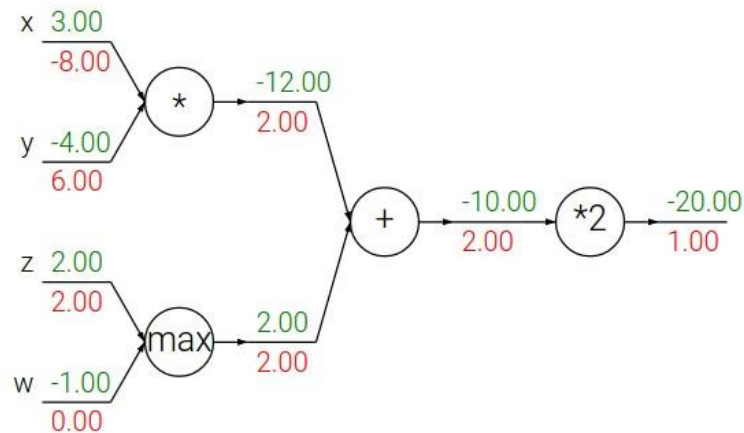


Patterns in backward flow

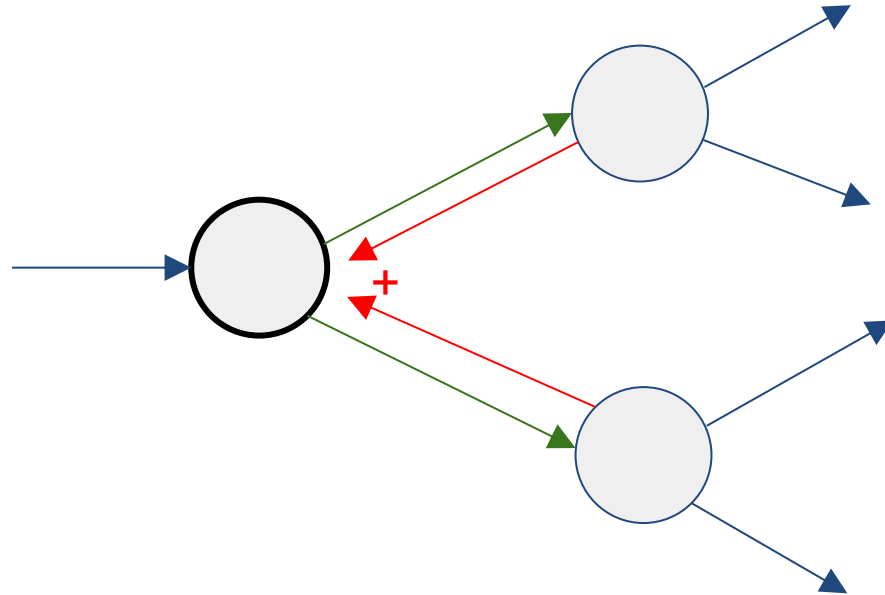
add gate: gradient distributor

max gate: gradient router

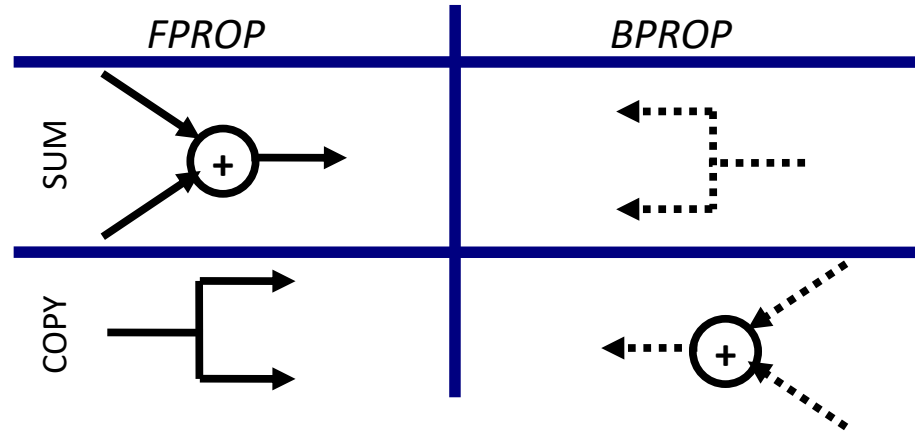
mul gate: gradient switcher



Gradients add at branches



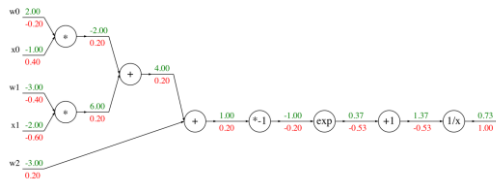
Duality in Fprop and Bprop



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

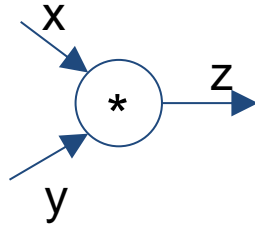
Modularized implementation: forward / backward API



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
```

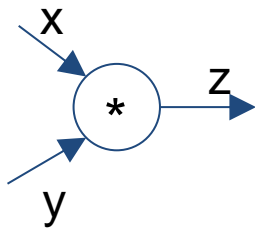
```
        return [dx, dy]
```

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

Diagram showing the forward and backward passes of a multiplication gate. The forward pass calculates $z = x * y$. The backward pass calculates the gradients $\frac{\partial L}{\partial x}$ and $\frac{\partial L}{\partial y}$ from the upstream gradient $\frac{\partial L}{\partial z}$. Arrows indicate the flow of gradients from the backward pass output back to the inputs of the forward pass.

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

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 BGDneeload_hdfs_fix		Latest commit e6b7a71 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	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			

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

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

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

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

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

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

$$\begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1m} & b1 \\ w_{21} & w_{22} & \cdots & w_{2m} & b2 \\ w_{31} & w_{32} & \cdots & w_{3m} & b3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \\ 1 \end{bmatrix}$$

W

x

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

Where c is number of classes

d is dimensionality of input

Closer Look at a Linear Classifier

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 $\mathbf{M} \in \mathbb{R}^{k \times \ell}$

