# Lecture 4:
# Neural Networks and
# Backpropagation

# Announcements: Assignment 1

**Assignment 1** due **Fri 4/16 at 11:59pm**

# Administrative: Project Proposal

Due **Mon 4/19**

TA expertise are posted on the webpage.

(http://cs231n.stanford.edu/office_hours.html)
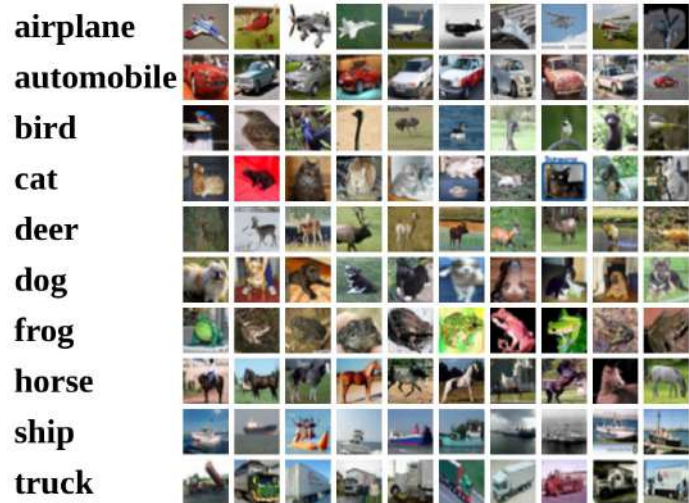
# Administrative: Discussion Section

Discussion section tomorrow:
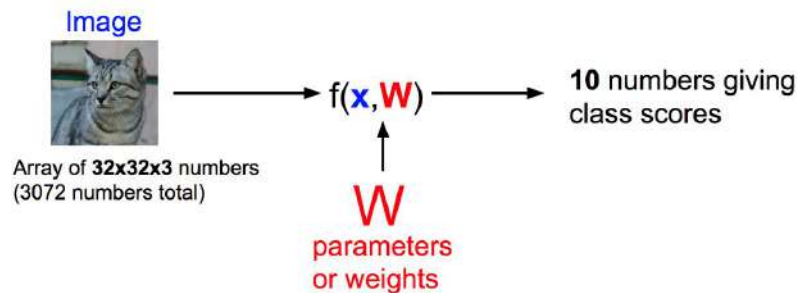
Backpropagation

# Administrative: Midterm Updates

- Tues, **May 4** and is worth **15%** of your grade.
- available for **24 hours** on Gradescope from May 4, **12PM** PDT to May 5, 11:59 AM PDT.
- **3-hour** consecutive timeframe
- Exam will be designed for **1.5 hours**.
- Open book and open internet but no collaboration
- Only make private posts during those 24 hours
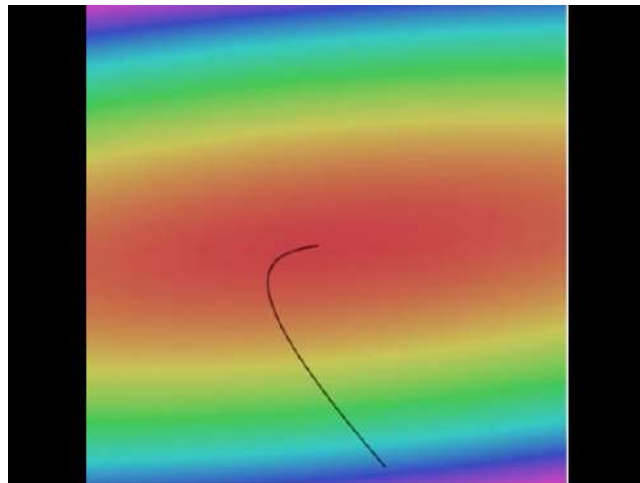
# Recap: from last time



$$f(x,W) = Wx + b$$



**Image**

Array of **32x32x3** numbers
(3072 numbers total)

$f(\mathbf{x},\mathbf{W})$

**10** numbers giving class scores

**W**
parameters or weights

# Recap: loss functions

$$s = f(x; W) = Wx$$ Linear score function

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

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + \lambda \sum_k W_k^2$$ data loss + regularization

# Finding the best W: Optimize with Gradient Descent





```
# Vanilla Gradient Descent

while True:
  weights_grad = evaluate_gradient(loss_fun, data, weights)
  weights += - step_size * weights_grad # perform parameter update
```

# Gradient descent

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

**Numerical gradient**: slow :(, approximate :(, easy to write :)
**Analytic gradient**: fast :), exact :), error-prone :(

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

# Stochastic Gradient Descent (SGD)

$$L(W) = \frac{1}{N} \sum_{i=1}^{N} L_i(x_i, y_i, W) + \lambda R(W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^{N} \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

Full sum expensive when N is large!

Approximate sum using a **minibatch** of examples
32 / 64 / 128 common

```
# Vanilla Minibatch Gradient Descent

while True:
    data_batch = sample_training_data(data, 256) # sample 256 examples
    weights_grad = evaluate_gradient(loss_fun, data_batch, weights)
    weights += - step_size * weights_grad # perform parameter update
```

# What we are going to discuss today!

$$s = f(x; W) = Wx$$ Linear score function

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

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + \lambda \sum_k W_k^2$$ data loss + regularization

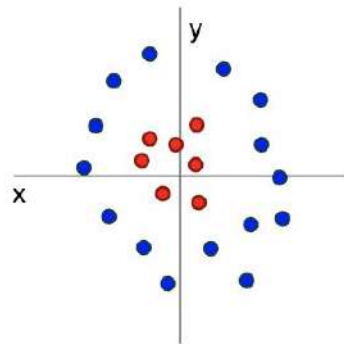How to find the best W? $\boxed{\nabla_W L}$

# Problem: Linear Classifiers are not very powerful

## Visual Viewpoint



Linear classifiers learn
one template per class

## Geometric Viewpoint



Linear classifiers
can only draw linear
decision boundaries

# Pixel Features



f(x) = Wx

Class scores

plane    car    bird    cat    deer

dog    frog    horse    ship    truck

# Image Features



$$f(x) = Wx$$

Feature Representation
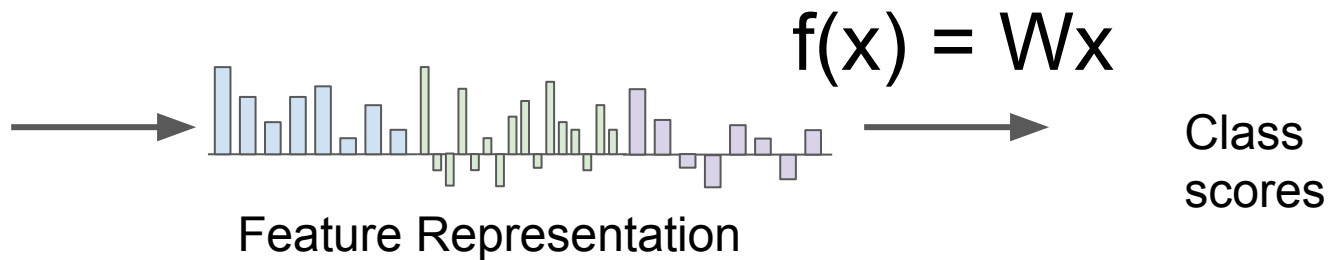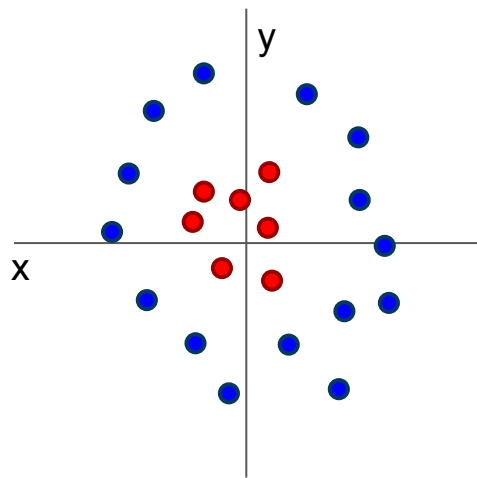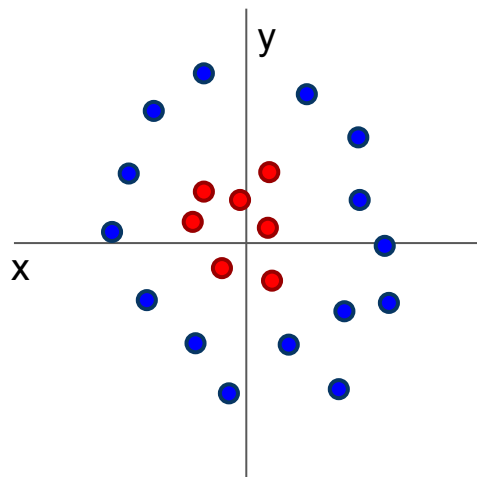
Class scores

# Image Features: Motivation



Cannot separate red
and blue points with
linear classifier

# Image Features: Motivation



$$f(x, y) = (r(x, y), \theta(x, y))$$
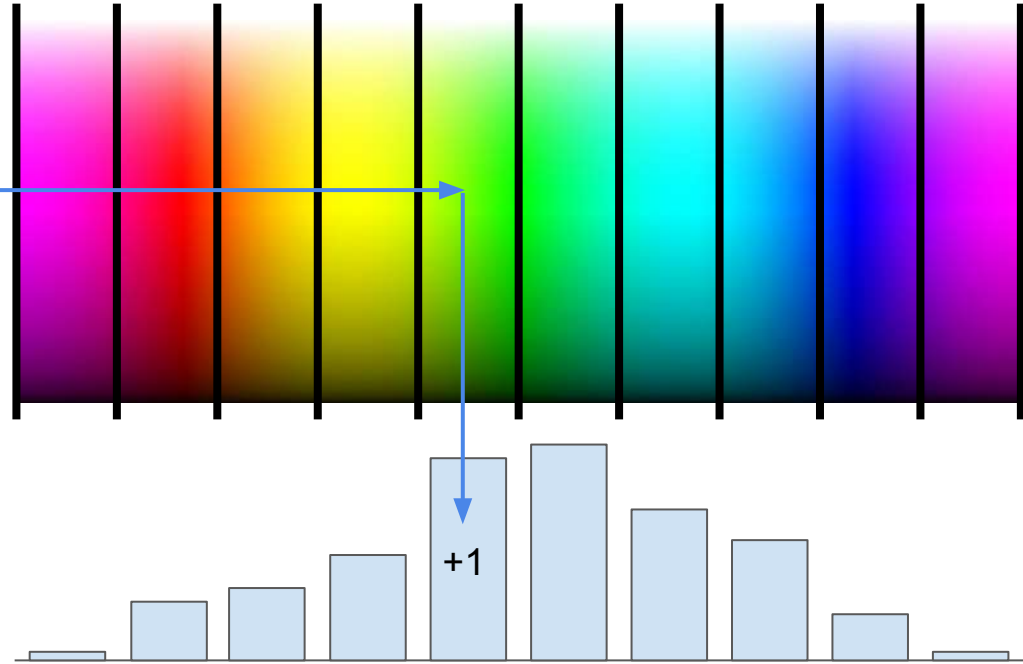
Cannot separate red
and blue points with
linear classifier

After applying feature
transform, points can
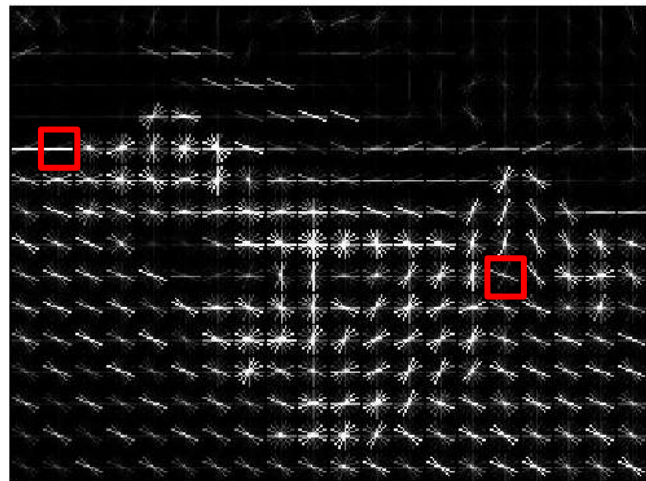be separated by linear
classifier

# Example: Color Histogram



+1

# Example: Histogram of Oriented Gradients (HoG)



Divide image into 8x8 pixel regions
Within each region quantize edge
direction into 9 bins

Example: 320x240 image gets divided
into 40x30 bins; in each bin there are
9 numbers so feature vector has
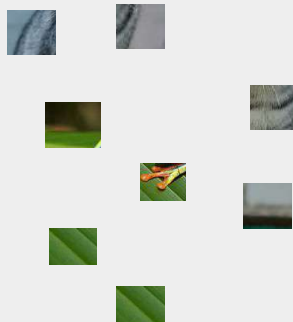30*40*9 = 10,800 numbers

Lowe, "Object recognition from local scale-invariant features", ICCV 1999
Dalal and Triggs, "Histograms of oriented gradients for human detection," CVPR 2005

# Example: Bag of Words



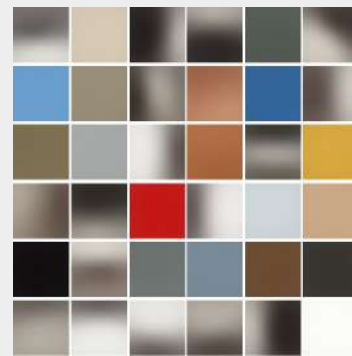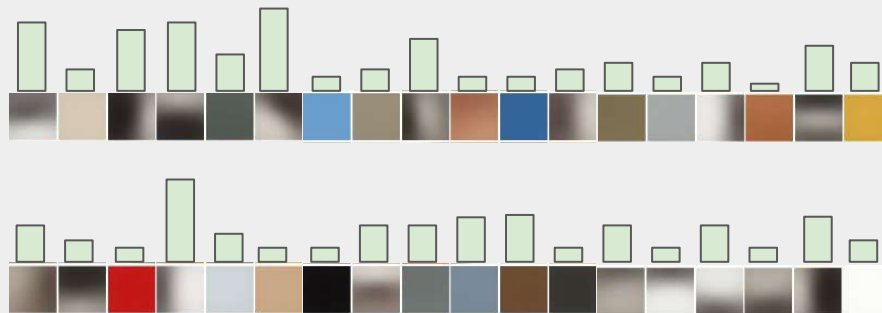**Step 1: Build codebook**

Extract random patches

Cluster patches to form "codebook" of "visual words"

**Step 2: Encode images**

Fei-Fei and Perona, "A bayesian hierarchical model for learning natural scene categories", CVPR 2005

# Image Features

# Image features vs ConvNets



**f**

**10** numbers giving scores for classes

training

Krizhevsky, Sutskever, and Hinton, "Imagenet classification with deep convolutional neural networks", NIPS 2012. Figure copyright Krizhevsky, Sutskever, and Hinton, 2012. Reproduced with permission.

**10** numbers giving scores for classes

training

# One Solution: Feature Transformation



$f(x, y) = (r(x, y), \theta(x, y))$

Transform data with a cleverly chosen **feature transform** f, then apply linear classifier

Color Histogram

Histogram of Oriented Gradients (HoG)

# Today: Neural Networks

# Neural networks: the original linear classifier

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

$$x \in \mathbb{R}^D, W \in \mathbb{R}^{C \times D}$$

# Neural networks: 2 layers

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

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

$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$$

(In practice we will usually add a learnable bias at each layer as well)

# Neural networks: also called fully connected network

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

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

$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$$

"Neural Network" is a very broad term; these are more accurately called "fully-connected networks" or sometimes "multi-layer perceptrons" (MLP)

(In practice we will usually add a learnable bias at each layer as well)

# Neural networks: 3 layers

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

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

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

$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H_1 \times D}, W_2 \in \mathbb{R}^{H_2 \times H_1}, W_3 \in \mathbb{R}^{C \times H_2}$$

(In practice we will usually add a learnable bias at each layer as well)

# Neural networks: hierarchical computation

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

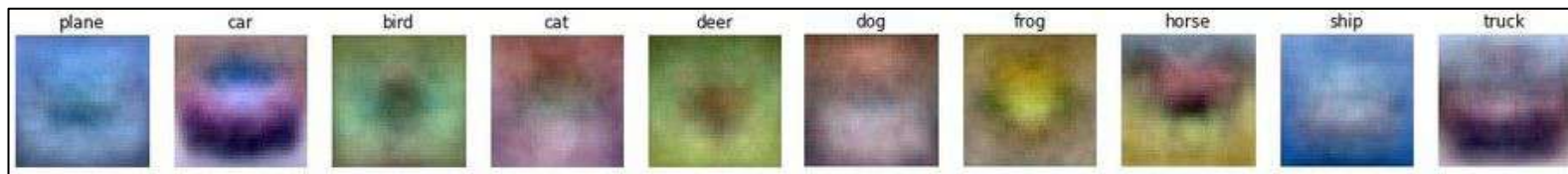(**Now**) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$



$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$$

# Neural networks: learning 100s of templates

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

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



3072     x   W1   h   W2   s   10

100

Learn 100 templates instead of 10.　　　　Share templates between classes

# Neural networks: why is max operator important?

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

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

The function $\max(0, z)$ is called the **activation function.**
**Q:** What if we try to build a neural network without one?

$$f = W_2 W_1 x$$

# Neural networks: why is max operator important?

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

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

The function $\max(0, z)$ is called the **activation function.**
**Q:** What if we try to build a neural network without one?

