

CS 360: Machine Learning

Sara Mathieson, Sorelle Friedler

Spring 2024



HVERFORD
COLLEGE

Admin

- **Lab 7** check in today
 - Should be finished with the fully connected network
- **Project proposal** due April 8 (short)
- I may receive a call

Outline for April 2

- Finish CNNs
- Neural network regularization
- Backpropagation

Outline for April 2

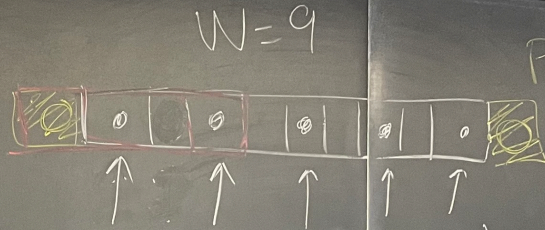
- **Finish CNNs**
- Neural network regularization
- Backpropagation

CNN filters

- W = input width
- F = filter size
(almost always odd)
- P = padding on one side
- S = stride

Typical
 $\Rightarrow S=1$

how much to shift the filter each pass



$$W - F + 1$$

$$9 - 3 + 1 = 7$$

$$\boxed{W - F + 2P + 1}$$

new input shape

$P=1$



output size?

$$\boxed{\frac{W - F + 2P}{S} + 1}$$

without padding

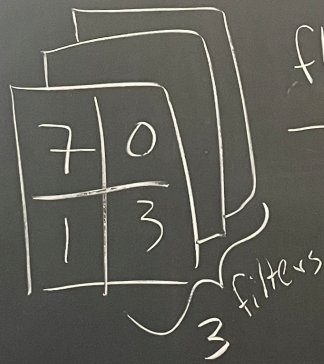
formula for output shape

$$9 - 3 + 2 \cdot 1 + 1 = \boxed{9} \text{ with padding}$$

$$S = 2$$

$$\frac{9 - 3 + 2 \cdot 1}{2} + 1 \Rightarrow \boxed{5} \text{ output shape}$$

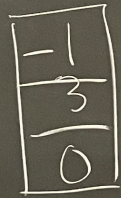
FC layers
after conv
layers



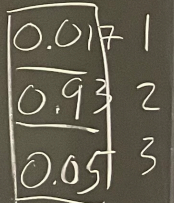
flatten



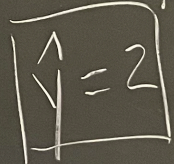
3 classes



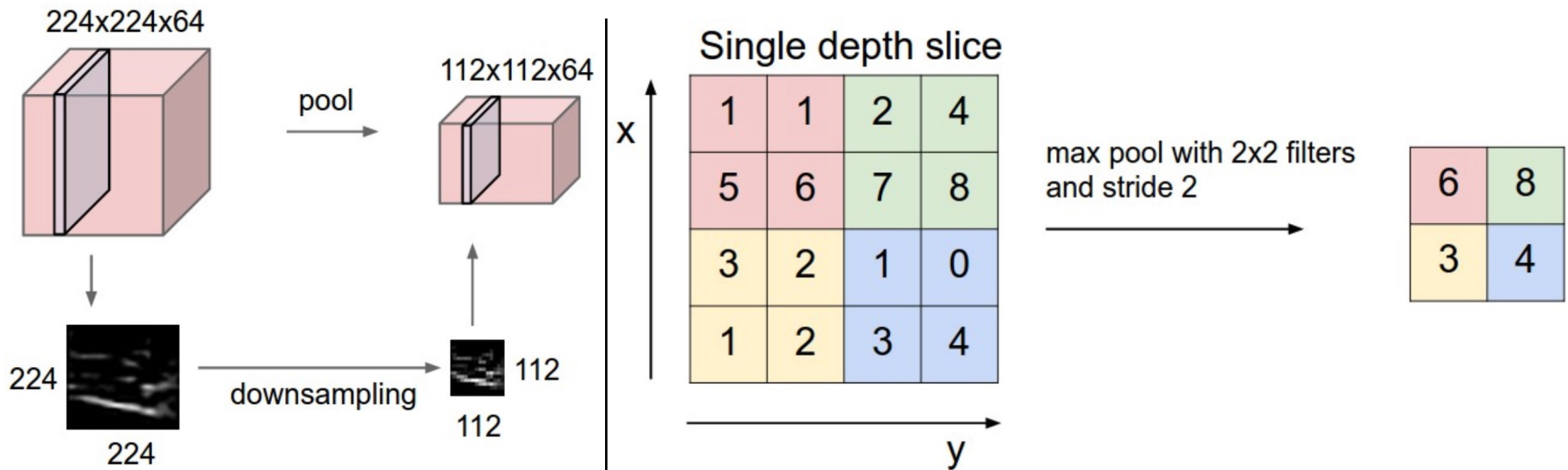
softmax



Output probs



Pooling



Pooling layer downsamples the volume spatially, independently in each depth slice of the input volume. **Left:** In this example, the input volume of size $[224 \times 224 \times 64]$ is pooled with filter size 2, stride 2 into output volume of size $[112 \times 112 \times 64]$. Notice that the volume depth is preserved. **Right:** The most common downsampling operation is max, giving rise to **max pooling**, here shown with a stride of 2. That is, each max is taken over 4 numbers (little 2×2 square).

Handout 18

- (a) Which steps require parameter learning? (out of CONV, RELU, POOL, FLATTEN, FC)
- (b) First layer params
- (c) Second layer params
- (d) Third layer params
- (e) Total # params

Handout 18

(a) Which steps require parameter learning? (out of CONV, RELU, POOL, FLATTEN, FC)

CONV, FC

(b) First layer params

(c) Second layer params

(d) Third layer params

(e) Total # params

Handout 18

(a) Which steps require parameter learning? (out of CONV, RELU, POOL, FLATTEN, FC)

CONV, FC

(b) First layer params $5*5*3*20 + 20 = 1520$

(c) Second layer params $3*3*20*10 + 10 = 1810$

(d) Third layer params $8*8*10*10 + 10 = 6410$

(e) Total # params 9740

Handout 18

(a) Which steps require parameter learning? (out of CONV, RELU, POOL, FLATTEN, FC)

CONV, FC

(b) First layer params $5*5*3*20 + 20 = 1520$

(c) Second layer params $3*3*20*10 + 10 = 1810$

(d) Third layer params $8*8*10*10 + 10 = 6410$

(e) Total # params 9740

If we had a FC with $p_1=100$ and $p_2=50$, we would have 312,860 params to learn (check this after class). CNN is much better!