	S $\begin{bmatrix} \end{bmatrix}$	V $\begin{bmatrix} \end{bmatrix}$	M $\begin{bmatrix} \end{bmatrix}$
S	$\frac{\partial s_1}{\partial s_2}$ $\begin{bmatrix} \end{bmatrix}$	$\frac{\partial s}{\partial v}$ $\begin{bmatrix} \end{bmatrix}$	$\frac{\partial s}{\partial M}$ $\begin{bmatrix} \end{bmatrix}$
V	$\frac{\partial v}{\partial s}$ $\begin{bmatrix} \end{bmatrix}$	$\frac{\partial v_1}{\partial v_2}$ $\begin{bmatrix} \end{bmatrix}$	Tensors
M	$\frac{\partial M}{\partial s}$ $\begin{bmatrix} \end{bmatrix}$		

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 $\mathbf{M} \in \mathbb{R}^{k \times \ell}$

- 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)

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

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

Conventions:

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

$$\begin{array}{c} \text{Col } j \\ \text{Row } i \end{array} \left[\begin{array}{ccccc} \frac{\partial v_1^1}{\partial v_1^2} & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \frac{\partial v_i^1}{\partial v_1^2} & \dots & \frac{\partial v_i^1}{\partial v_j^2} & \dots & \frac{\partial v_i^1}{\partial v_{m_2}^2} \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \end{array} \right]_{m_1 \times m_2}$$

- 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:

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

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

Example 1:

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

Example 2:

$$\begin{aligned} y &= w^T x = \sum_k w_k x_k \\ \frac{\partial y}{\partial x} &= \left[\frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_m} \right] \\ &= [w_1, \dots, w_m] \quad \text{because} \quad \frac{\partial (\sum_k w_k x_k)}{\partial x_i} = w_i \\ &= w^T \end{aligned}$$

Example 3:

$$y = Wx$$

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

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

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 

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

- 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**