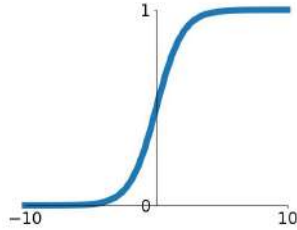$$f = W_2 W_1 x \qquad W_3 = W_2 W_1 \in \mathbb{R}^{C \times H}, f = W_3 x$$

**A**: We end up with a linear classifier again!

# Activation functions

**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**

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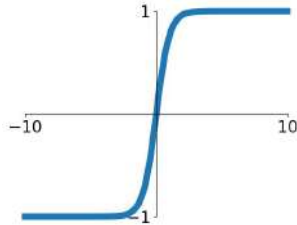
# Activation functions

**Sigmoid**

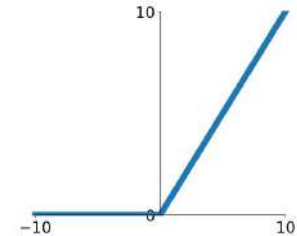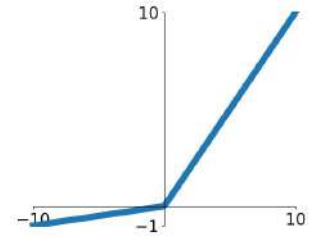$\sigma(x) = \frac{1}{1+e^{-x}}$

**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**

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

# Neural networks: Architectures



"2-layer Neural Net", or
"1-hidden-layer Neural Net"

"3-layer Neural Net", or
"2-hidden-layer Neural Net"

**"Fully-connected" layers**

# Example feed-forward computation of a neural network



```
# forward-pass of a 3-layer neural network:
f = lambda x: 1.0/(1.0 + np.exp(-x)) # activation function (use sigmoid)
x = np.random.randn(3, 1) # random input vector of three numbers (3x1)
h1 = f(np.dot(W1, x) + b1) # calculate first hidden layer activations (4x1)
h2 = f(np.dot(W2, h1) + b2) # calculate second hidden layer activations (4x1)
out = np.dot(W3, h2) + b3 # output neuron (1x1)
```

# Full implementation of training a 2-layer Neural Network needs ~20 lines:

```python
import numpy as np
from numpy.random import randn

N, D_in, H, D_out = 64, 1000, 100, 10
x, y = randn(N, D_in), randn(N, D_out)
w1, w2 = randn(D_in, H), randn(H, D_out)

for t in range(2000):
  h = 1 / (1 + np.exp(-x.dot(w1)))
  y_pred = h.dot(w2)
  loss = np.square(y_pred - y).sum()
  print(t, loss)

  grad_y_pred = 2.0 * (y_pred - y)
  grad_w2 = h.T.dot(grad_y_pred)
  grad_h = grad_y_pred.dot(w2.T)
  grad_w1 = x.T.dot(grad_h * h * (1 - h))

  w1 -= 1e-4 * grad_w1
  w2 -= 1e-4 * grad_w2
```

# Full implementation of training a 2-layer Neural Network needs ~20 lines:

```python
1   import numpy as np
2   from numpy.random import randn
3
4   N, D_in, H, D_out = 64, 1000, 100, 10
5   x, y = randn(N, D_in), randn(N, D_out)
6   w1, w2 = randn(D_in, H), randn(H, D_out)
7
8   for t in range(2000):
9     h = 1 / (1 + np.exp(-x.dot(w1)))
10    y_pred = h.dot(w2)
11    loss = np.square(y_pred - y).sum()
12    print(t, loss)
13
14    grad_y_pred = 2.0 * (y_pred - y)
15    grad_w2 = h.T.dot(grad_y_pred)
16    grad_h = grad_y_pred.dot(w2.T)
17    grad_w1 = x.T.dot(grad_h * h * (1 - h))
18
19    w1 -= 1e-4 * grad_w1
20    w2 -= 1e-4 * grad_w2
```

Define the network

# Full implementation of training a 2-layer Neural Network needs ~20 lines:

```python
import numpy as np
from numpy.random import randn


N, D_in, H, D_out = 64, 1000, 100, 10
x, y = randn(N, D_in), randn(N, D_out)
w1, w2 = randn(D_in, H), randn(H, D_out)

for t in range(2000):
  h = 1 / (1 + np.exp(-x.dot(w1)))
  y_pred = h.dot(w2)
  loss = np.square(y_pred - y).sum()
  print(t, loss)

  grad_y_pred = 2.0 * (y_pred - y)
  grad_w2 = h.T.dot(grad_y_pred)
  grad_h = grad_y_pred.dot(w2.T)
  grad_w1 = x.T.dot(grad_h * h * (1 - h))

  w1 -= 1e-4 * grad_w1
  w2 -= 1e-4 * grad_w2
```

Define the network

Forward pass

# Full implementation of training a 2-layer Neural Network needs ~20 lines:

```
1   import numpy as np
2   from numpy.random import randn
3
4   N, D_in, H, D_out = 64, 1000, 100, 10
5   x, y = randn(N, D_in), randn(N, D_out)
6   w1, w2 = randn(D_in, H), randn(H, D_out)
7
8   for t in range(2000):
9     h = 1 / (1 + np.exp(-x.dot(w1)))
10    y_pred = h.dot(w2)
11    loss = np.square(y_pred - y).sum()
12    print(t, loss)
13
14    grad_y_pred = 2.0 * (y_pred - y)
15    grad_w2 = h.T.dot(grad_y_pred)
16    grad_h = grad_y_pred.dot(w2.T)
17    grad_w1 = x.T.dot(grad_h * h * (1 - h))
18
19    w1 -= 1e-4 * grad_w1
20    w2 -= 1e-4 * grad_w2
```

Define the network

Forward pass

Calculate the analytical gradients

# Full implementation of training a 2-layer Neural Network needs ~20 lines:

```python
import numpy as np
from numpy.random import randn

N, D_in, H, D_out = 64, 1000, 100, 10
x, y = randn(N, D_in), randn(N, D_out)
w1, w2 = randn(D_in, H), randn(H, D_out)

for t in range(2000):
    h = 1 / (1 + np.exp(-x.dot(w1)))
    y_pred = h.dot(w2)
    loss = np.square(y_pred - y).sum()
    print(t, loss)

    grad_y_pred = 2.0 * (y_pred - y)
    grad_w2 = h.T.dot(grad_y_pred)
    grad_h = grad_y_pred.dot(w2.T)
    grad_w1 = x.T.dot(grad_h * h * (1 - h))

    w1 -= 1e-4 * grad_w1
    w2 -= 1e-4 * grad_w2
```

Define the network

Forward pass

Calculate the analytical gradients

Gradient descent

# Setting the number of layers and their sizes



more neurons = more capacity

Do not use size of neural network as a regularizer. Use stronger regularization instead:



$\lambda = 0.001$     $\lambda = 0.01$     $\lambda = 0.1$

(Web demo with ConvNetJS:
http://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html)

$$L(W) = \frac{1}{N} \sum_{i=1}^{N} L_i(f(x_i, W), y_i) + \lambda R(W)$$

This image by Fotis Bobolas is
licensed under CC-BY 2.0

Impulses carried toward cell body

dendrite

presynaptic
terminal

axon

cell body

Impulses carried away
from cell body

Impulses carried toward cell body

dendrite

presynaptic terminal

axon

cell body

Impulses carried away from cell body

This image by Felipe Perucho is licensed under CC-BY 3.0

$x_0$          $w_0$

axon from a neuron          synapse

$w_0 x_0$

dendrite

cell body

$w_1 x_1$

$$\sum_i w_i x_i + b$$          $f$

$$f\left(\sum_i w_i x_i + b\right)$$

output axon

activation function

$w_2 x_2$

Impulses carried toward cell body

dendrite

presynaptic terminal

axon

cell body

Impulses carried away from cell body

This image by Felipe Perucho is licensed under CC-BY 3.0

sigmoid activation function

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

$x_0$

axon from a neuron

$w_0$ synapse

$w_0 x_0$

dendrite

$w_1 x_1$

cell body

$\sum_i w_i x_i + b$ $f$

$f\left(\sum_i w_i x_i + b\right)$

output axon

activation function

$w_2 x_2$

Impulses carried toward cell body

dendrite

presynaptic
terminal

axon

cell body

Impulses carried away
from cell body

```
class Neuron:
    # ...
    def neuron_tick(inputs):
        """ assume inputs and weights are 1-D numpy arrays and bias is a number """
        cell_body_sum = np.sum(inputs * self.weights) + self.bias
        firing_rate = 1.0 / (1.0 + math.exp(-cell_body_sum)) # sigmoid activation func
        return firing_rate
```

$x_0$

$w_0$

axon from a neuron

synapse

$w_0 x_0$

dendrite

$w_1 x_1$

cell body

$\sum_i w_i x_i + b$  $f$

$f\left(\sum_i w_i x_i + b\right)$

output axon

activation
function

$w_2 x_2$

Biological Neurons:
Complex connectivity patterns

Neurons in a neural network:
Organized into regular layers for computational efficiency

Biological Neurons:
Complex connectivity patterns

But neural networks with random connections can work too!

Xie et al, "Exploring Randomly Wired Neural Networks for Image Recognition", arXiv 2019

# Be very careful with your brain analogies!

**Biological Neurons:**
- Many different types
- Dendrites can perform complex non-linear computations
- Synapses are not a single weight but a complex non-linear dynamical system

[Dendritic Computation. London and Hausser]

# Plugging in neural networks with loss functions

$$s = f(x; W_1, W_2) = W_2 \max(0, W_1 x)$$    Nonlinear score function

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

$$R(W) = \sum_k W_k^2$$    Regularization

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + \lambda R(W_1) + \lambda R(W_2)$$    Total loss: data loss + regularization

# Problem: How to compute gradients?

$$s = f(x; W_1, W_2) = W_2 \max(0, W_1 x)$$   Nonlinear score function

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

$$R(W) = \sum_k W_k^2$$   Regularization

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + \lambda R(W_1) + \lambda R(W_2)$$   Total loss: data loss + regularization

If we can compute $\dfrac{\partial L}{\partial W_1}, \dfrac{\partial L}{\partial W_2}$ then we can learn $W_1$ and $W_2$

# (Bad) Idea: Derive $\nabla_W L$ on paper

$$s = f(x; W) = Wx$$

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

$$= \sum_{j \neq y_i} \max(0, W_{j,:} \cdot x + W_{y_i,:} \cdot x + 1)$$

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + \lambda \sum_k W_k^2$$

$$= \frac{1}{N} \sum_{i=1}^{N} \sum_{j \neq y_i} \max(0, W_{j,:} \cdot x + W_{y_i,:} \cdot x + 1) + \lambda \sum_k W_k^2$$

$$\nabla_W L = \nabla_W \left( \frac{1}{N} \sum_{i=1}^{N} \sum_{j \neq y_i} \max(0, W_{j,:} \cdot x + W_{y_i,:} \cdot x + 1) + \lambda \sum_k W_k^2 \right)$$

**Problem**: Very tedious: Lots of matrix calculus, need lots of paper

**Problem**: What if we want to change loss? E.g. use softmax instead of SVM? Need to re-derive from scratch =(

**Problem**: Not feasible for very complex models!

# Better Idea: Computational graphs + Backpropagation

$$f = Wx$$

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

# Convolutional network (AlexNet)

input image

weights

loss



Figure copyright Alex Krizhevsky, Ilya Sutskever, and
Geoffrey Hinton, 2012. Reproduced with permission.

# Really complex neural networks!!

input image

loss



Figure reproduced with permission from a Twitter post by Andrej Karpathy.

# Neural Turing Machine



Figure reproduced with permission from a [Twitter post](#) by Andrej Karpathy.

# Solution: Backpropagation

# 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

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

Backpropagation: a simple example

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

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

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

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

Backpropagation: a simple example

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

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



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

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

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

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

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

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



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

Backpropagation: a simple example

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

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

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

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

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



$$\frac{\partial f}{\partial f}$$
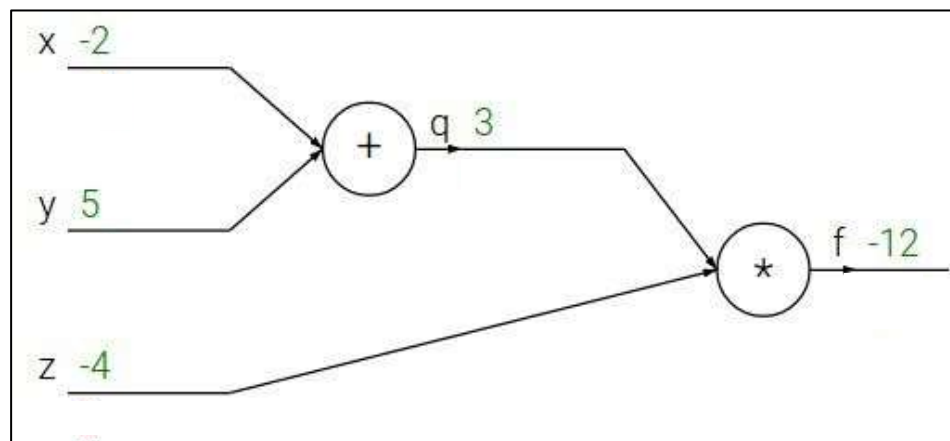
Backpropagation: a simple example

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

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

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

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

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



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

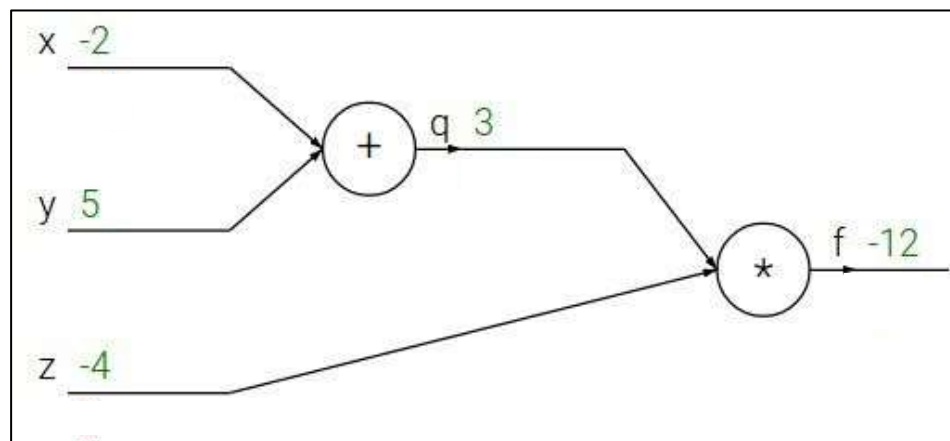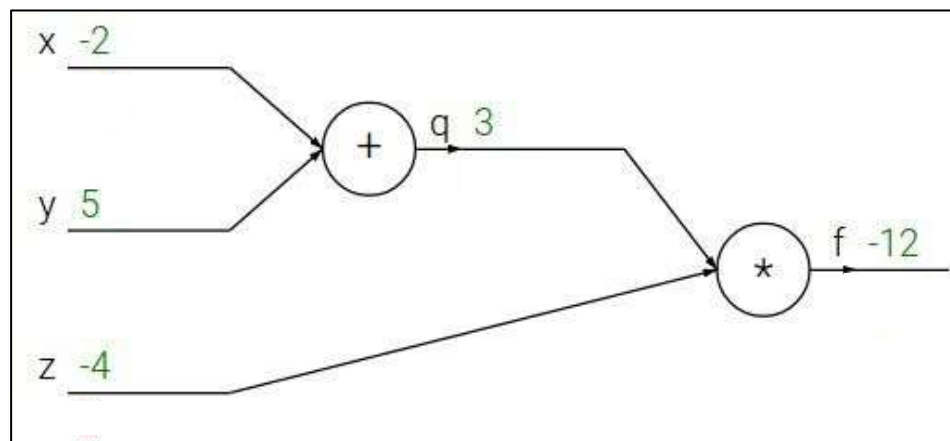Backpropagation: a simple example