Handout 19, #1

(a) $W=10$, $F=7$, $P=3$, $S=3$

(b) Draw padding

(c) Shade units where cross-correlation is performed

Handout 19, #1

(a) $W=10, F=7, P=3, S=3$

$$(10-7+6)/3 + 1 = 4 \quad (\text{output size})$$

(b) Draw padding

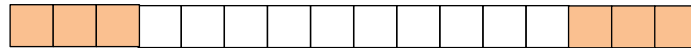
(c) Shade units where cross-correlation is performed

Handout 19, #1

(a) $W=10, F=7, P=3, S=3$

$$(10-7+6)/3 + 1 = 4 \quad (\text{output size})$$

(b) Draw padding



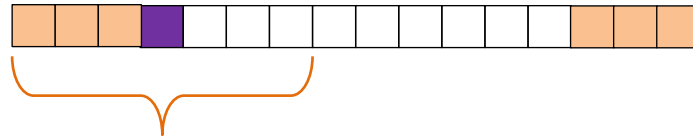
(c) Shade units where cross-correlation is performed

Handout 19, #1

(a) $W=10, F=7, P=3, S=3$

$$(10-7+6)/3 + 1 = 4 \quad (\text{output size})$$

(b) Draw padding



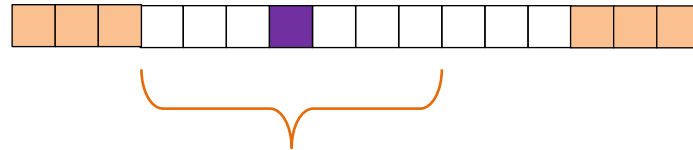
(c) Shade units where cross-correlation is performed

Handout 19, #1

(a) $W=10, F=7, P=3, S=3$

$$(10-7+6)/3 + 1 = 4 \quad (\text{output size})$$

(b) Draw padding



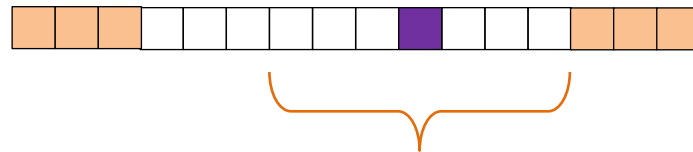
(c) Shade units where cross-correlation is performed

Handout 19, #1

(a) $W=10, F=7, P=3, S=3$

$$(10-7+6)/3 + 1 = 4 \quad (\text{output size})$$

(b) Draw padding



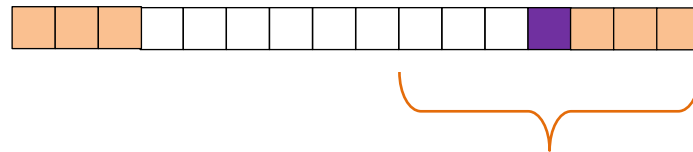
(c) Shade units where cross-correlation is performed

Handout 19, #1

(a) $W=10, F=7, P=3, S=3$

$$(10-7+6)/3 + 1 = 4 \quad (\text{output size})$$

(b) Draw padding

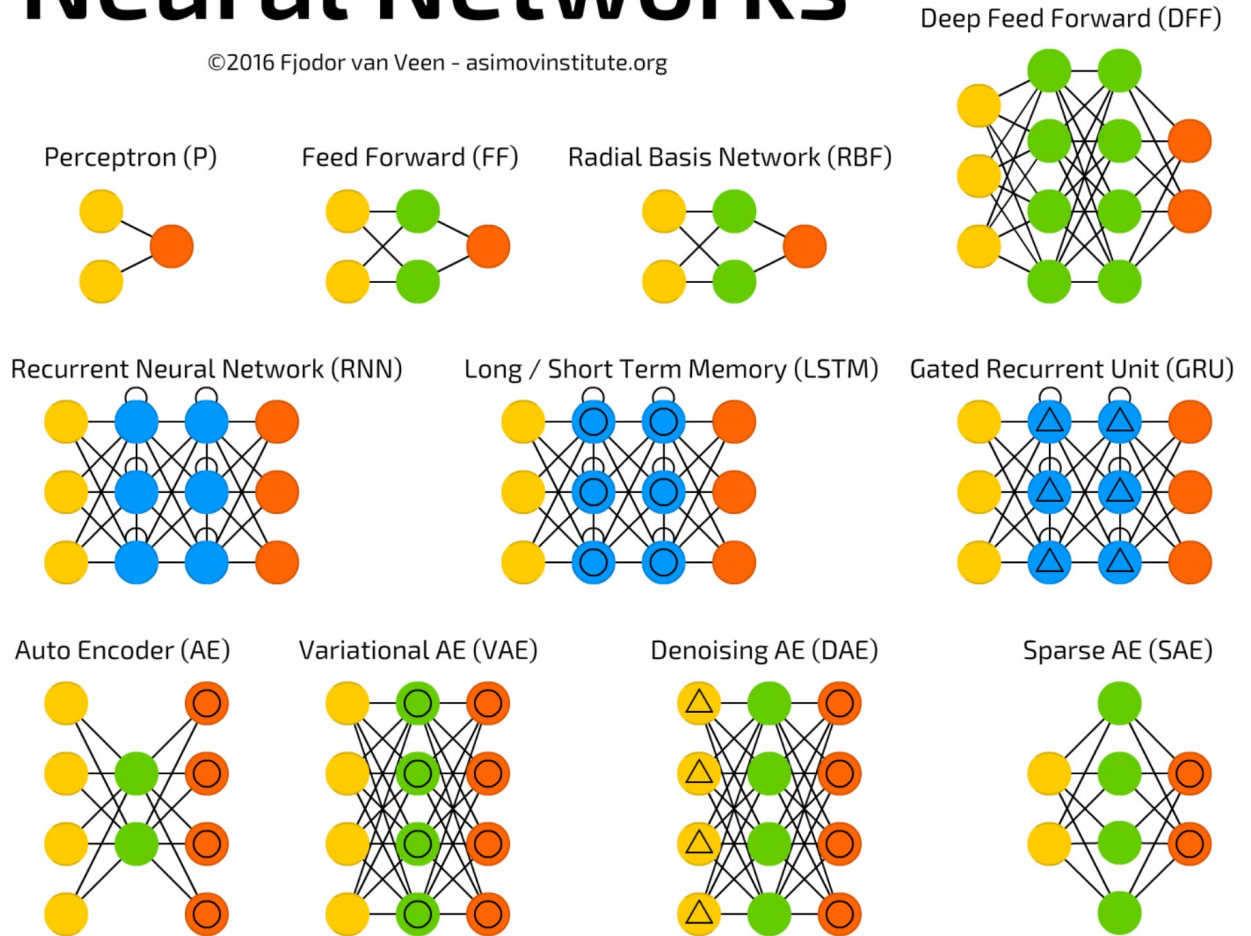


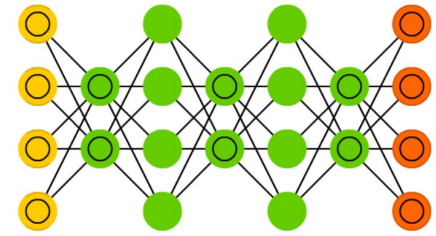
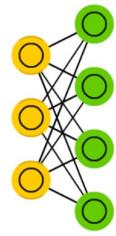
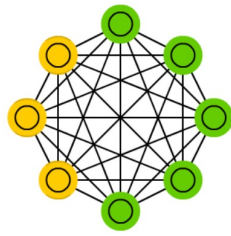
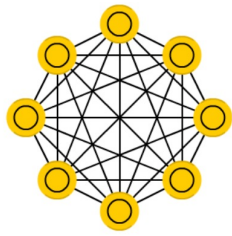
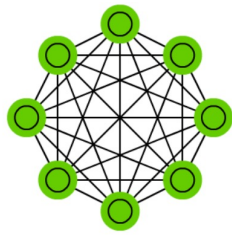
(c) Shade units where cross-correlation is performed

A mostly complete chart of Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org

- Backfed Input Cell
- Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probablistic Hidden Cell
- △ Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Different Memory Cell
- Kernel
- Convolution or Pool