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

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



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

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

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

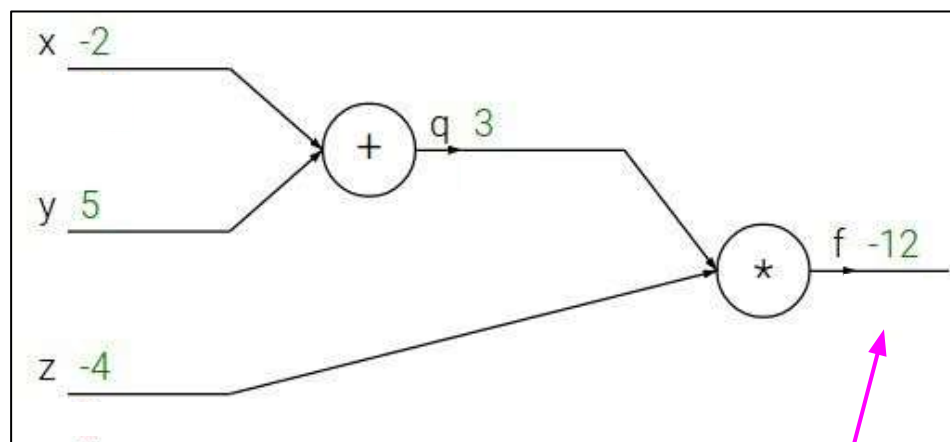$$\frac{\partial f}{\partial z}$$

Backpropagation: a simple example

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

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

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

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

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

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

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

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



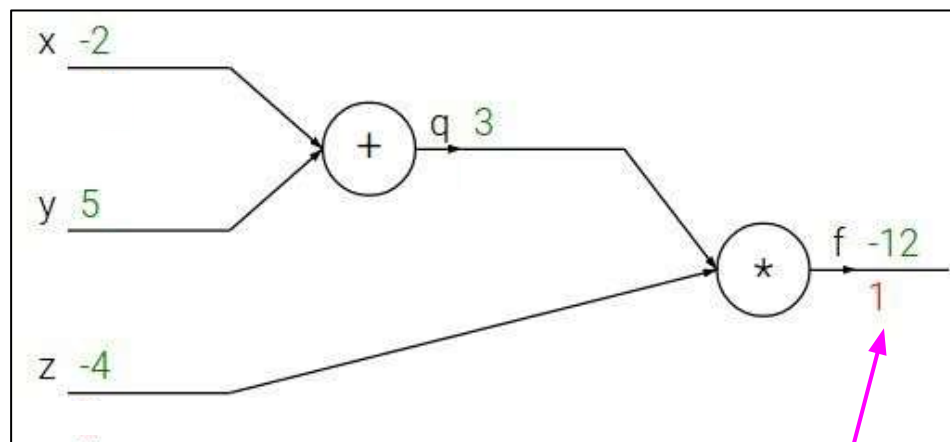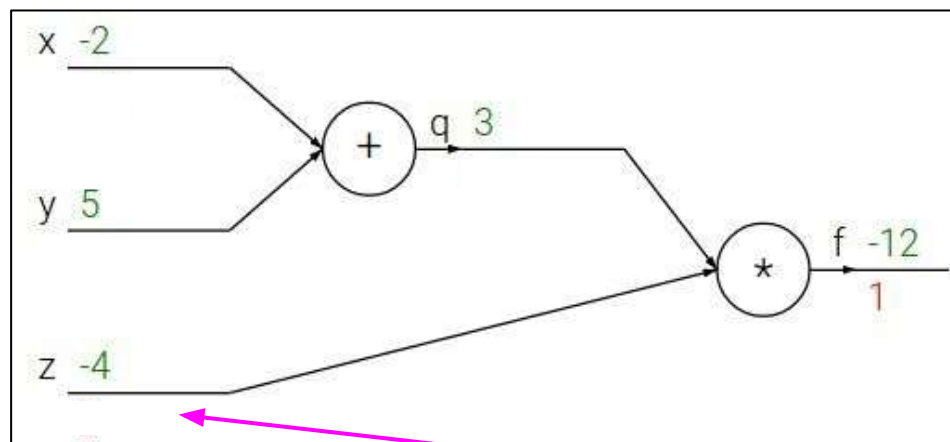$\frac{\partial f}{\partial q}$

Backpropagation: a simple example

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

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

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

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

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



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

Chain rule:

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

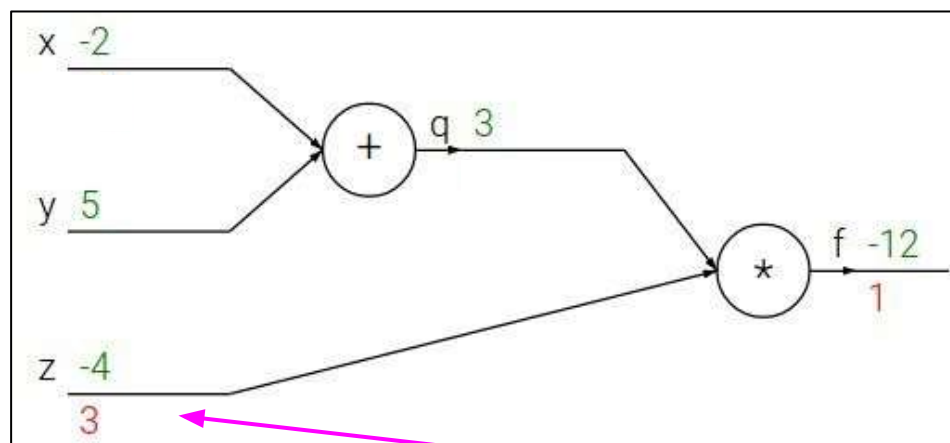Upstream gradient    Local gradient

Backpropagation: a simple example

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

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

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

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

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       Local
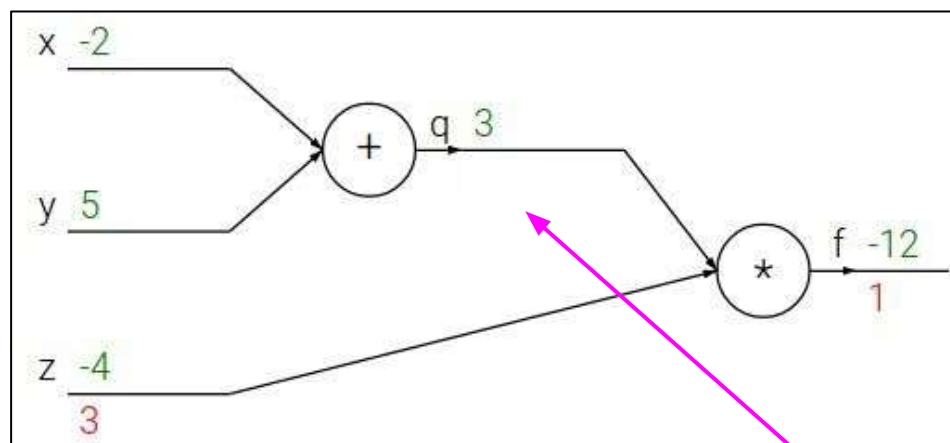gradient       gradient

Backpropagation: a simple example

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

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

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

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

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



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

Chain rule:

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

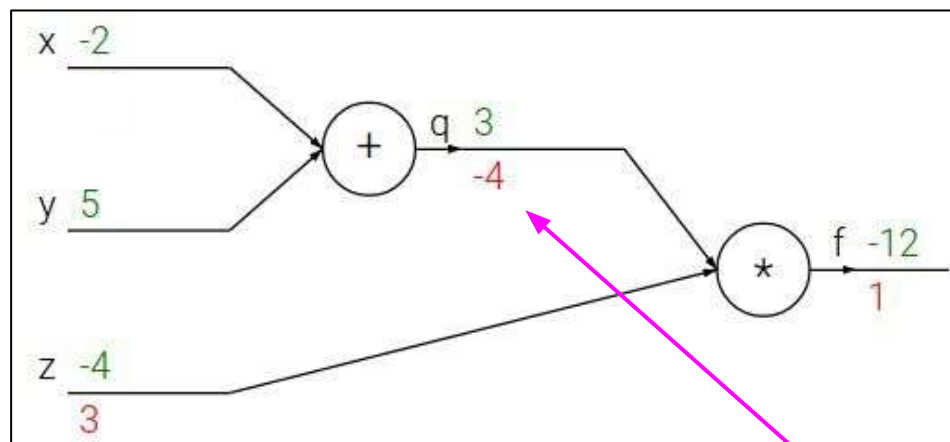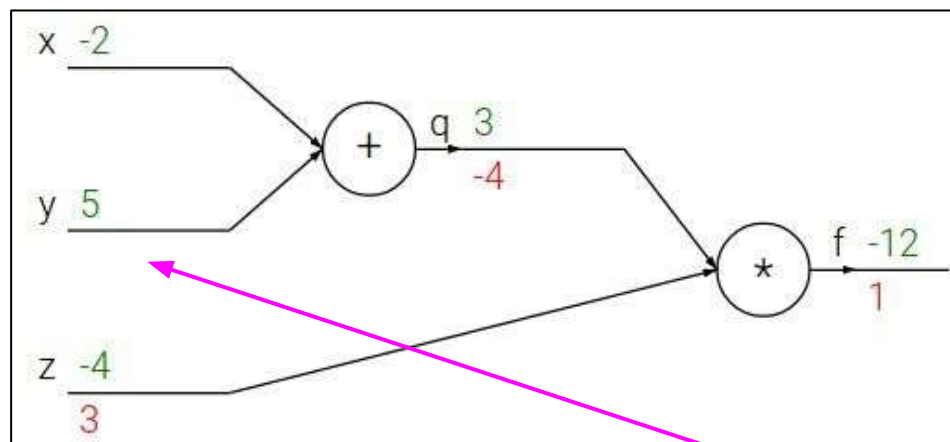Upstream gradient    Local gradient

Backpropagation: a simple example

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

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

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

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

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



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

Chain rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$
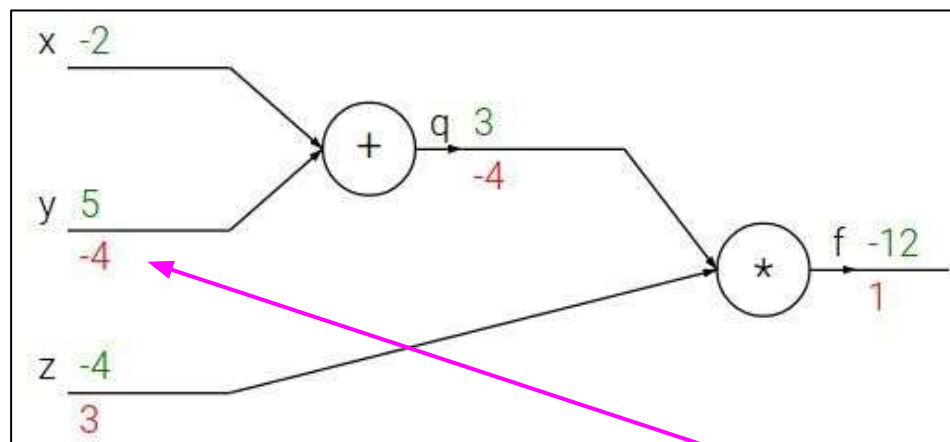
Upstream gradient    Local gradient

"local gradient"

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

$$\frac{\partial z}{\partial y}$$

$x$

$y$

$z$

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

f

"Upstream gradient"

$$x$$

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

"Downstream gradients"

$$y$$

"local gradient"

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

$$f$$

$$\frac{\partial z}{\partial y}$$

$$z$$

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

"Upstream gradient"

$x$

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

$\dfrac{\partial z}{\partial x}$

"local gradient"

$z$

"Downstream gradients"

$y$

$\dfrac{\partial L}{\partial y} = \dfrac{\partial L}{\partial z}\dfrac{\partial z}{\partial y}$

$\dfrac{\partial z}{\partial y}$

f

$\dfrac{\partial L}{\partial z}$

"Upstream gradient"

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

Another example:
$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



w0 2.00

x0 -1.00

w1 -3.00

x1 -2.00

w2 -3.00

Another example: $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$

Another example: $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



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

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

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

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

Another example: $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



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

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

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

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

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



Upstream gradient    Local gradient

$$(1.00)(\frac{-1}{1.37^2}) = -0.53$$
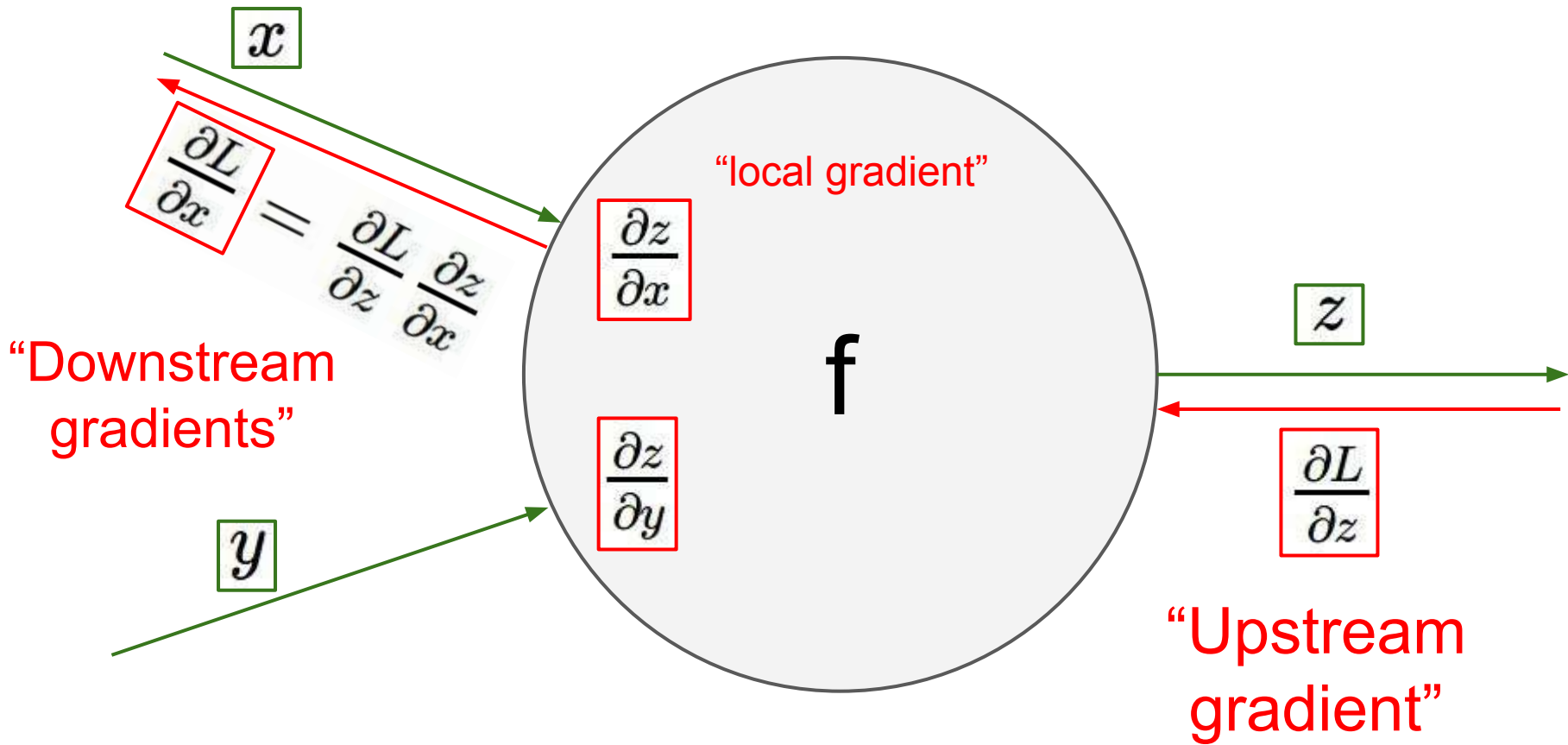
$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \qquad \Big| \qquad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \Big| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$
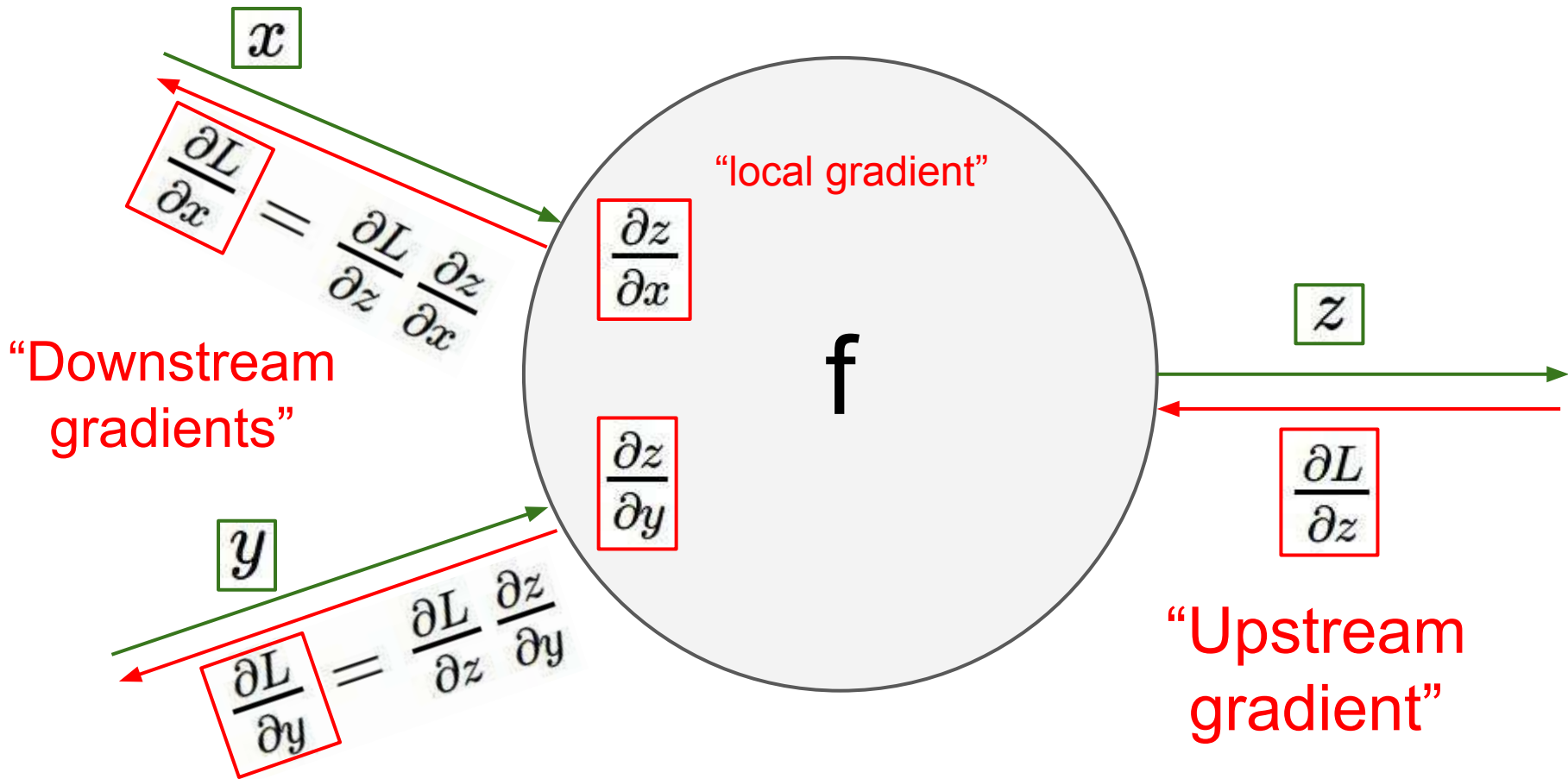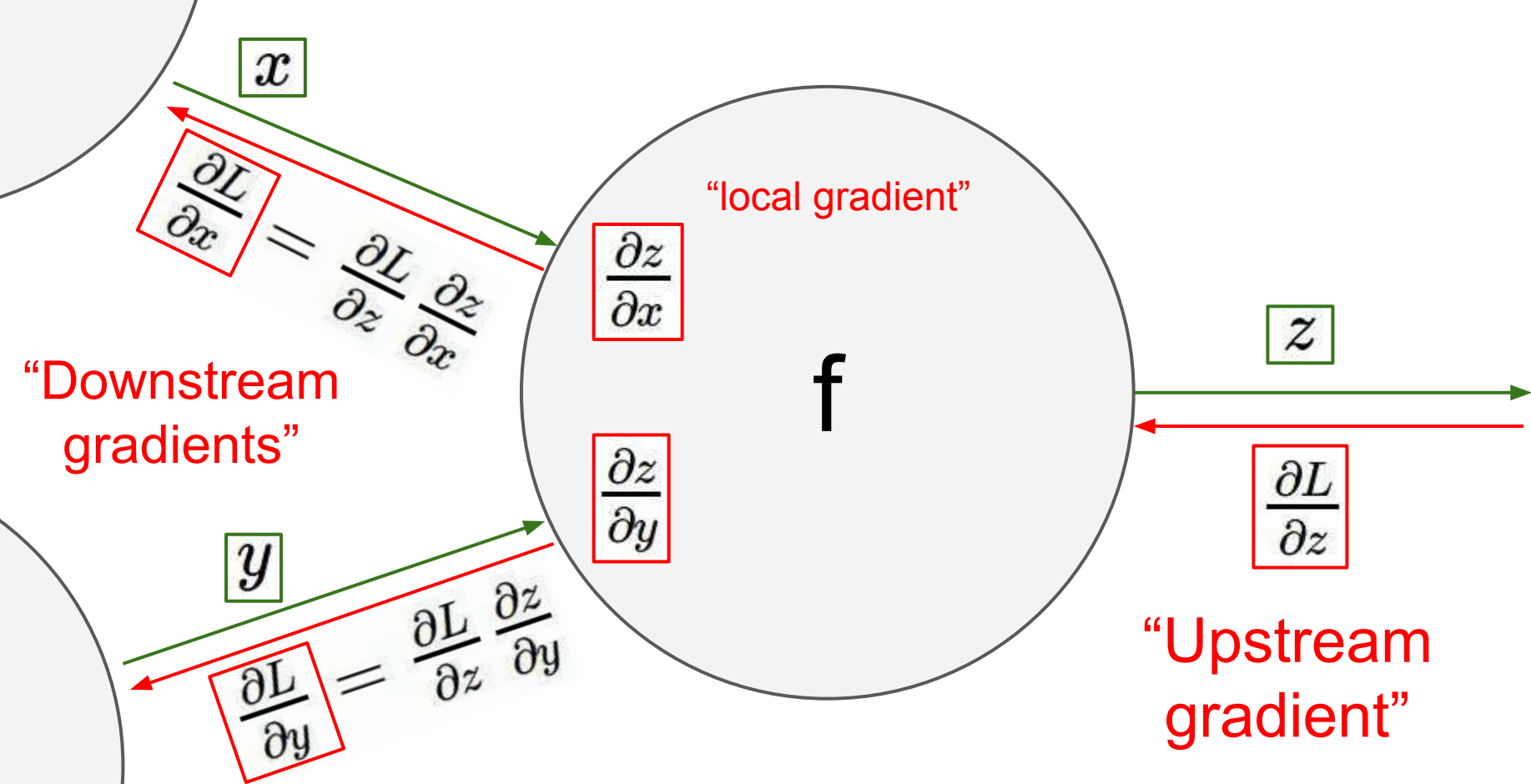
Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



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

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

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

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

Another example: $f(w,x) = \dfrac{1}{1+e^{-(w_0x_0+w_1x_1+w_2)}}$



Upstream gradient    Local gradient

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

| | |
|---|---|
| $f(x) = e^x$ | $\rightarrow$ | $\dfrac{df}{dx} = e^x$ |
| $f_a(x) = ax$ | $\rightarrow$ | $\dfrac{df}{dx} = a$ |

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

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

Another example:
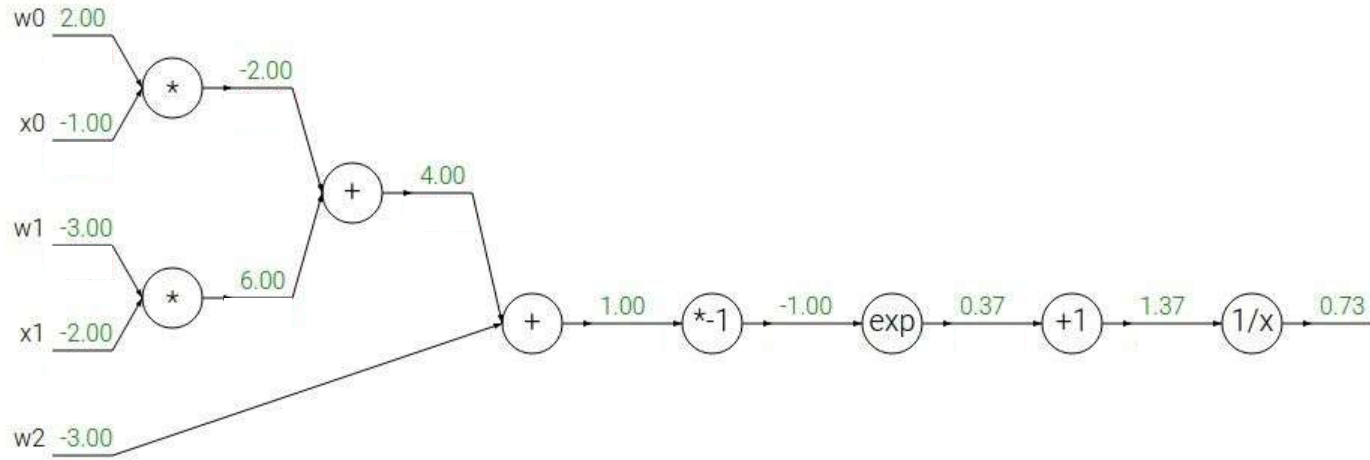
$$f(w,x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



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

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

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

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

Another example:

$$f(w,x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



Upstream gradient    Local gradient

$$(-0.53)(e^{-1}) = -0.20$$

| | | |
|---|---|---|
| $f(x) = e^x$ | $\rightarrow$ | $\frac{df}{dx} = e^x$ |
| $f_a(x) = ax$ | $\rightarrow$ | $\frac{df}{dx} = a$ |

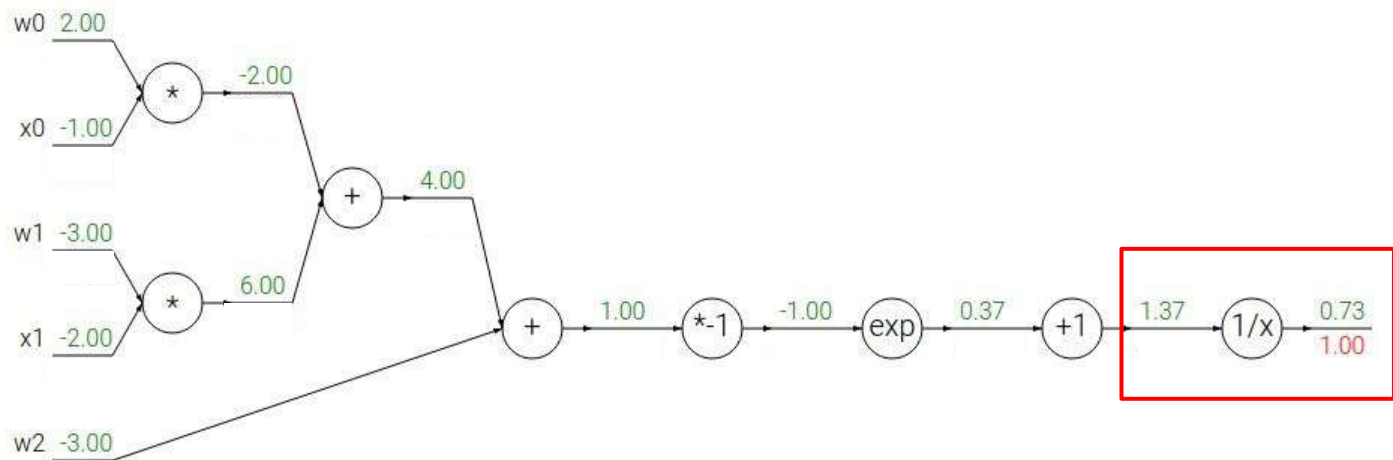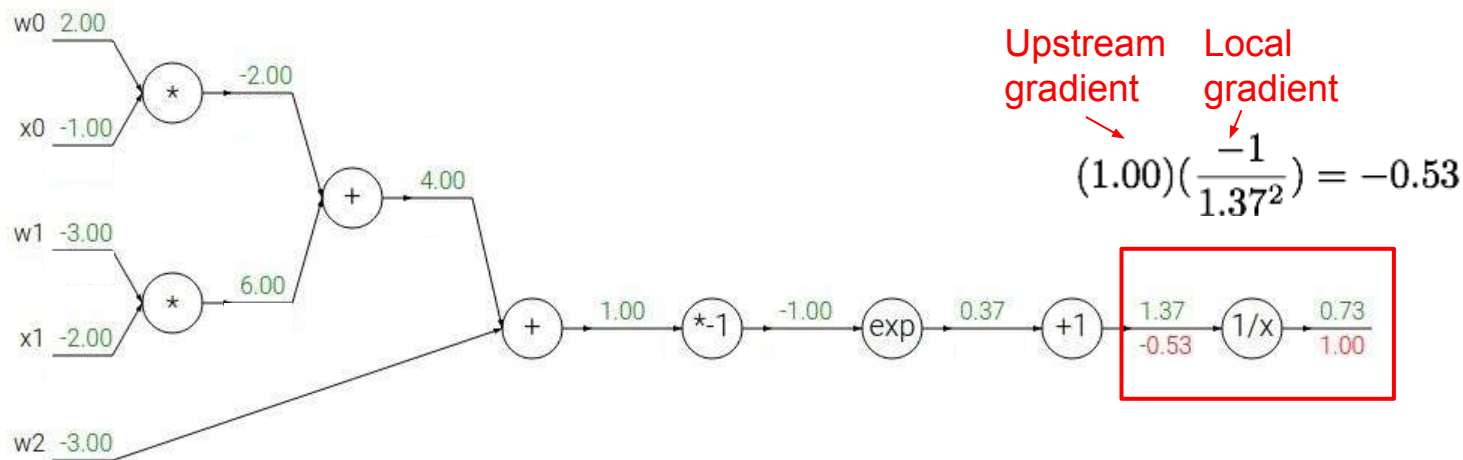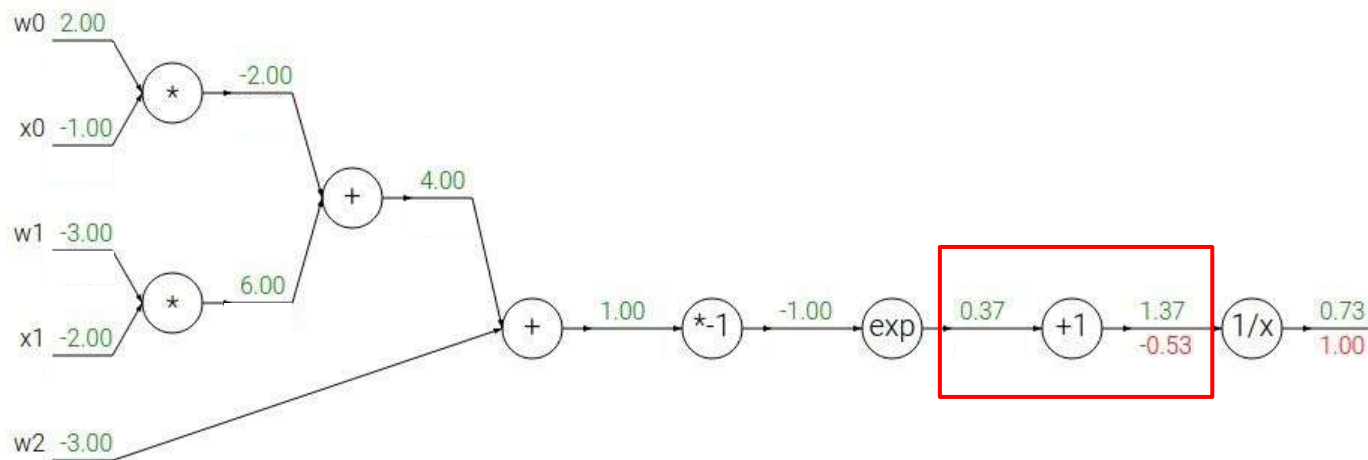| | | |
|---|---|---|
| $f(x) = \frac{1}{x}$ | $\rightarrow$ | $\frac{df}{dx} = -1/x^2$ |
| $f_c(x) = c + x$ | $\rightarrow$ | $\frac{df}{dx} = 1$ |

Another example:

$$f(w,x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



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

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

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

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

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



Upstream gradient    Local gradient

$$(-0.20)(-1) = 0.20$$

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

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

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

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

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



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

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \Big| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

w0 2.00
x0 -1.00
-2.00

w1 -3.00
x1 -2.00
6.00

4.00
0.20

1.00
0.20

*-1  -1.00
-0.20

exp  0.37
-0.53

+1  1.37
-0.53

1/x  0.73
1.00

w2 -3.00
0.20

[upstream gradient] x [local gradient]
[0.2] x [1] = 0.2
[0.2] x [1] = 0.2  (both inputs!)