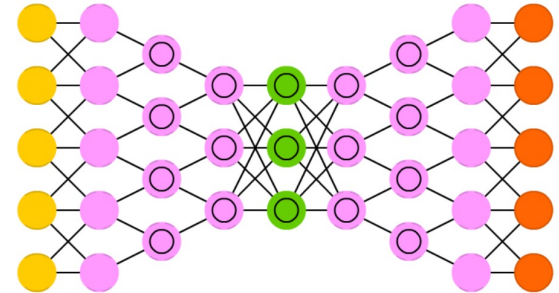
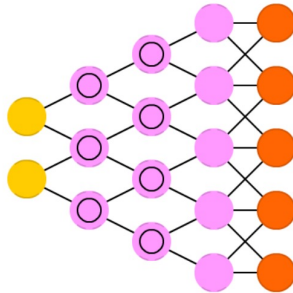
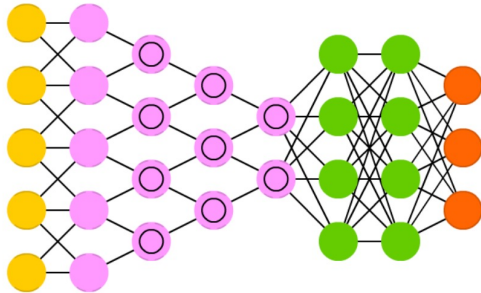




Deep Convolutional Network (DCN)

Deconvolutional Network (DN)

Deep Convolutional Inverse Graphics Network (DCIGN)

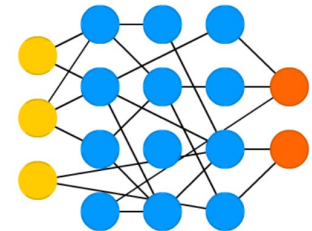
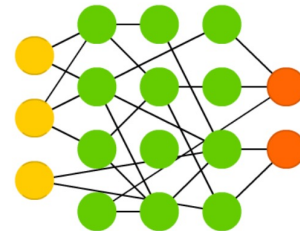
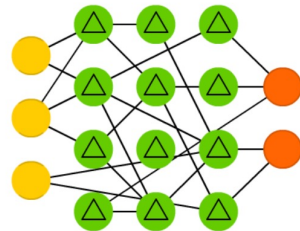
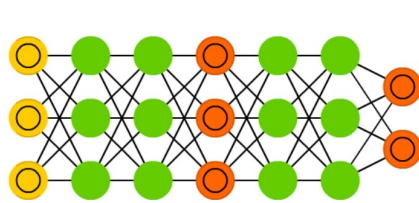


Generative Adversarial Network (GAN)

Liquid State Machine (LSM)

Extreme Learning Machine (ELM)

Echo State Network (ESN)

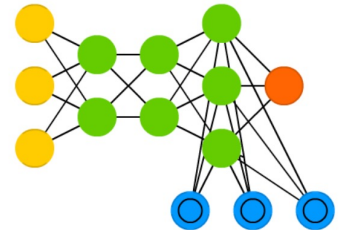
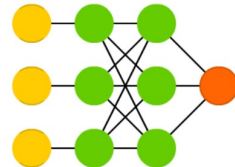
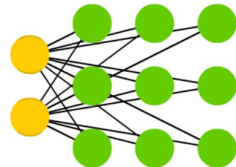
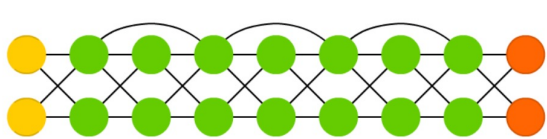


Deep Residual Network (DRN)

Kohonen Network (KN)

Support Vector Machine (SVM)

Neural Turing Machine (NTM)