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

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

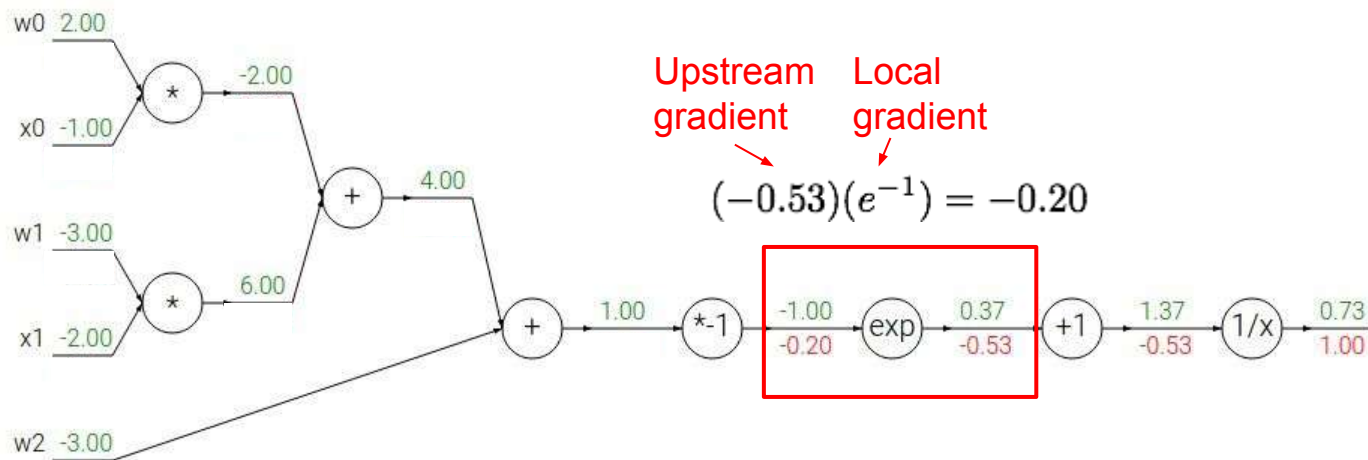$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$
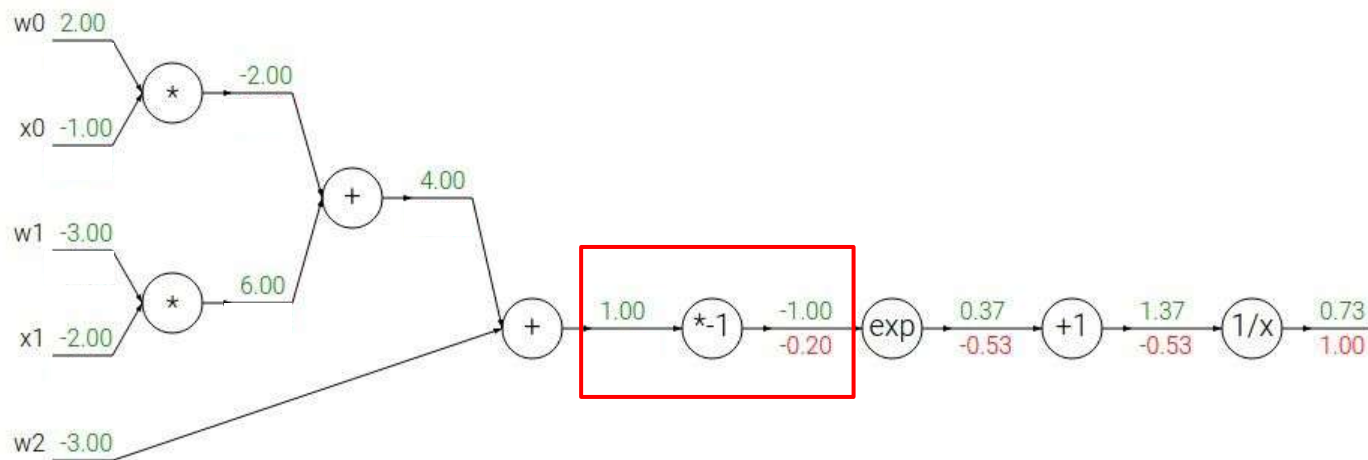
$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



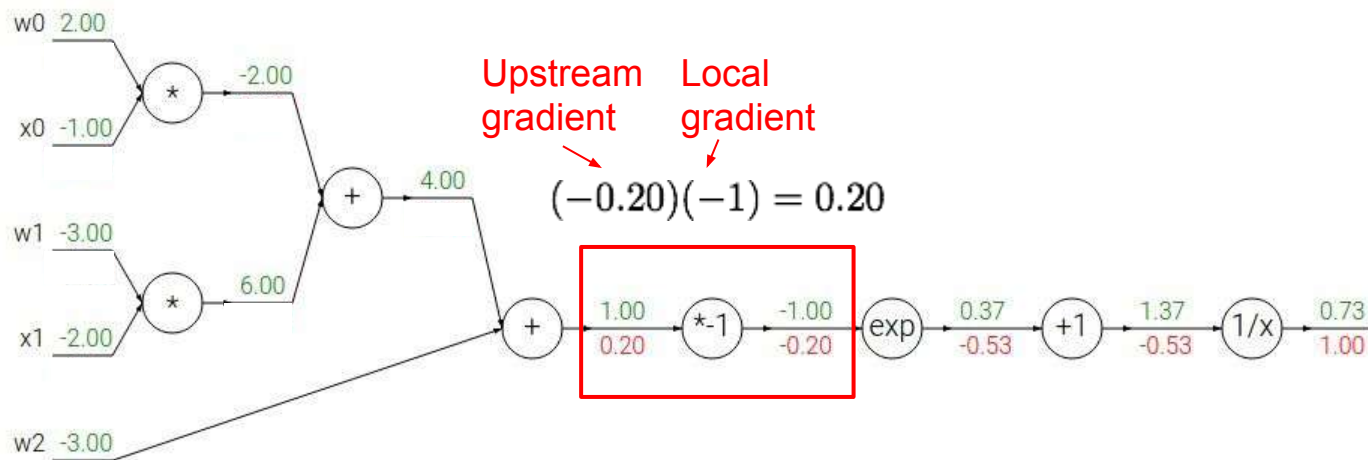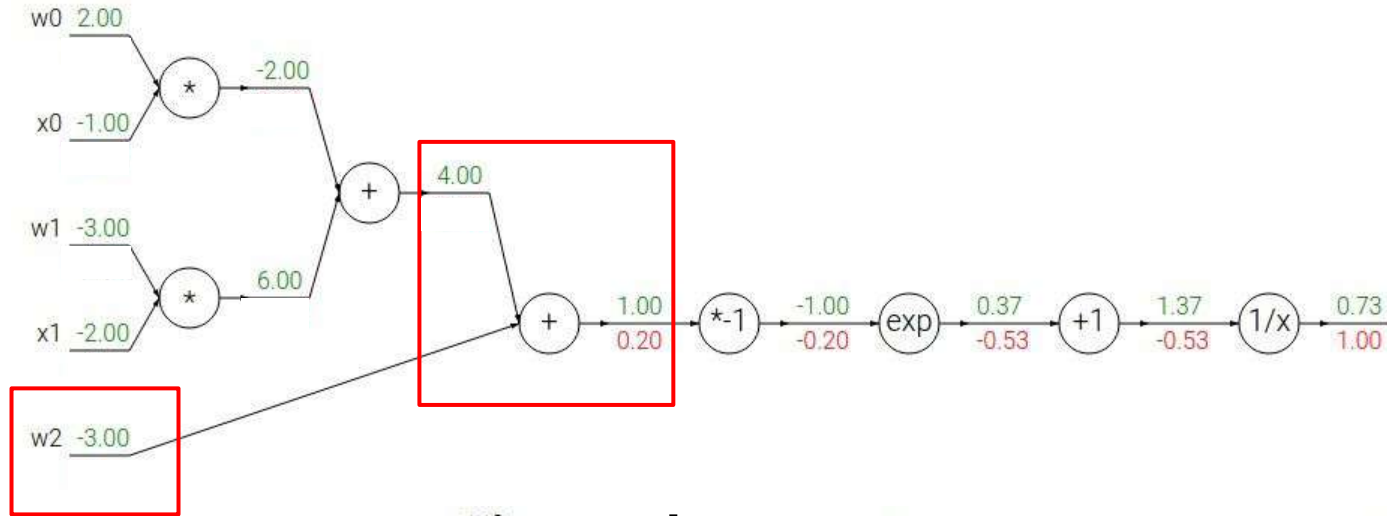$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \qquad\Big|\qquad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad\Big|\qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



[upstream gradient] x [local gradient]
w0: [0.2] x [-1] = -0.2
x0: [0.2] x [2] = 0.4

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

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

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

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

# Another example:

$$f(w,x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

Computational graph representation may not be unique. Choose one where local gradients at each node can be easily expressed!

Sigmoid function

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



Sigmoid

w0  2.00
    -0.20

x0  -1.00
    0.40

w1  -3.00

x1  -2.00

w2  -3.00
    0.20

*   -2.00
    0.20

*   6.00
    0.20

+   4.00
    0.20

+   1.00
    0.20

*-1   -1.00
      -0.20

exp   0.37
      -0.53

+1    1.37
      -0.53

1/x   0.73
      1.00

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

Sigmoid function

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



Sigmoid

Sigmoid local gradient:

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

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

Sigmoid function

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



Sigmoid

[upstream gradient] x [local gradient]
[1.00] x [(1 - 1/(1+e$^1$)) (1/(1+e$^1$))] = 0.2

Sigmoid local gradient:

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

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

Sigmoid function

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



Sigmoid

[upstream gradient] x [local gradient]
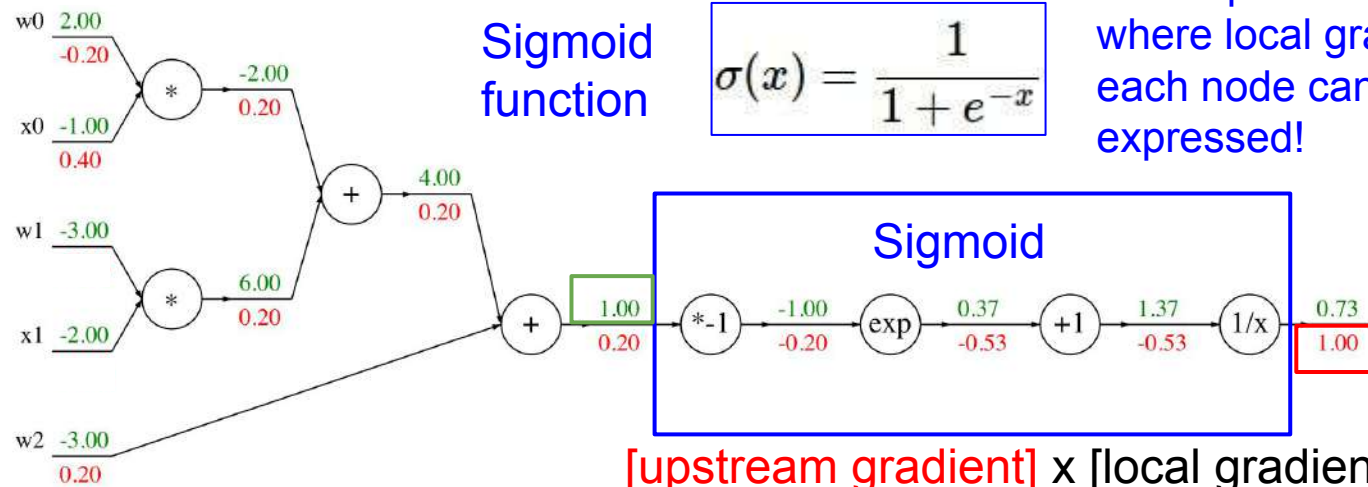[1.00] x [(1 - 0.73) (0.73)] = 0.2

Sigmoid local gradient:

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

# Patterns in gradient flow

**add** gate: gradient distributor

# Patterns in gradient flow

**add** gate: gradient distributor



**mul** gate: "swap multiplier"

# Patterns in gradient flow

**add** gate: gradient distributor



**mul** gate: "swap multiplier"



**copy** gate: gradient adder

# Patterns in gradient flow

**add** gate: gradient distributor



**mul** gate: "swap multiplier"



**copy** gate: gradient adder



**max** gate: gradient router

# Backprop Implementation: "Flat" code



```
def f(w0, x0, w1, x1, w2):
    s0 = w0 * x0
    s1 = w1 * x1
    s2 = s0 + s1
    s3 = s2 + w2
    L = sigmoid(s3)
```

Forward pass:
Compute output

Backward pass:
Compute grads

```
grad_L = 1.0
grad_s3 = grad_L * (1 - L) * L
grad_w2 = grad_s3
grad_s2 = grad_s3
grad_s0 = grad_s2
grad_s1 = grad_s2
grad_w1 = grad_s1 * x1
grad_x1 = grad_s1 * w1
grad_w0 = grad_s0 * x0
grad_x0 = grad_s0 * w0
```

# Backprop Implementation: "Flat" code



Forward pass:
Compute output

```
def f(w0, x0, w1, x1, w2):
    s0 = w0 * x0
    s1 = w1 * x1
    s2 = s0 + s1
    s3 = s2 + w2
    L = sigmoid(s3)
```

Base case

```
grad_L = 1.0
grad_s3 = grad_L * (1 - L) * L
grad_w2 = grad_s3
grad_s2 = grad_s3
grad_s0 = grad_s2
grad_s1 = grad_s2
grad_w1 = grad_s1 * x1
grad_x1 = grad_s1 * w1
grad_w0 = grad_s0 * x0
grad_x0 = grad_s0 * w0
```

# Backprop Implementation: "Flat" code



Forward pass:
Compute output

```
def f(w0, x0, w1, x1, w2):
    s0 = w0 * x0
    s1 = w1 * x1
    s2 = s0 + s1
    s3 = s2 + w2
    L = sigmoid(s3)
```

Sigmoid

```
grad_L = 1.0
grad_s3 = grad_L * (1 - L) * L
grad_w2 = grad_s3
grad_s2 = grad_s3
grad_s0 = grad_s2
grad_s1 = grad_s2
grad_w1 = grad_s1 * x1
grad_x1 = grad_s1 * w1
grad_w0 = grad_s0 * x0
grad_x0 = grad_s0 * w0
```

# Backprop Implementation: "Flat" code



Forward pass: Compute output

```
def f(w0, x0, w1, x1, w2):
    s0 = w0 * x0
    s1 = w1 * x1
    s2 = s0 + s1
    s3 = s2 + w2
    L = sigmoid(s3)
```

Add gate

```
grad_L = 1.0
grad_s3 = grad_L * (1 - L) * L
grad_w2 = grad_s3
grad_s2 = grad_s3
grad_s0 = grad_s2
grad_s1 = grad_s2
grad_w1 = grad_s1 * x1
grad_x1 = grad_s1 * w1
grad_w0 = grad_s0 * x0
grad_x0 = grad_s0 * w0
```

# Backprop Implementation: "Flat" code



Forward pass: Compute output

```
def f(w0, x0, w1, x1, w2):
    s0 = w0 * x0
    s1 = w1 * x1
    s2 = s0 + s1
    s3 = s2 + w2
    L = sigmoid(s3)
```

```
grad_L = 1.0
grad_s3 = grad_L * (1 - L) * L
grad_w2 = grad_s3
grad_s2 = grad_s3
```

Add gate

```
grad_s0 = grad_s2
grad_s1 = grad_s2
grad_w1 = grad_s1 * x1
grad_x1 = grad_s1 * w1
grad_w0 = grad_s0 * x0
grad_x0 = grad_s0 * w0
```

# Backprop Implementation: "Flat" code

```python
def f(w0, x0, w1, x1, w2):
    s0 = w0 * x0
    s1 = w1 * x1
    s2 = s0 + s1
    s3 = s2 + w2
    L = sigmoid(s3)

    grad_L = 1.0
    grad_s3 = grad_L * (1 - L) * L
    grad_w2 = grad_s3
    grad_s2 = grad_s3
    grad_s0 = grad_s2
    grad_s1 = grad_s2
    grad_w1 = grad_s1 * x1
    grad_x1 = grad_s1 * w1
    grad_w0 = grad_s0 * x0
    grad_x0 = grad_s0 * w0
```

Multiply gate

# Backprop Implementation: "Flat" code



Forward pass: Compute output

```
def f(w0, x0, w1, x1, w2):
    s0 = w0 * x0
    s1 = w1 * x1
    s2 = s0 + s1
    s3 = s2 + w2
    L = sigmoid(s3)
```
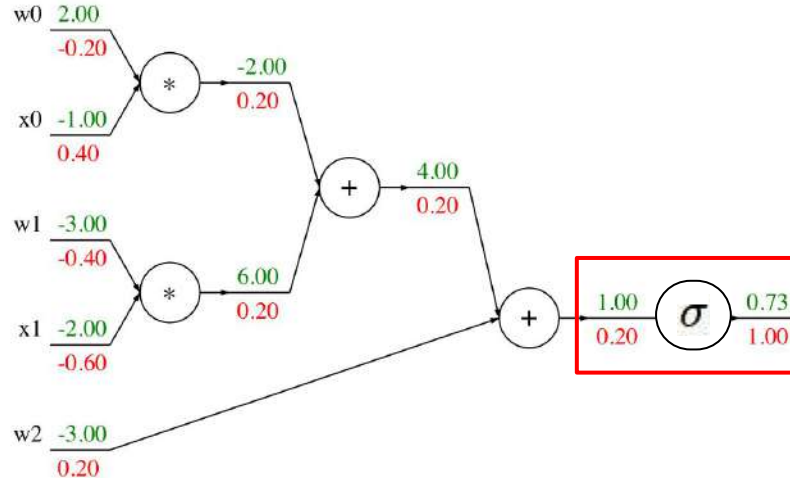
```
grad_L = 1.0
grad_s3 = grad_L * (1 - L) * L
grad_w2 = grad_s3
grad_s2 = grad_s3
grad_s0 = grad_s2
grad_s1 = grad_s2
grad_w1 = grad_s1 * x1
grad_x1 = grad_s1 * w1
grad_w0 = grad_s0 * x0
grad_x0 = grad_s0 * w0
```

Multiply gate

# "Flat" Backprop: Do this for assignment 1!

## Stage your forward/backward computation!

E.g. for the SVM:

margins

```
# receive W (weights), X (data)
# forward pass (we have 6 lines)
scores = #...
margins = #...
data_loss = #...
reg_loss = #...
loss = data_loss + reg_loss
# backward pass (we have 5 lines)
dmargins = # ... (optionally, we go direct to dscores)
dscores = #...
dW = #...
```

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

x

W

\*

s (scores)

hinge loss

+

L

R

$R(W)$

# "Flat" Backprop: Do this for assignment 1!

E.g. for two-layer neural net:

```
# receive W1,W2,b1,b2 (weights/biases), X (data)
# forward pass:
h1 = #... function of X,W1,b1
scores = #... function of h1,W2,b2
loss = #... (several lines of code to evaluate Softmax loss)
# backward pass:
dscores = #...
dh1,dW2,db2 = #...
dW1,db1 = #...
```

# Backprop Implementation: Modularized API

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



```python
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

Gate / Node / Function object: Actual PyTorch code

x

z

\*

y

(x,y,z are scalars)

```python
class Multiply(torch.autograd.Function):
    @staticmethod
    def forward(ctx, x, y):
        ctx.save_for_backward(x, y)
        z = x * y
        return z
    @staticmethod
    def backward(ctx, grad_z):
        x, y = ctx.saved_tensors
        grad_x = y * grad_z    # dz/dx * dL/dz
        grad_y = x * grad_z    # dz/dy * dL/dz
        return grad_x, grad_y
```

Need to stash some values for use in backward

Upstream gradient

Multiply upstream and local gradients

# Example: PyTorch operators

# PyTorch sigmoid layer

```
1   #ifndef TH_GENERIC_FILE
2   #define TH_GENERIC_FILE "THNN/generic/Sigmoid.c"
3   #else
4
5   void THNN_(Sigmoid_updateOutput)(
6            THNNState *state,
7            THTensor *input,
8            THTensor *output)
9   {
10    THTensor_(sigmoid)(output, input);
11  }
12
13  void THNN_(Sigmoid_updateGradInput)(
14            THNNState *state,
15            THTensor *gradOutput,
16            THTensor *gradInput,
17            THTensor *output)
18  {
19    THNN_CHECK_NELEMENT(output, gradOutput);
20    THTensor_(resizeAs)(gradInput, output);
21    TH_TENSOR_APPLY3(scalar_t, gradInput, scalar_t, gradOutput, scalar_t, output,
22      scalar_t z = *output_data;
23      *gradInput_data = *gradOutput_data * (1. - z) * z;
24    );
25  }
26
27  #endif
```

**Forward**

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

[Source](#)

```
1   #ifndef TH_GENERIC_FILE
2   #define TH_GENERIC_FILE "THNN/generic/Sigmoid.c"
3   #else
4
5   void THNN_(Sigmoid_updateOutput)(
6           THNNState *state,
7           THTensor *input,
8           THTensor *output)
9   {
10    THTensor_(sigmoid)(output, input);
11  }
12
13  void THNN_(Sigmoid_updateGradInput)(
14          THNNState *state,
15          THTensor *gradOutput,
16          THTensor *gradInput,
17          THTensor *output)
18  {
19    THNN_CHECK_NELEMENT(output, gradOutput);
20    THTensor_(resizeAs)(gradInput, output);
21    TH_TENSOR_APPLY3(scalar_t, gradInput, scalar_t, gradOutput, scalar_t, output,
22      scalar_t z = *output_data;
23      *gradInput_data = *gradOutput_data * (1. - z) * z;
24    );
25  }
26
27  #endif
```

Forward

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

```
static void sigmoid_kernel(TensorIterator& iter) {
  AT_DISPATCH_FLOATING_TYPES(iter.dtype(), "sigmoid_cpu", [&]() {
    unary_kernel_vec(
      iter,
      [=](scalar_t a) -> scalar_t { return (1 / (1 + std::exp((-a)))); },
      [=](Vec256<scalar_t> a) {
        a = Vec256<scalar_t>((scalar_t)(0)) - a;
        a = a.exp();
        a = Vec256<scalar_t>((scalar_t)(1)) + a;
        a = a.reciprocal();
        return a;
      });
  });
}
```

Forward actually defined elsewhere...

```
return (1 / (1 + std::exp((-a))));
```

Source

# PyTorch sigmoid layer

```
1   #ifndef TH_GENERIC_FILE
2   #define TH_GENERIC_FILE "THNN/generic/Sigmoid.c"
3   #else
4
5   void THNN_(Sigmoid_updateOutput)(
6           THNNState *state,
7           THTensor *input,
8           THTensor *output)
9   {
10    THTensor_(sigmoid)(output, input);
11  }
12
13  void THNN_(Sigmoid_updateGradInput)(
14          THNNState *state,
15          THTensor *gradOutput,
16          THTensor *gradInput,
17          THTensor *output)
18  {
19    THNN_CHECK_NELEMENT(output, gradOutput);
20    THTensor_(resizeAs)(gradInput, output);
21    TH_TENSOR_APPLY3(scalar_t, gradInput, scalar_t, gradOutput, scalar_t, output,
22      scalar_t z = *output_data;
23      *gradInput_data = *gradOutput_data * (1. - z) * z;
24    );
25  }
26
27  #endif
```

**Forward**

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

**Backward**

$$(1 - \sigma(x))\,\sigma(x)$$