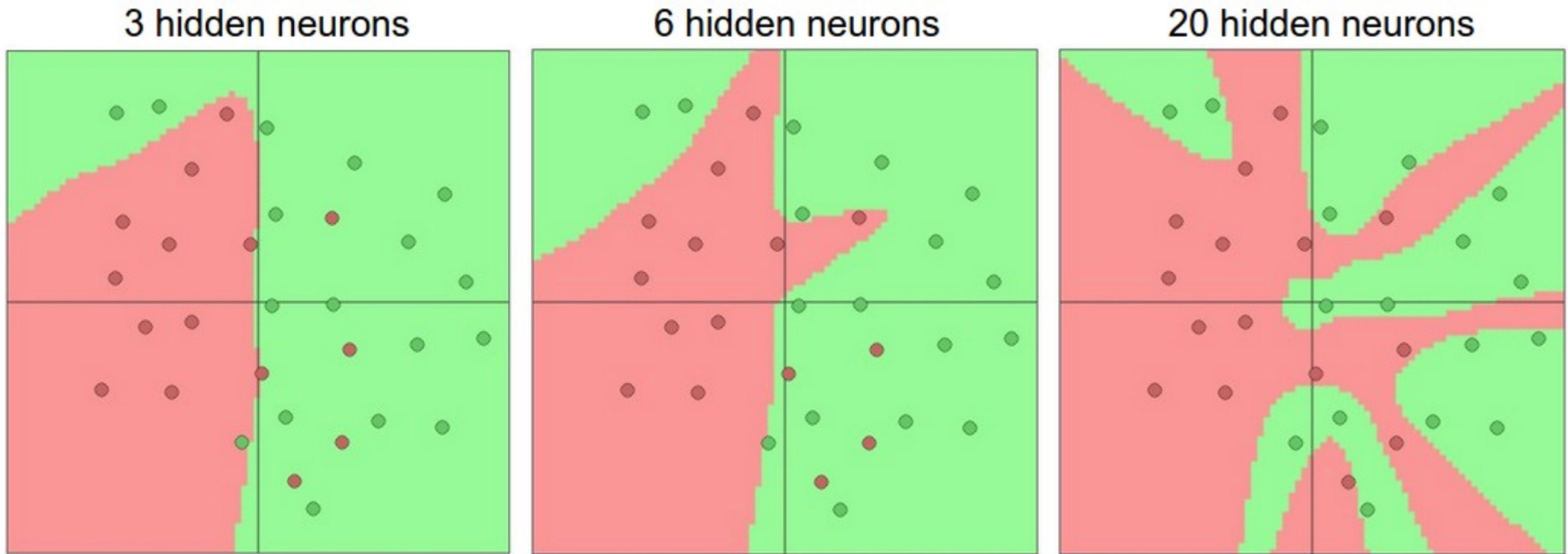
Outline for April 2

- Finish CNNs
- **Neural network regularization**
- Backpropagation

Weight initialization

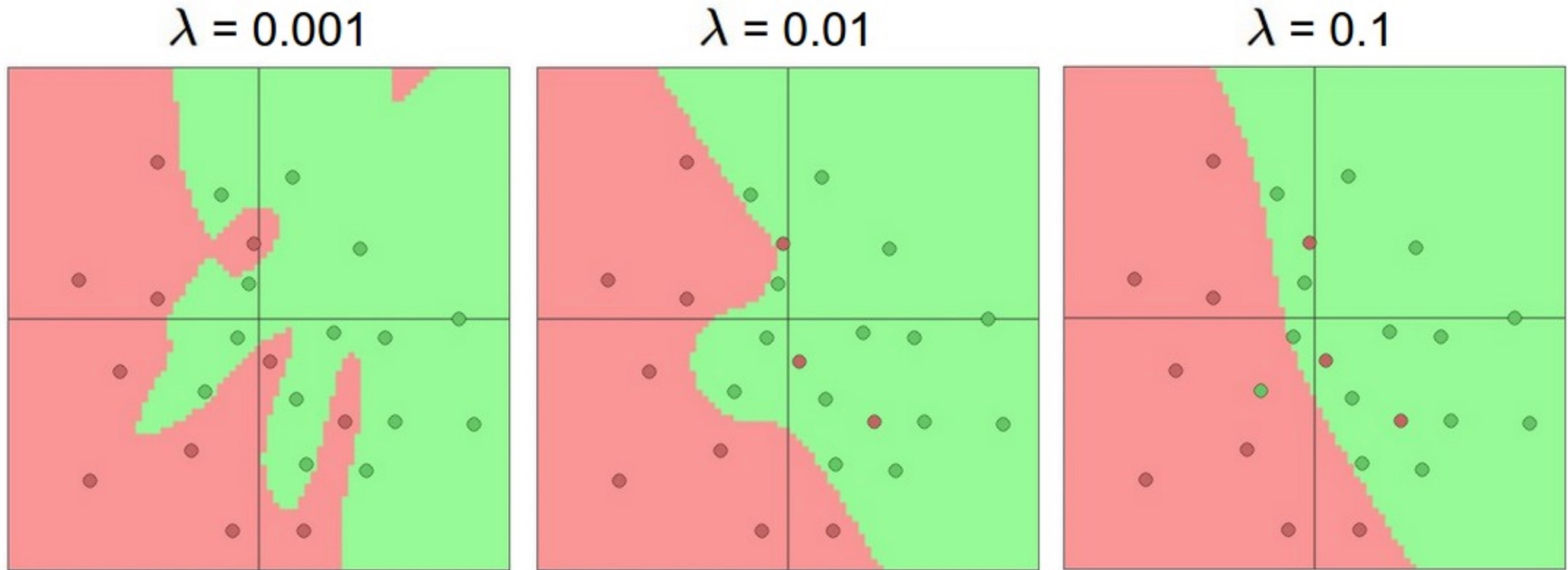
- We still have to initialize the pre-training
- All 0's initialization is bad! Causes nodes to compute the same outputs, so then the weights go through the same updates during gradient descent
- Need asymmetry! => usually use small random values

More hidden units can contribute to overfitting



Larger Neural Networks can represent more complicated functions. The data are shown as circles colored by their class, and the decision regions by a trained neural network are shown underneath. You can play with these examples in this [ConvNetsJS demo](http://cs231n.github.io/neural-networks-1/).

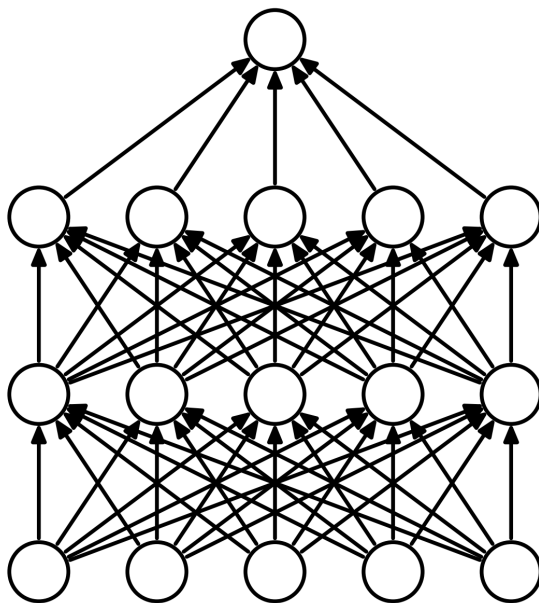
However! It is always better to use a more expressive network and regularize in other ways



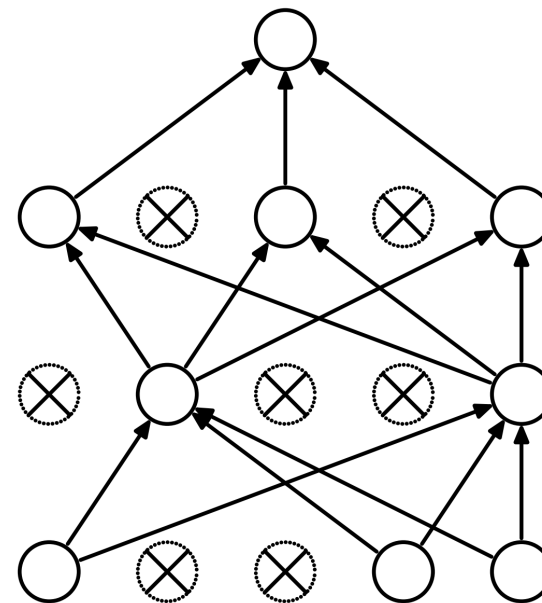
The effects of regularization strength: Each neural network above has 20 hidden neurons, but changing the regularization strength makes its final decision regions smoother with a higher regularization. You can play with these examples in this [ConvNetsJS demo](#).

One regularization approach: dropout

- Idea: keep a neuron active with some probability p , otherwise, do not send its output forward to the next layer



(a) Standard Neural Net



(b) After applying dropout.

Image and more information: “Dropout: A Simple Way to Prevent Neural Networks from Overfitting”

<http://www.cs.toronto.edu/~rsalakhu/papers/srivastava14a.pdf>

Outline for April 2

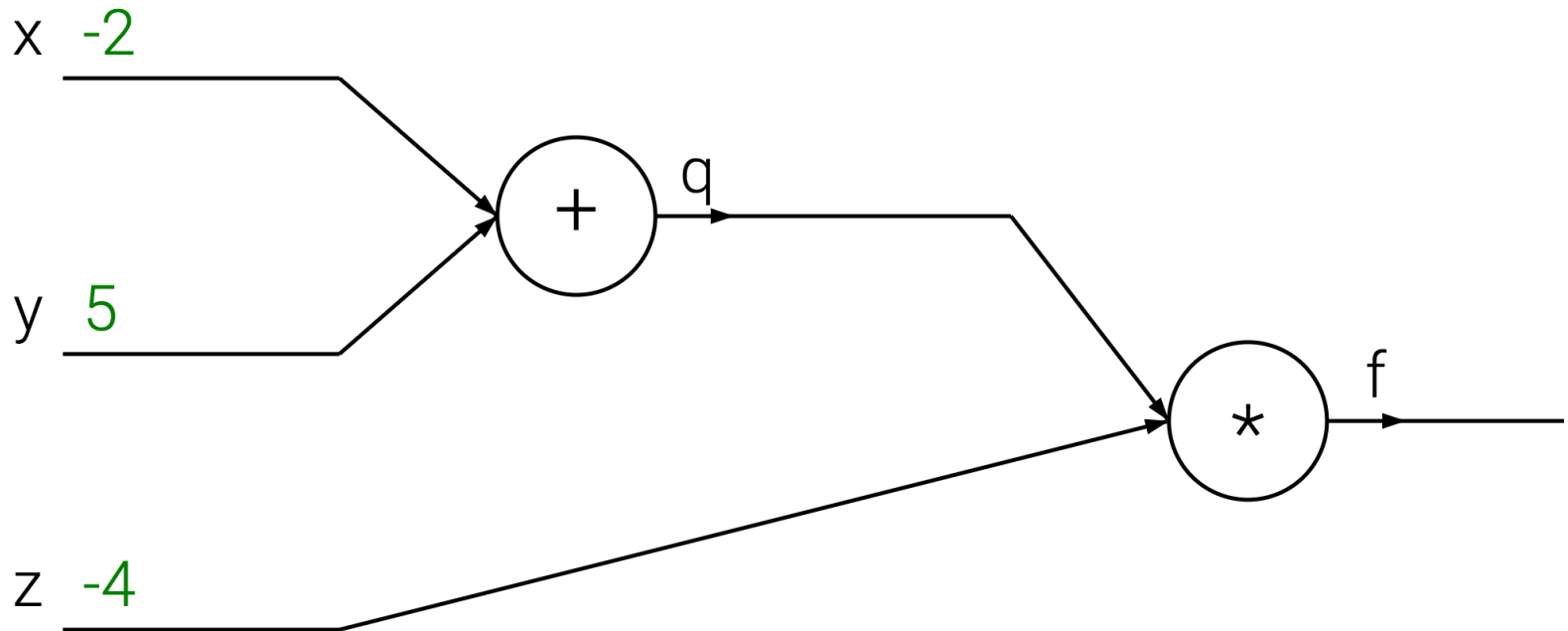
- Finish CNNs
- Neural network regularization
- **Backpropagation**