```
static void sigmoid_kernel(TensorIterator& iter) {
  AT_DISPATCH_FLOATING_TYPES(iter.dtype(), "sigmoid_cpu", [&]() {
    unary_kernel_vec(
      iter,
      [=](scalar_t a) -> scalar_t { return (1 / (1 + std::exp((-a)))); },
      [=](Vec256<scalar_t> a) {
        a = Vec256<scalar_t>((scalar_t)(0)) - a;
        a = a.exp();
        a = Vec256<scalar_t>((scalar_t)(1)) + a;
        a = a.reciprocal();
        return a;
      });
  });
}
```

Forward actually
defined elsewhere...

Source

# Summary for today:

- **(Fully-connected) Neural Networks** are stacks of linear functions and nonlinear activation functions; they have much more representational power than linear classifiers
- **backpropagation** = recursive application of the chain rule along a computational graph to compute the gradients of all inputs/parameters/intermediates
- implementations maintain a graph structure, where the nodes implement the **forward**() / **backward**() API
- **forward**: compute result of an operation and save any intermediates needed for gradient computation in memory
- **backward**: apply the chain rule to compute the gradient of the loss function with respect to the inputs

So far: backprop with scalars

Next time: vector-valued functions!

# Next Time: Convolutional Networks!

# Recap: Vector derivatives

Scalar to Scalar

$x \in \mathbb{R}, y \in \mathbb{R}$

Regular derivative:

$$\frac{\partial y}{\partial x} \in \mathbb{R}$$

If x changes by a
small amount, how
much will y change?

# Recap: Vector derivatives

**Scalar to Scalar**

$x \in \mathbb{R}, y \in \mathbb{R}$

Regular derivative:

$$\frac{\partial y}{\partial x} \in \mathbb{R}$$

If x changes by a small amount, how much will y change?

**Vector to Scalar**

$x \in \mathbb{R}^N, y \in \mathbb{R}$

Derivative is **Gradient**:

$$\frac{\partial y}{\partial x} \in \mathbb{R}^N \quad \left( \frac{\partial y}{\partial x} \right)_n = \frac{\partial y}{\partial x_n}$$

For each element of x, if it changes by a small amount then how much will y change?

# Recap: Vector derivatives

| Scalar to Scalar | Vector to Scalar | Vector to Vector |
|---|---|---|
| $x \in \mathbb{R}, y \in \mathbb{R}$ | $x \in \mathbb{R}^N, y \in \mathbb{R}$ | $x \in \mathbb{R}^N, y \in \mathbb{R}^M$ |
| Regular derivative: | Derivative is **Gradient**: | Derivative is **Jacobian**: |
| $\frac{\partial y}{\partial x} \in \mathbb{R}$ | $\frac{\partial y}{\partial x} \in \mathbb{R}^N \quad \left(\frac{\partial y}{\partial x}\right)_n = \frac{\partial y}{\partial x_n}$ | $\frac{\partial y}{\partial x} \in \mathbb{R}^{N \times M} \quad \left(\frac{\partial y}{\partial x}\right)_{n,m} = \frac{\partial y_m}{\partial x_n}$ |
| If x changes by a small amount, how much will y change? | For each element of x, if it changes by a small amount then how much will y change? | For each element of x, if it changes by a small amount then how much will each element of y change? |

# Backprop with Vectors

$x$

$y$

**f**

$z$

Loss L still a scalar!

# Backprop with Vectors

$D_x$ $x$

$D_y$ $y$

**f**

Loss L still a scalar!

$z$ $D_z$

# Backprop with Vectors

$D_x$ $x$

$D_y$ $y$

f

Loss L still a scalar!

$z$ $D_z$

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

"Upstream gradient"

# Backprop with Vectors

$D_x$ $x$

$D_y$ $y$

$f$

Loss L still a scalar!

$z$ $D_z$

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

"Upstream gradient"

For each element of z, how much does it influence L?

# Backprop with Vectors



"local gradients"

Loss L still a scalar!

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

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

$$\frac{\partial z}{\partial y}$$

$$D_x \quad x$$

$$z \quad D_z$$

"Downstream gradients"

$$D_y \quad y$$

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

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

"Upstream gradient"

For each element of z, how much does it influence L?

# Backprop with Vectors



$D_x$ $x$

"Downstream gradients"

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

$D_y$ $y$

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

"local gradients"

$\frac{\partial z}{\partial x}$ $[D_x \times D_z]$

$f$

$\frac{\partial z}{\partial y}$ $[D_y \times D_z]$

Jacobian matrices

Loss L still a scalar!

$z$ $D_z$

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

"Upstream gradient"

For each element of z, how much does it influence L?

# Backprop with Vectors



$D_x$ $\boxed{x}$

$D_x$ $\boxed{\dfrac{\partial L}{\partial x}} = \dfrac{\partial z}{\partial x}\dfrac{\partial L}{\partial z}$

"Downstream gradients"

Matrix-vector multiply

$D_y$ $\boxed{y}$

$\boxed{\dfrac{\partial L}{\partial y}} = \dfrac{\partial z}{\partial y}\dfrac{\partial L}{\partial z}$

$D_y$

"local gradients"

$\boxed{\dfrac{\partial z}{\partial x}}$ $[D_x \times D_z]$

**f**

$\boxed{\dfrac{\partial z}{\partial y}}$ $[D_y \times D_z]$

Jacobian matrices

Loss L still a scalar!

$\boxed{z}$ $D_z$

$\boxed{\dfrac{\partial L}{\partial z}}$ $D_z$

"Upstream gradient"

For each element of z, how much does it influence L?

Gradients of variables wrt loss have same dims as the original variable



$D_x$ $x$

$D_x$ $\dfrac{\partial L}{\partial x}$

$D_y$ $y$

$D_y$ $\dfrac{\partial L}{\partial y}$

f

Loss L still a scalar!

$z$ $D_z$

$\dfrac{\partial L}{\partial z}$ $D_z$

"Upstream gradient"

For each element of z, how much does it influence L?

# Backprop with Vectors

4D input x:

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

f(x) = max(0,x)
*(elementwise)*

4D output z:

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

# Backprop with Vectors

4D input x:

[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

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

4D output z:

[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

4D dL/dz:

[ 4 ]
[ -1 ]
[ 5 ]
[ 9 ]

Upstream gradient

# Backprop with Vectors

4D input x:

[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

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

4D output z:

[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

Jacobian dz/dx

[ 1 0 0 0 ]
[ 0 0 0 0 ]
[ 0 0 1 0 ]
[ 0 0 0 0 ]

4D dL/dz:

[ 4 ]
[ -1 ]
[ 5 ]
[ 9 ]

Upstream
gradient

# Backprop with Vectors

4D input x:

[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

f(x) = max(0,x)
*(elementwise)*

4D output z:

[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

[dz/dx] [dL/dz]

[ 1 0 0 0 ] [ 4 ]
[ 0 0 0 0 ] [ -1 ]
[ 0 0 1 0 ] [ 5 ]
[ 0 0 0 0 ] [ 9 ]

4D dL/dz:

[ 4 ]
[ -1 ]
[ 5 ]
[ 9 ]

Upstream
gradient

# Backprop with Vectors

4D input x:

[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

f(x) = max(0,x)
*(elementwise)*

4D output z:

[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

4D dL/dx:

[ 4 ]
[ 0 ]
[ 5 ]
[ 0 ]

[dz/dx] [dL/dz]

[ 1 0 0 0 ] [ 4 ]
[ 0 0 0 0 ] [ -1 ]
[ 0 0 1 0 ] [ 5 ]
[ 0 0 0 0 ] [ 9 ]

4D dL/dz:

[ 4 ]
[ -1 ]
[ 5 ]
[ 9 ]

Upstream gradient

# Backprop with Vectors

4D input x:

[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

**f(x) = max(0,x)**
*(elementwise)*

4D output z:

[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

Jacobian is **sparse**: off-diagonal entries always zero! Never **explicitly** form Jacobian -- instead use **implicit** multiplication

4D dL/dx:          [dz/dx] [dL/dz]          4D dL/dz:

[ 4 ]     ←     [ 1 0 0 0 ] [ 4 ]     ←     [ 4 ]     ←
[ 0 ]     ←     [ 0 0 0 0 ] [ -1 ]     ←     [ -1 ]     ←
[ 5 ]     ←     [ 0 0 1 0 ] [ 5 ]     ←     [ 5 ]     ←
[ 0 ]     ←     [ 0 0 0 0 ] [ 9 ]     ←     [ 9 ]     ←

Upstream gradient

# Backprop with Vectors

4D input x:

[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

f(x) = max(0,x)
*(elementwise)*

4D output z:

[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

Jacobian is **sparse**: off-diagonal entries always zero! Never **explicitly** form Jacobian -- instead use **implicit** multiplication

4D dL/dx:

[ 4 ] ←
[ 0 ] ←
[ 5 ] ←
[ 0 ] ←

[dz/dx] [dL/dz]

$$\left(\frac{\partial L}{\partial x}\right)_i = \begin{cases} \left(\frac{\partial L}{\partial z}\right)_i & \text{if } x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

4D dL/dz:

← [ 4 ] ←
← [ -1 ] ←
← [ 5 ] ←
← [ 9 ] ←

Upstream gradient

# Backprop with Matrices (or Tensors)

dL/dx always has the
same shape as x!

$[D_x \times M_x]$ $x$

$z$ $[D_z \times M_z]$

**Matrix-vector
multiply**

$[D_y \times M_y]$ $y$

**f**

Jacobian
matrices

# Backprop with Matrices (or Tensors)

Loss L still a scalar!

dL/dx always has the same shape as x!

$[D_x \times M_x]$ $x$

$[D_x \times M_x]$ $\dfrac{\partial L}{\partial x} = \dfrac{\partial z}{\partial x}\dfrac{\partial L}{\partial x}$

"Downstream gradients"

Matrix-vector multiply

$[D_y \times M_y]$ $y$

$\dfrac{\partial L}{\partial y} = \dfrac{\partial z}{\partial y}\dfrac{\partial L}{\partial z}$

$[D_y \times M_y]$

f

Jacobian matrices

$z$ $[D_z \times M_z]$

$\dfrac{\partial L}{\partial z}$ $[D_z \times M_z]$

"Upstream gradient"
For each element of z, how much does it influence L?

# Backprop with Matrices (or Tensors)

$[D_x \times M_x]$ $x$

$[D_x \times M_x]$ $\dfrac{\partial L}{\partial x} = \dfrac{\partial z}{\partial x}\dfrac{\partial L}{\partial x}$

"Downstream gradients"

Matrix-vector multiply

$[D_y \times M_y]$ $y$

$\dfrac{\partial L}{\partial y} = \dfrac{\partial z}{\partial y}\dfrac{\partial L}{\partial z}$

$[D_y \times M_y]$

"local gradients"

$\dfrac{\partial z}{\partial x}$

$\dfrac{\partial z}{\partial y}$

Jacobian matrices

For each element of y, how much does it influence each element of z?

Loss L still a scalar!

dL/dx always has the same shape as x!

$z$ $[D_z \times M_z]$

$\dfrac{\partial L}{\partial z}$ $[D_z \times M_z]$

"Upstream gradient"
For each element of z, how much does it influence L?

# Backprop with Matrices (or Tensors)

Loss L still a scalar!

dL/dx always has the same shape as x!

$[D_x \times M_x]$  $x$

$[D_x \times M_x]$  $\dfrac{\partial L}{\partial x} = \dfrac{\partial z}{\partial x}\dfrac{\partial L}{\partial x}$

"Downstream gradients"

Matrix-vector multiply

"local gradients"

$\dfrac{\partial z}{\partial x}$  $[(D_x \times M_x) \times (D_z \times M_z)]$

$z$  $[D_z \times M_z]$

$\dfrac{\partial z}{\partial y}$  $[(D_y \times M_y) \times (D_z \times M_z)]$

$[D_y \times M_y]$  $y$

$\dfrac{\partial L}{\partial y} = \dfrac{\partial z}{\partial y}\dfrac{\partial L}{\partial z}$

$[D_y \times M_y]$

Jacobian matrices

$\dfrac{\partial L}{\partial z}$  $[D_z \times M_z]$

"Upstream gradient"

For each element of z, how much does it influence L?

For each element of y, how much does it influence each element of z?

# Backprop with Matrices

x: [N×D]

[ 2  **1**  -3 ]
[ -3  4   2 ]

w: [D×M]

[ 3  2  1  -1]
[ 2  1  3  2]
[ 3  2  1  -2]

Matrix Multiply

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

y: [N×M]

[ **13**  **9**  **-2**  **-6** ]
[  5  2  17  1 ]

dL/dy: [N×M]

[  2  3 -3  9 ]
[ -8  1  4  6 ]

Also see derivation in the course notes:
http://cs231n.stanford.edu/handouts/linear-backprop.pdf

# Backprop with Matrices

x: [N×D]

[ 2  **1**  -3 ]
[ -3  4   2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

**Jacobians**:
dy/dx: [(N×D)×(N×M)]
dy/dw: [(D×M)×(N×M)]

y: [N×M]

[**13  9  -2  -6** ]
[ 5  2  17  1 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[ -8  1  4  6 ]

For a neural net we may have
N=64, D=M=4096
Each Jacobian takes 256 GB of memory!
Must work with them implicitly!

# Backprop with Matrices

x: [N×D]

[ 2  **1**  -3 ]
[ -3  4   2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

**Q**: What parts of y are affected by one element of x?

y: [N×M]

[ **13  9  -2  -6** ]
[  5  2  17  1 ]

dL/dy: [N×M]

[  2  3 -3  9 ]
[ -8  1  4  6 ]

# Backprop with Matrices

x: [N×D]

[ 2  **1** -3 ]
[-3  4  2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

**Matrix Multiply**

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

**Q**: What parts of y are affected by one element of x?

**A**: $\boxed{x_{n,d}}$ affects the whole row $y_{n,\cdot}$

$$\frac{\partial L}{\partial x_{n,d}} = \sum_m \frac{\partial L}{\partial y_{n,m}} \frac{\partial y_{n,m}}{\partial x_{n,d}}$$

y: [N×M]

[**13  9  -2  -6**]
[ 5  2  17  1 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[-8  1  4  6 ]

# Backprop with Matrices

x: [N×D]

[ 2  **1** -3 ]
[ -3  4   2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

y: [N×M]

[ **13**  **9**  **-2**  **-6** ]
[  5   2   17   1 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[ -8  1  4  6 ]

**Q**: What parts of y are affected by one element of x?

**A**: $x_{n,d}$ affects the whole row $y_{n,\cdot}$

**Q**: How much does $x_{n,d}$ affect $y_{n,m}$?

$$\frac{\partial L}{\partial x_{n,d}} = \sum_m \frac{\partial L}{\partial y_{n,m}} \frac{\partial y_{n,m}}{\partial x_{n,d}}$$

# Backprop with Matrices

x: [N×D]

[ 2  **1** -3 ]
[ -3  4  2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

**Matrix Multiply**

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

y: [N×M]

[ **13  9** -2 -6 ]
[ 5  2  17  1 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[ -8  1  4  6 ]

**Q**: What parts of y are affected by one element of x?
**A**: $x_{n,d}$ affects the whole row $y_{n,\cdot}$

**Q**: How much does $x_{n,d}$ affect $y_{n,m}$?
**A**: $w_{d,m}$

$$\frac{\partial L}{\partial x_{n,d}} = \sum_m \frac{\partial L}{\partial y_{n,m}} \frac{\partial y_{n,m}}{\partial x_{n,d}} = \sum_m \frac{\partial L}{\partial y_{n,m}} w_{d,m}$$

# Backprop with Matrices

x: [N×D]

[ 2  **1** -3 ]
[ -3  4   2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

[N×D]  [N×M] [M×D]

$$\frac{\partial L}{\partial x} = \left(\frac{\partial L}{\partial y}\right) w^T$$

## Matrix Multiply

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

y: [N×M]

[ **13**  **9**  **-2**  **-6** ]
[  5   2   17   1 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[ -8  1  4  6 ]

**Q**: What parts of y are affected by one element of x?

**A**: $x_{n,d}$ affects the whole row $y_{n,\cdot}$

**Q**: How much does $x_{n,d}$ affect $y_{n,m}$?

**A**: $w_{d,m}$

$$\frac{\partial L}{\partial x_{n,d}} = \sum_m \frac{\partial L}{\partial y_{n,m}} \frac{\partial y_{n,m}}{\partial x_{n,d}} = \sum_m \frac{\partial L}{\partial y_{n,m}} w_{d,m}$$

# Backprop with Matrices

x: [N×D]

[ 2  **1** -3 ]
[ -3  4  2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

By similar logic:

y: [N×M]

[ **13  9** -2 -6 ]
[ 5  2  17  1 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[ -8  1  4  6 ]

[N×D]  [N×M] [M×D]

$$\frac{\partial L}{\partial x} = \left(\frac{\partial L}{\partial y}\right) w^T$$

[D×M]  [D×N] [N×M]

$$\frac{\partial L}{\partial w} = x^T \left(\frac{\partial L}{\partial y}\right)$$

These formulas are easy to remember: they are the only way to make shapes match up!

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$

$\Downarrow \quad \Downarrow$

$\in \mathbb{R}^n \in \mathbb{R}^{n \times n}$

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} W$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} x$$



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

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

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} \text{W}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} \text{x}$$

$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$*$  →  L2  →  0.116

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

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

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} \mathbf{W}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$$

$\mathbf{x}$

$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$*$

L2

0.116

1.00

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

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

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} \mathbf{W}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} \mathbf{x}$$



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

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

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

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

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix}_{\mathbf{W}}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}_{\mathbf{x}}$$

$$*$$

$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

$$L2$$

0.116

1.00

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

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

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

$\nabla_q f = 2q$

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} \mathbf{W}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} \mathbf{x}$$



$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

$*$    L2

0.116
1.00

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

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

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

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} \text{W}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} \text{x}$$

$*$

$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

L2

0.116

1.00

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

$$\frac{\partial f}{\partial W_{i,j}} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial W_{i,j}}$$

$$= \sum_k (2q_k)(\mathbf{1}_{k=i} x_j)$$

$$= 2q_i x_j$$

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

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

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n}(W \cdot x)_i^2$

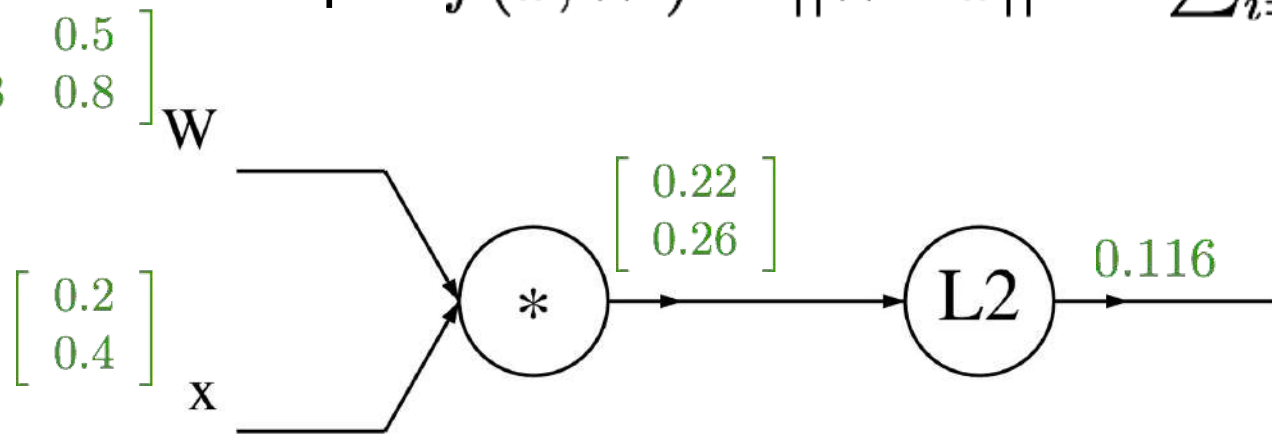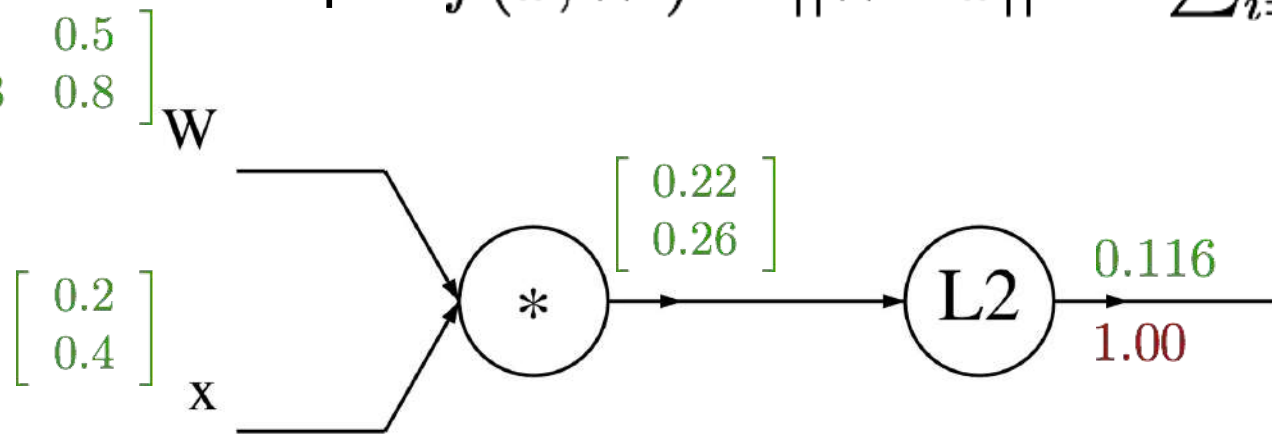$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} W$$

$$\begin{bmatrix} 0.088 & 0.176 \\ 0.104 & 0.208 \end{bmatrix}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$$

x

* 

$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

L2

0.116
1.00

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

$$\frac{\partial f}{\partial W_{i,j}} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial W_{i,j}}$$

$$= \sum_k (2q_k)(\mathbf{1}_{k=i} x_j)$$

$$= 2q_i x_j$$

$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$
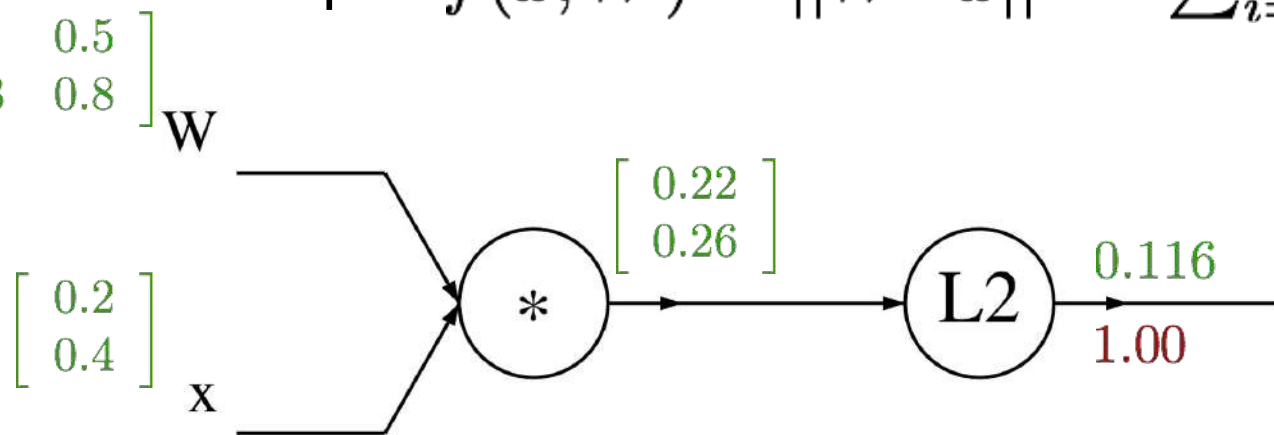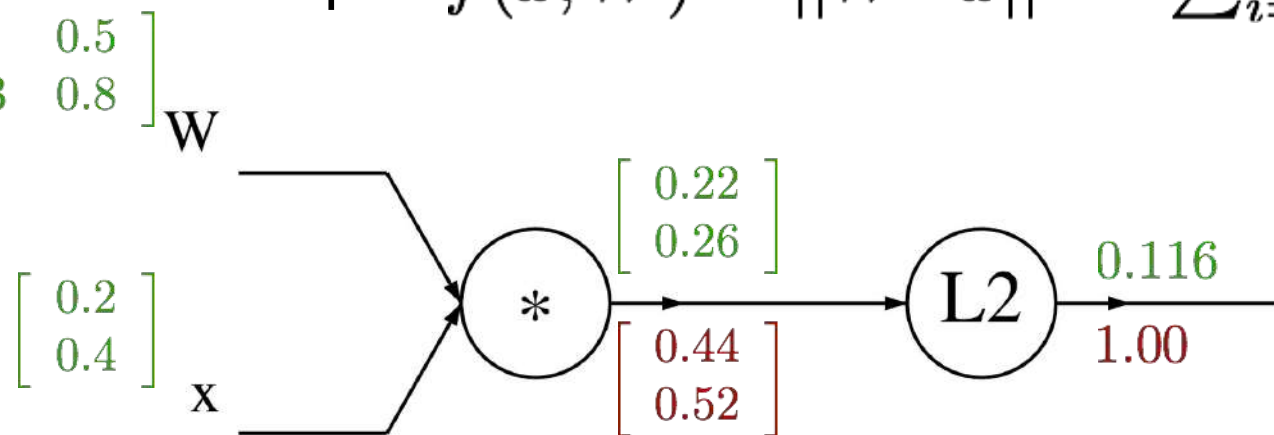
$$f(q) = ||q||^2 = q_1^2 + \cdots + q_n^2$$

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} \text{W}$$

$$\begin{bmatrix} 0.088 & 0.176 \\ 0.104 & 0.208 \end{bmatrix}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$$

x

$$\boxed{\nabla_W f = 2q \cdot x^T}$$

$*$

$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

L2

0.116

1.00

$$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

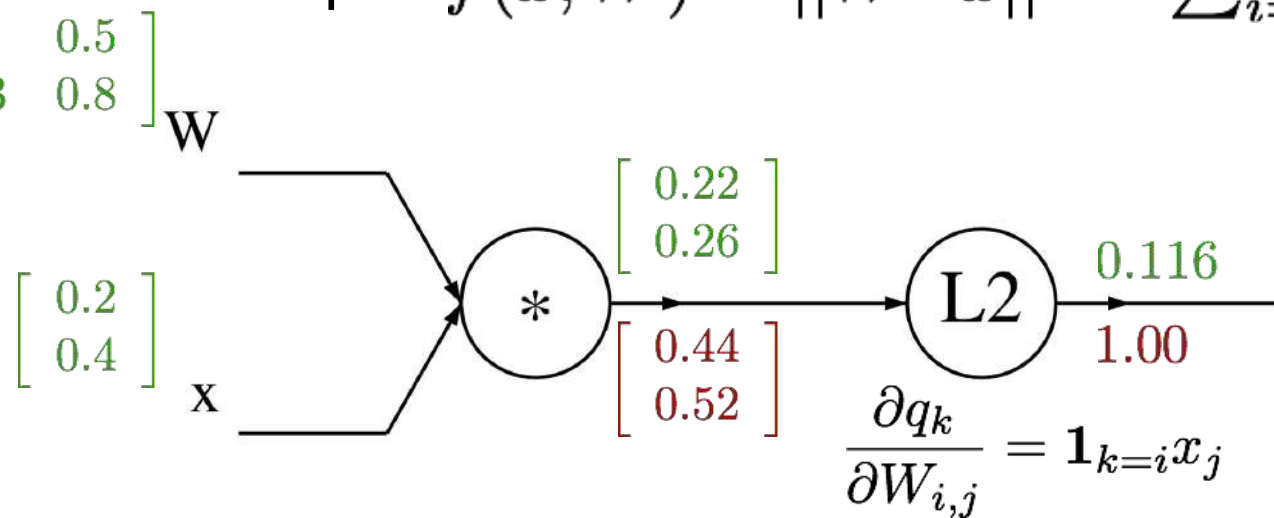$$\frac{\partial q_k}{\partial W_{i,j}} = \mathbf{1}_{k=i} x_j$$

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

$$f(q) = ||q||^2 = q_1^2 + \cdots + q_n^2$$
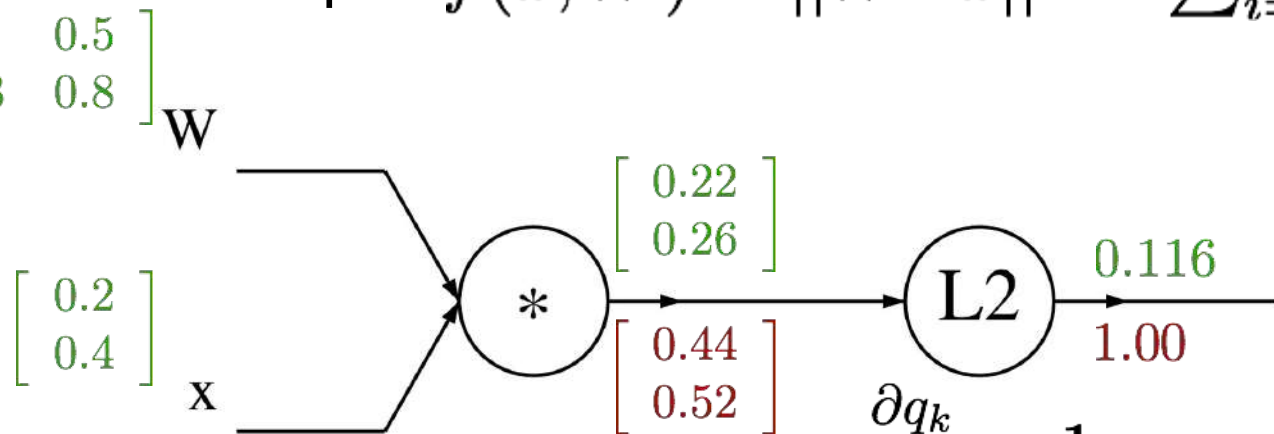
$$\frac{\partial f}{\partial W_{i,j}} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial W_{i,j}}$$

$$= \sum_k (2q_k)(\mathbf{1}_{k=i} x_j)$$

$$= 2q_i x_j$$

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} \text{W}$$

$$\begin{bmatrix} 0.088 & 0.176 \\ 0.104 & 0.208 \end{bmatrix}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$$

$$\text{x}$$

$\boxed{\nabla_W f = 2q \cdot x^T}$

$$* \qquad \begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

$$\text{L2} \qquad \begin{matrix} 0.116 \\ 1.00 \end{matrix}$$

Always check: The gradient with respect to a variable should have the same shape as the variable

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

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

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

$$\frac{\partial f}{\partial W_{i,j}} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial W_{i,j}}$$

$$= \sum_k (2q_k)(\mathbf{1}_{k=i} x_j)$$

$$= 2q_i x_j$$

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n}(W \cdot x)_i^2$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} \text{W}$$

$$\begin{bmatrix} 0.088 & 0.176 \\ 0.104 & 0.208 \end{bmatrix}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$$

x

*

$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

L2

0.116

1.00

$$\frac{\partial q_k}{\partial x_i} = W_{k,i}$$

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

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

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$
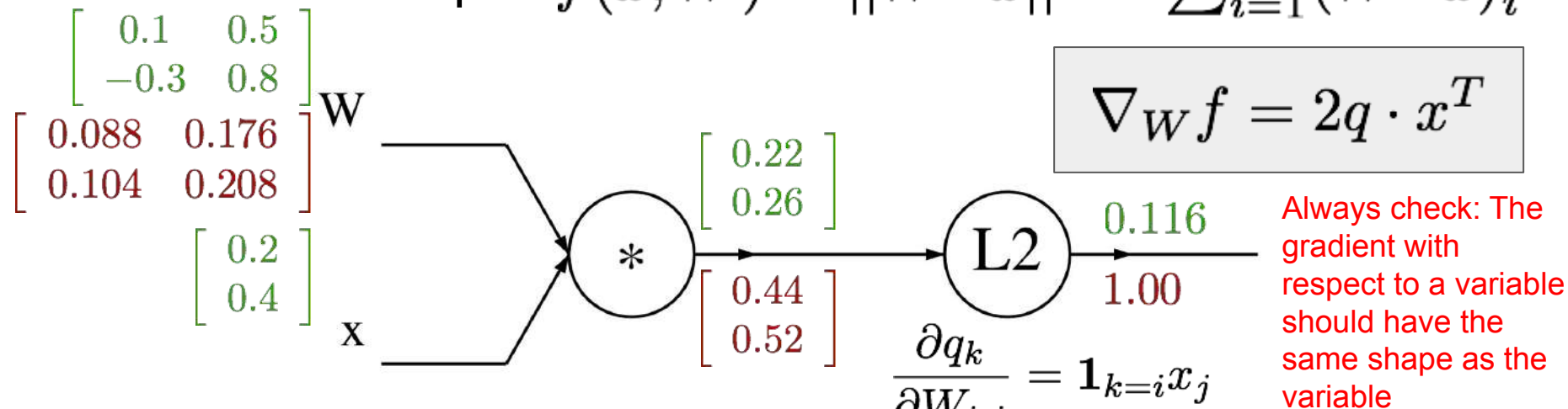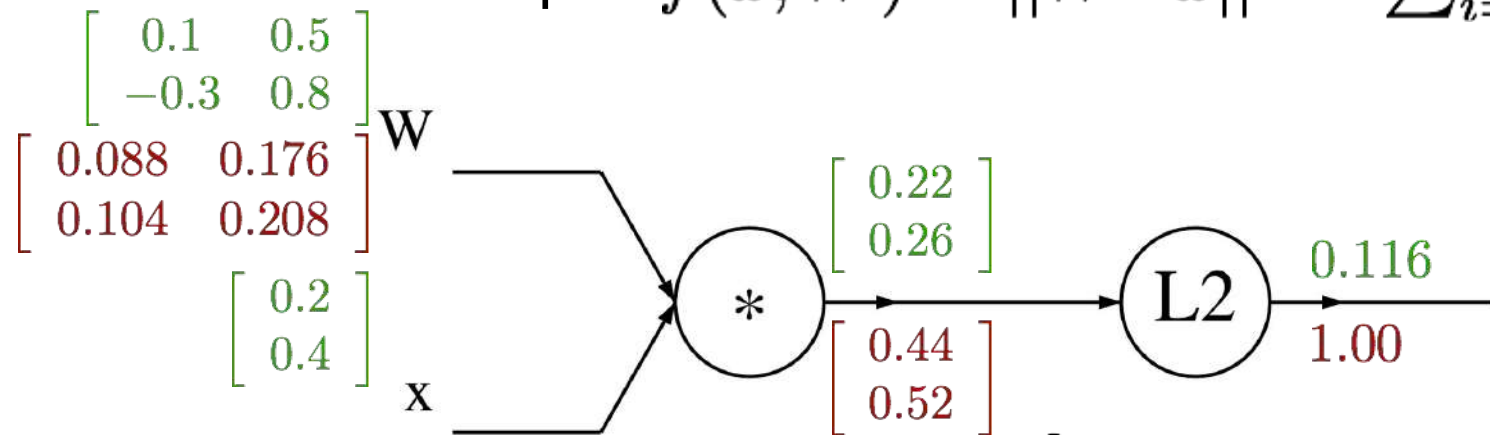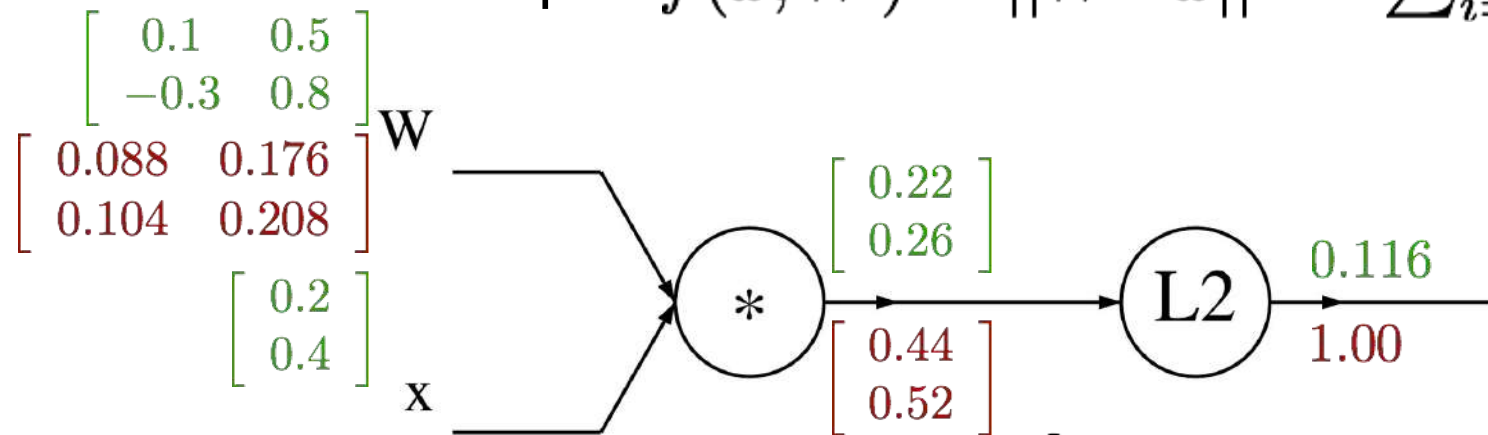
$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} \text{W}$$

$$\begin{bmatrix} 0.088 & 0.176 \\ 0.104 & 0.208 \end{bmatrix}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$$

x

$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

*

$$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

L2

0.116

1.00

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

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

$\dfrac{\partial q_k}{\partial x_i} = W_{k,i}$

$\dfrac{\partial f}{\partial x_i} = \sum_k \dfrac{\partial f}{\partial q_k} \dfrac{\partial q_k}{\partial x_i}$

$\quad\quad = \sum_k 2q_k W_{k,i}$

A vectorized example: $f(x, W) = ||W \cdot x||^2 = \sum_{i=1}^{n} (W \cdot x)_i^2$

$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix}$ W

$\begin{bmatrix} 0.088 & 0.176 \\ 0.104 & 0.208 \end{bmatrix}$

$\boxed{\nabla_x f = 2W^T \cdot q}$

$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$

$\begin{bmatrix} -0.112 \\ 0.636 \end{bmatrix}$ x

$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$

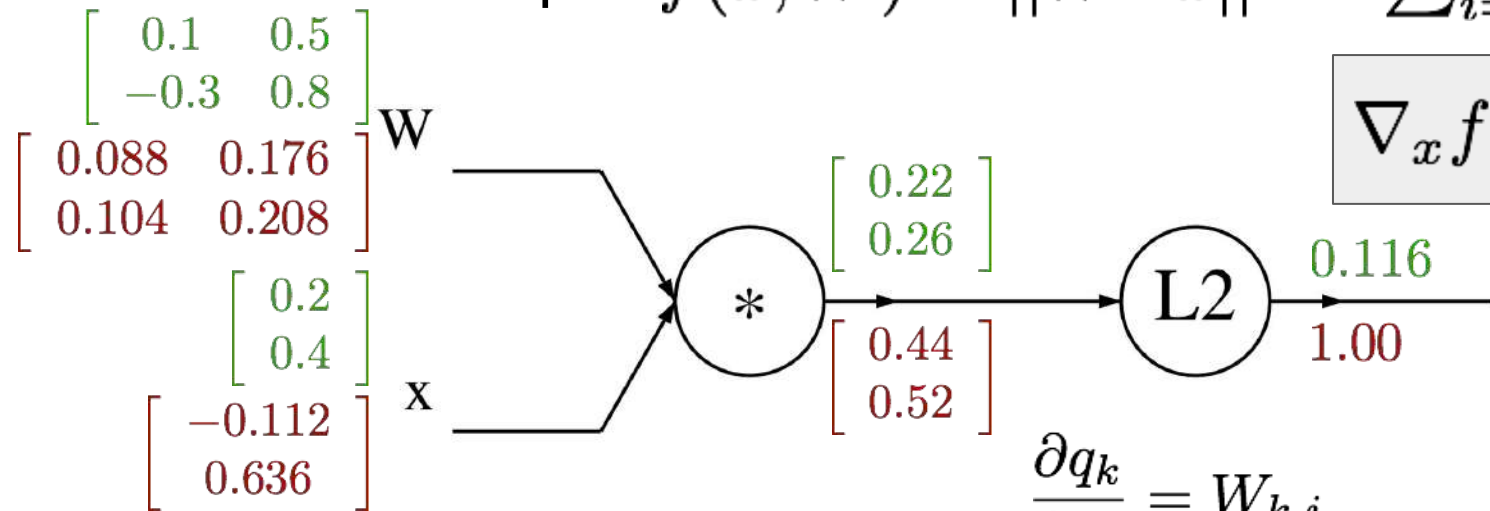$*$ → L2 → 0.116 / 1.00

$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$

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

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

$\dfrac{\partial q_k}{\partial x_i} = W_{k,i}$

$\dfrac{\partial f}{\partial x_i} = \sum_k \dfrac{\partial f}{\partial q_k} \dfrac{\partial q_k}{\partial x_i}$

$= \sum_k 2q_k W_{k,i}$

In discussion section: A matrix example...

$$z_1 = XW_1$$
$$h_1 = \text{ReLU}(z_1)$$
$$\hat{y} = h_1 W_2$$
$$L = \|\hat{y}\|_2^2$$

$$\frac{\partial L}{\partial W_2} = \quad \textbf{?}$$
$$\frac{\partial L}{\partial W_1} = \quad \textbf{?}$$