Backpropagation

- *High-level goal:* we want to know how the output depends on the input
- *Issue:* network is very complicated and overall gradient may be difficult to compute
- *Idea:* use the chain rule to compute local gradients throughout the network
- *Takeaway:* nodes can know about their value and local gradient without knowing about the network they are imbedded in

Backpropagation: Example

Forward pass: compute values

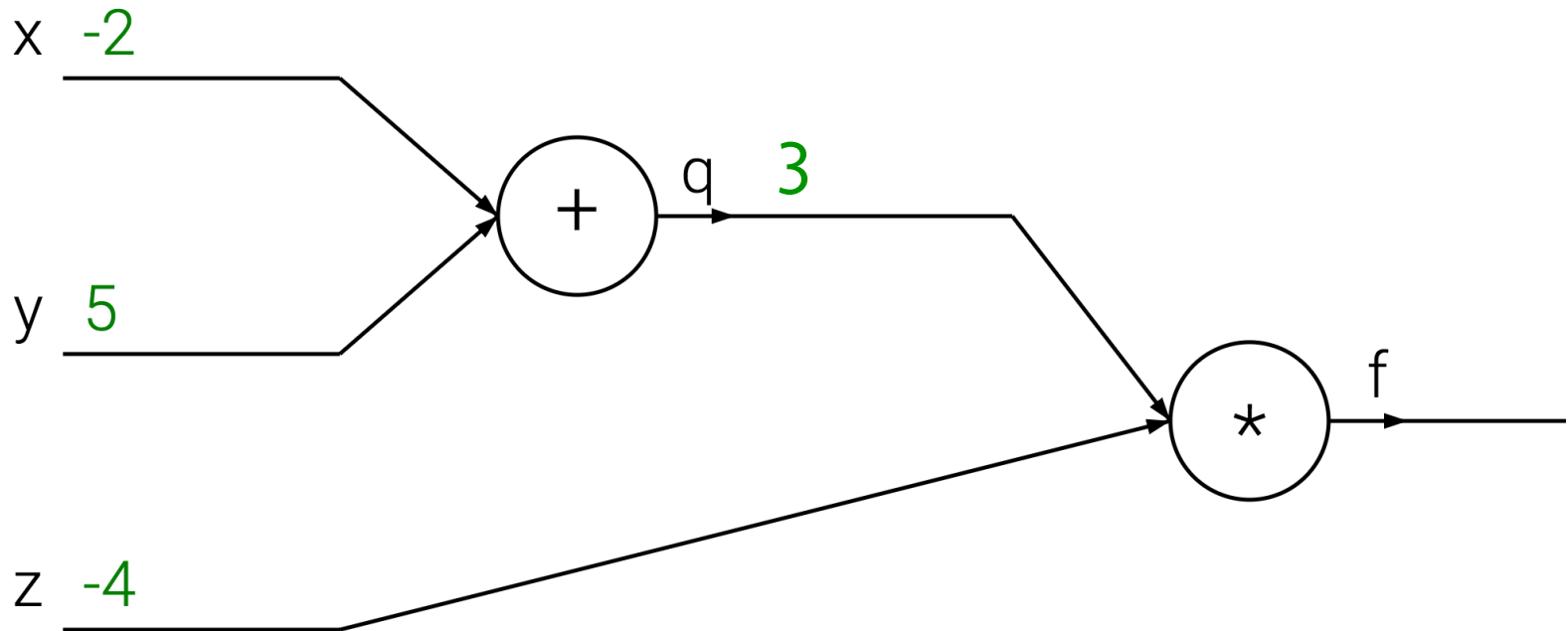


Example from:

<http://cs231n.github.io/optimization-2/>

Backpropagation: Example

Forward pass: compute values

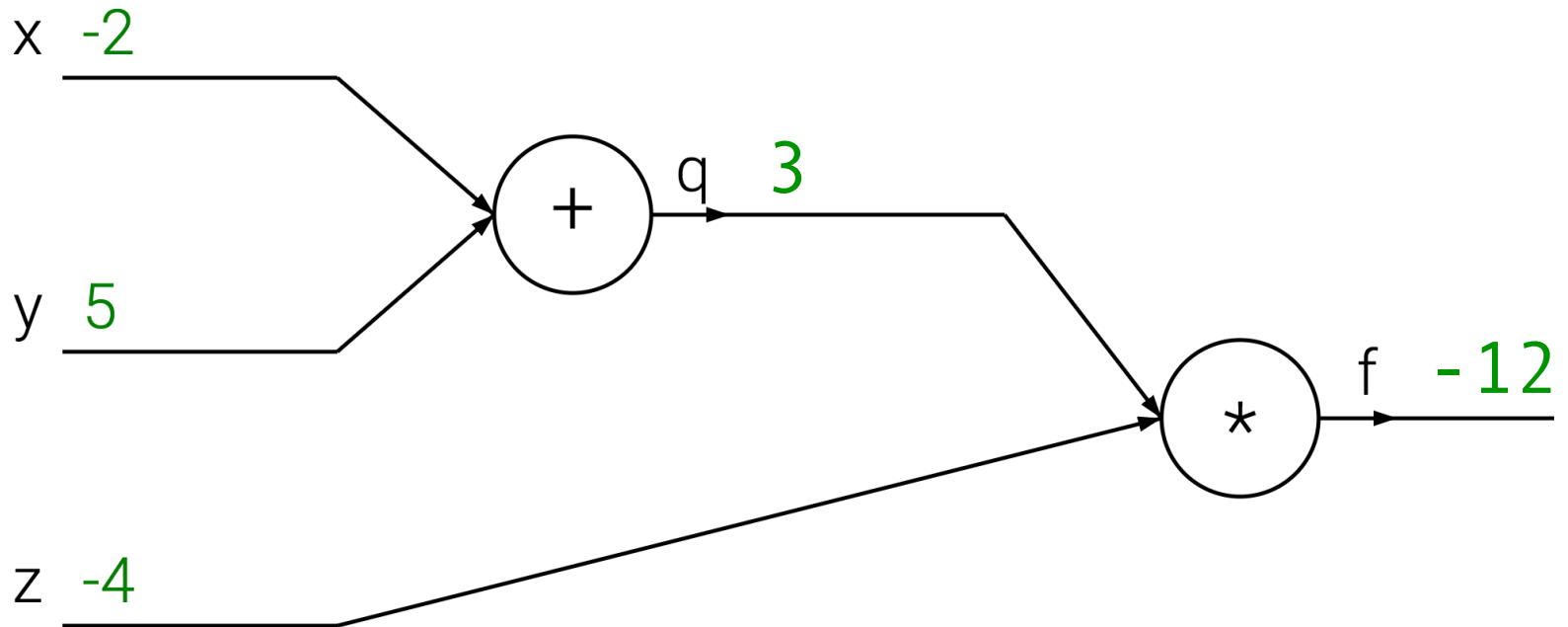


Example from:

<http://cs231n.github.io/optimization-2/>

Backpropagation: Example

Forward pass: compute values

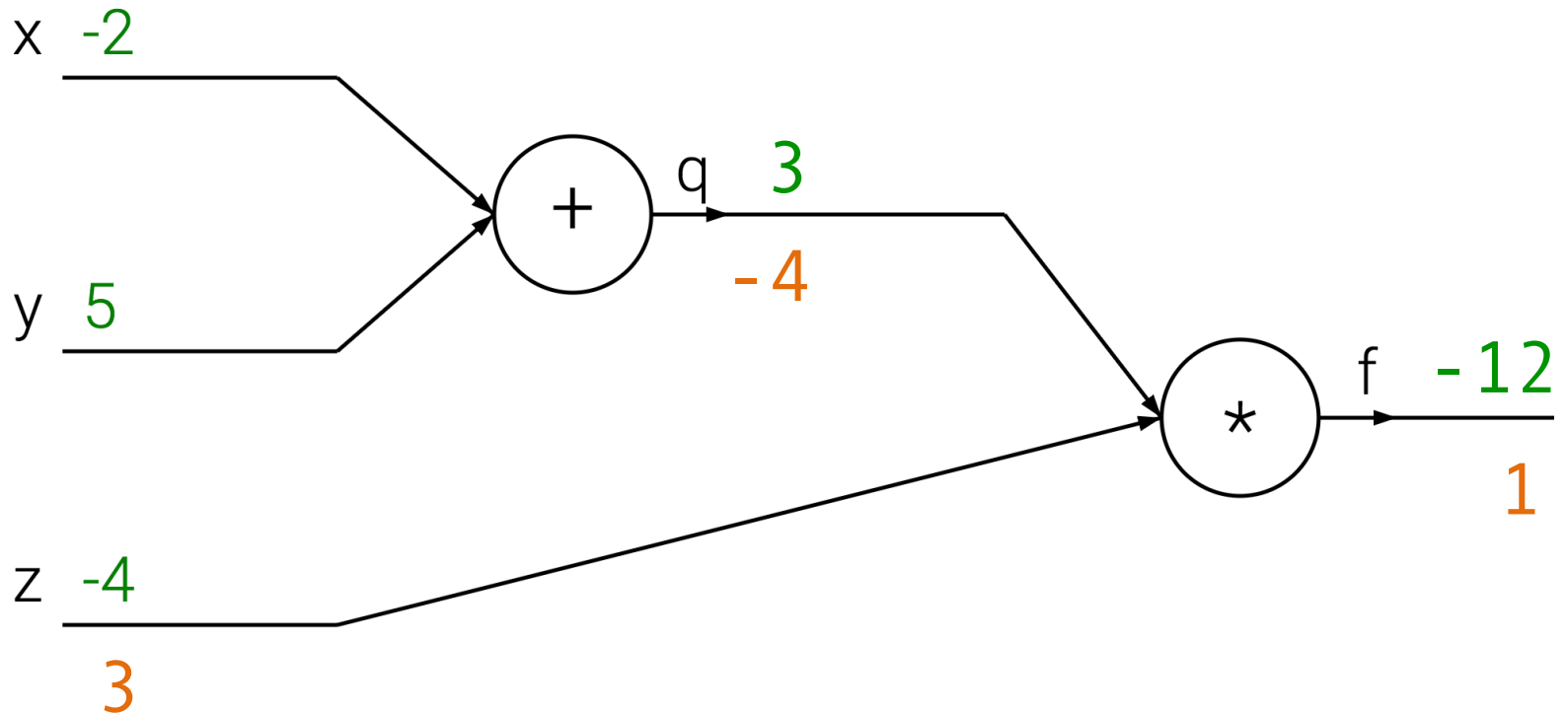


Example from:

<http://cs231n.github.io/optimization-2/>

Backpropagation: Example

Backward pass: compute local gradients

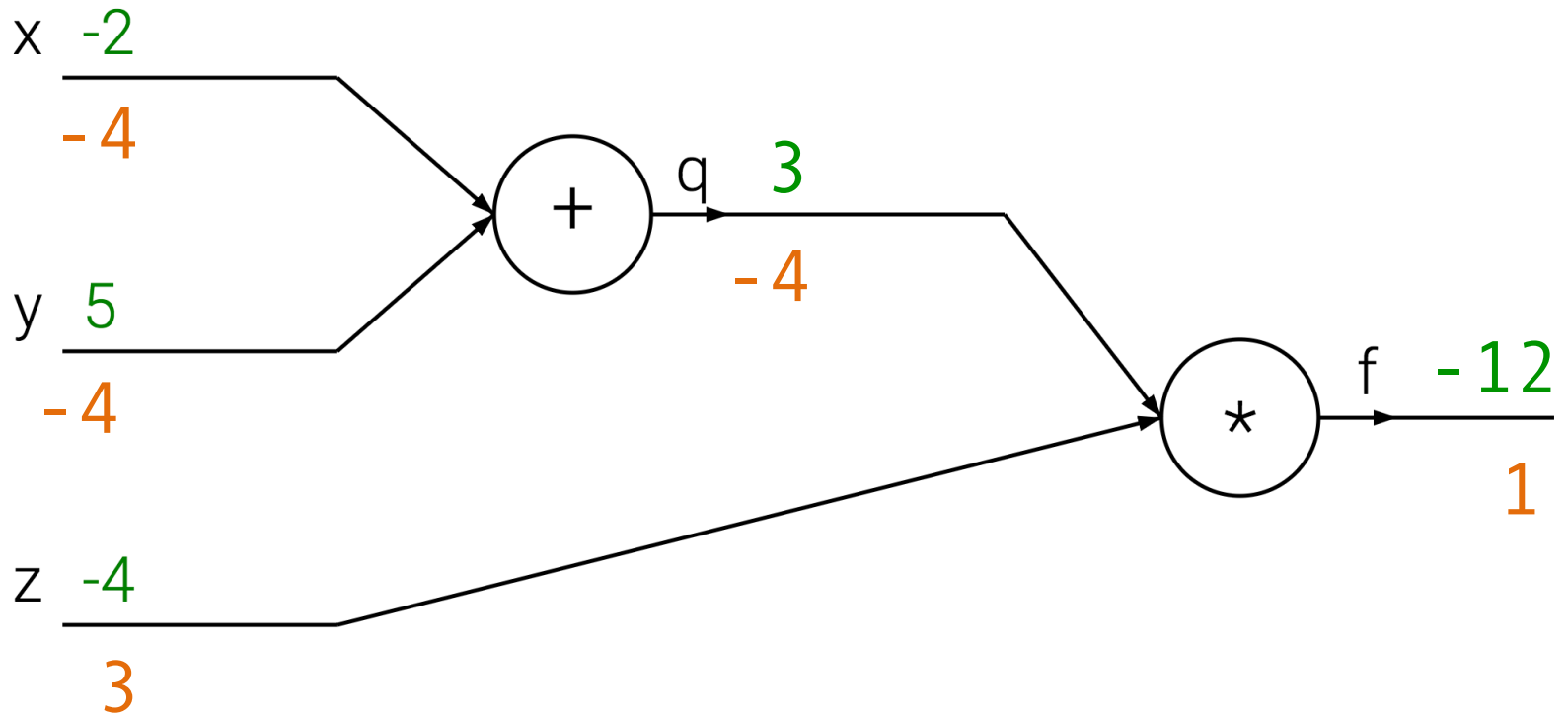


Example from:

<http://cs231n.github.io/optimization-2/>

Backpropagation: Example

Backward pass: compute local gradients

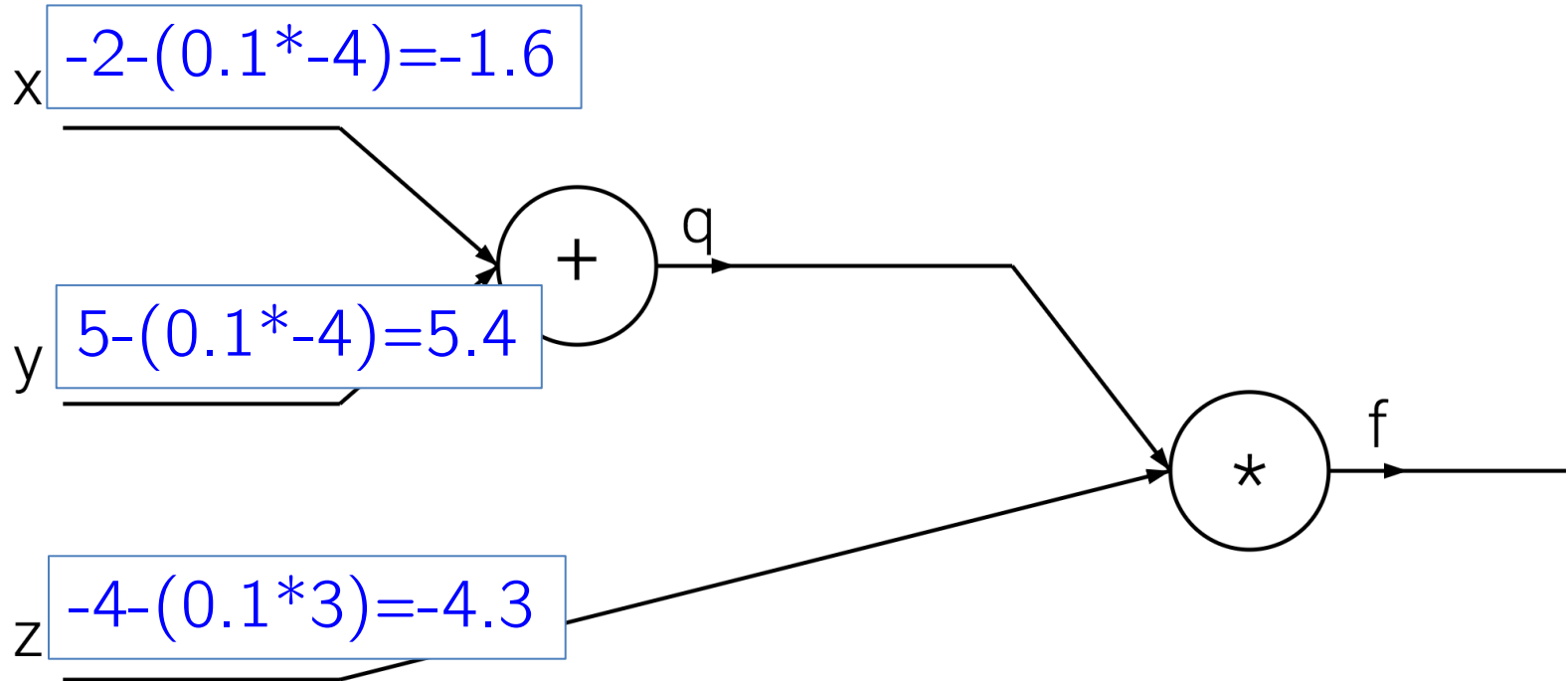


Example from:

<http://cs231n.github.io/optimization-2/>

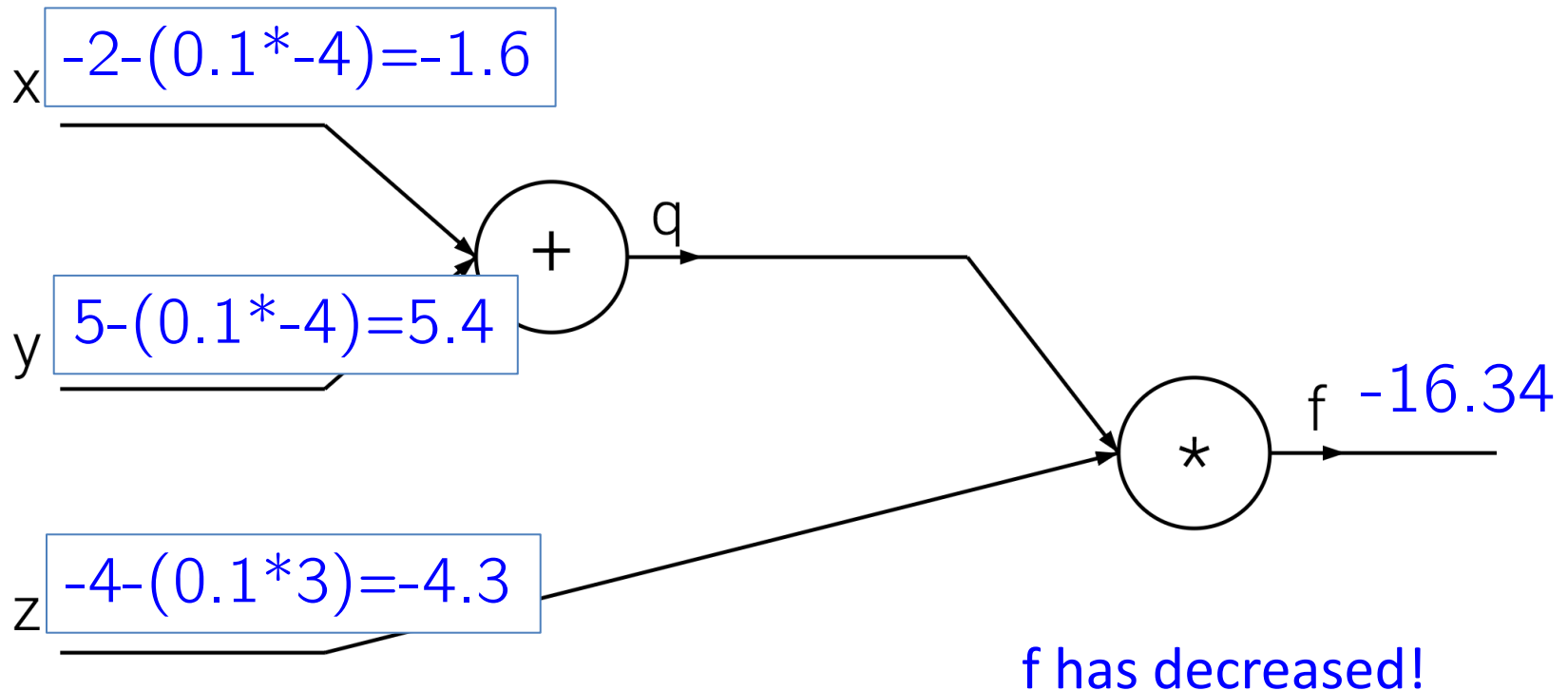
Backpropagation: Example

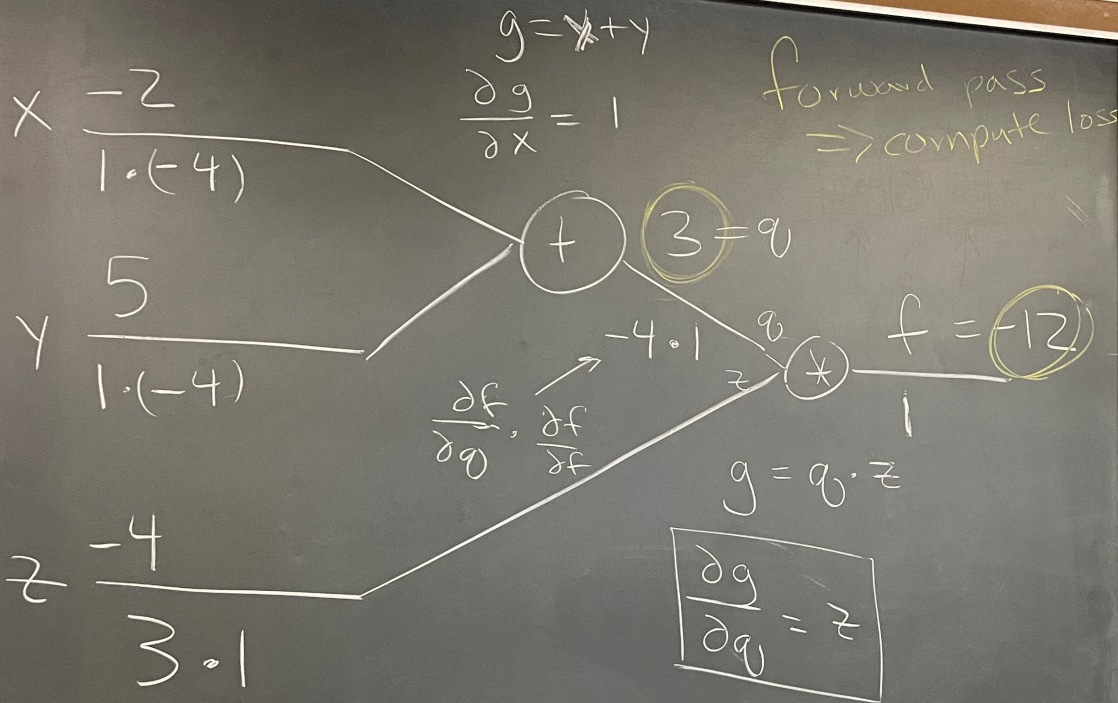
Now if we wanted to minimize $f \Rightarrow$ opposite direction of gradient



Backpropagation: Example

Now if we wanted to minimize $f \Rightarrow$ opposite direction of gradient





Backward pass
below lines

$$\frac{\partial g}{\partial z} = q$$

above lines

$$f(x, y, z) = \underbrace{(x+y)}_q z$$

$$f = qz$$

partial derivative of f with respect to x

$$\left\{ \frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \cdot \frac{\partial q}{\partial x} \right\} \text{chain rule}$$

$$= z \cdot 1$$

Minimize f

update trainable variables

$$x \leftarrow x - \eta \frac{\partial f}{\partial x}$$

$$-2 - 0.1(-4)$$

$$x \leftarrow -1.6$$

$$y \leftarrow y - \eta \frac{\partial f}{\partial y}$$

$$y \leftarrow 5 - 0.1(-4)$$

$$y \leftarrow 5.4$$

$$z \leftarrow -4 - 0.1(3)$$

$$z \leftarrow -4.3$$

$$f(x, y, z) = (-1.6 + 5.4)(-4.3)$$

$$= -16.34$$

lower than -12

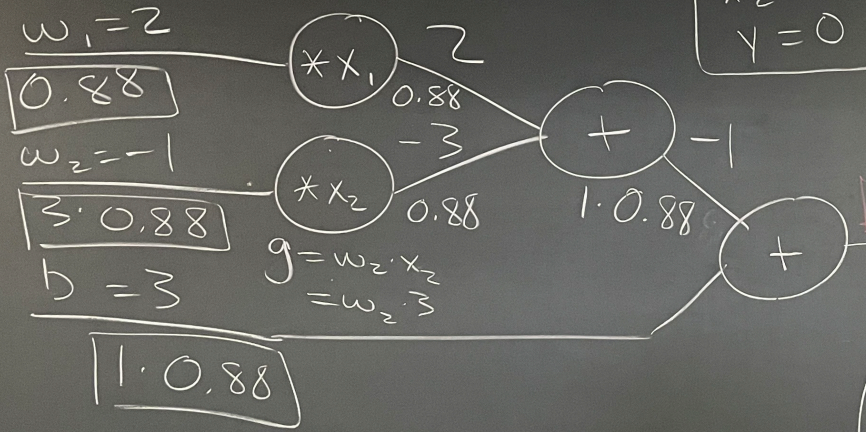


#4 and #5

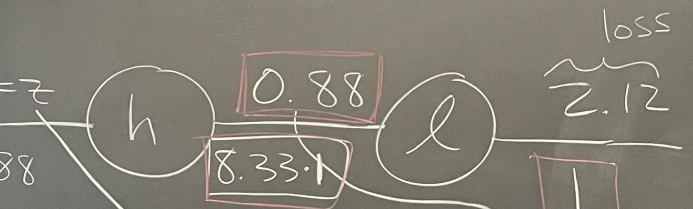
(#2 on your own)

③ logistic regression

$x_1 = 1$
 $x_2 = 3$
 $y = 0$ data



$$h(\vec{x}) = \frac{1}{1 + e^{-(w_1 x_1 + w_2 x_2 + b)}} = \frac{1}{1 + e^{-z}}$$



$$\frac{\partial h}{\partial z} = h(z)(1-h(z)) = h(z)(1-h(z))$$

$$\frac{\partial l}{\partial h} = \frac{1}{1-h(z)}$$

$$\frac{\partial l}{\partial h} = \frac{1}{1-0.88}$$

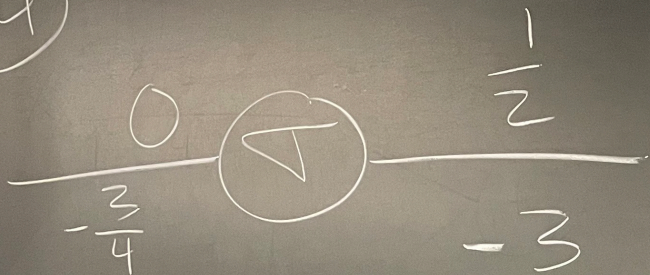
$$h(z)(1-h(z)) \cdot \frac{1}{1-h(z)}$$

$$\ell(h) = -\cancel{\gamma} \log h - (1-\cancel{\gamma}) \log(1-h)$$

$$\ell_0(h) = -\log(1-h)$$

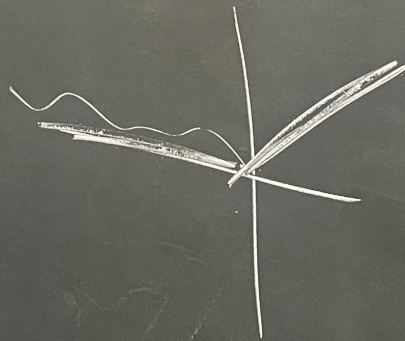
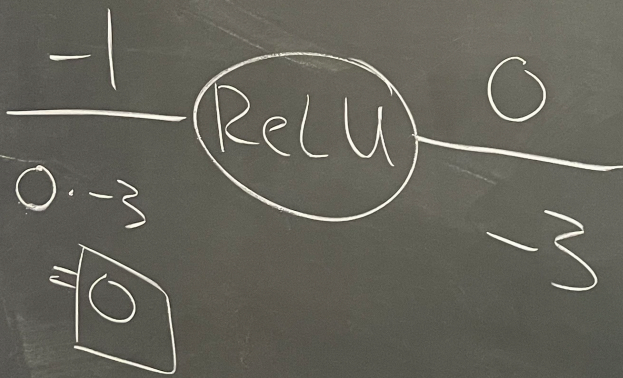
$$= \boxed{\frac{1}{1-h}}$$

④



$$\begin{aligned}\sigma'(z) &= \sigma(z)(1 - \sigma(z)) \\ &= \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}\end{aligned}$$

⑤



$$\sigma(0) =$$

$$\frac{1}{1 + e^{-0}} = \frac{1}{2}$$