#### **Deep Convolutional Neural Nets**

COMPSCI 371D — Machine Learning

・ 同 ト ・ ヨ ト ・ ヨ ト

## Outline

- 1 Why Neural Networks?
- 2 Circuits
- 3 Neurons, Layers, and Networks
- Orrelation and Convolution
- 6 AlexNet

< 🗇 🕨

→ Ξ →

## Why Neural Networks?

- Neural networks are very *expressive* (large  $\mathcal{H}$ )
- Can approximate any well-behaved function from a hypercube in ℝ<sup>d</sup> to an interval in ℝ within any ε > 0
- Universal approximators
- However
  - Complexity grows exponentially with d = dim(X)
  - *L<sub>T</sub>* is not convex (not even close)
  - Large  $\mathcal{H} \Rightarrow$  overfitting  $\Rightarrow$  *lots* of data!
- Amazon's Mechanical Turk made neural networks possible
- Even so, we cannot keep up with the curse of dimensionality!

## Why Neural Networks?

- Neural networks are data hungry
- Availability of lots of data is not a sufficient explanation
- There must be deeper reasons
- Special structure of image space (or audio space, or language)?
- Specialized network architectures?
- Regularization tricks and techniques?
- We don't really know. Stay tuned...
- Be prepared for some hand-waving and empirical statements

・ 同 ト ・ ヨ ト ・ ヨ ト

## Circuits

- Describe implementation of  $h : X \rightarrow Y$  on a computer
- Algorithm: A finite sequence of steps
- Circuit: Many gates of few types, wired together



- These are NAND gates. We'll use neurons
- Algorithms and circuits are equivalent
- Algorithm can simulate a circuit
- Computer is a circuit that runs algorithms!
- Computer really only computes Boolean functions...

・ 同 ト ・ ヨ ト ・ ヨ ト

#### Deep Neural Networks as Circuits

- Neural networks are typically described as circuits
- Nearly always implemented as algorithms
- One gate type, the neuron
- Many neurons that receive the same input form a layer
- A cascade of layers is a *network*
- A deep network has many layers
- Layers with a special constraint are called *convolutional*

#### The Neuron

- $y = \rho(a(\mathbf{x}))$  where  $a = \mathbf{v}^T \mathbf{x} + b$  $\mathbf{x} \in \mathbb{R}^d, \ y \in \mathbb{R}$
- **v** are the *gains*, *b* is the *bias*
- Together,  $\mathbf{w} = [\mathbf{v}, b]^T$  are the *weights*
- ρ(a) = max(0, a) (ReLU, Rectified Linear Unit)



38 N

# The Neuron as a Pattern Matcher (Almost)

• Left pattern is a drumbeat **g** (a pattern template):



- Which of the other two patterns **x** is a drumbeat?
- Normalize both  $\boldsymbol{g}$  and  $\boldsymbol{x}$  so that  $\|\boldsymbol{g}\|=\|\boldsymbol{x}\|=1$
- Then **g**<sup>*T*</sup>**x** is the cosine of the angle between the patterns
- If g<sup>T</sup>x ≥ −b for some threshold −b, output a = g<sup>T</sup>x + b (amount by which the cosine exceeds the threshold) otherwise, output 0
- $y = \rho(\mathbf{g}^T \mathbf{x} + b)$

## The Neuron as a Pattern Matcher (Almost)

- $y = \rho(\mathbf{g}^T \mathbf{x} + b)$
- A neuron is a pattern matcher, except for normalization and with a partial decision function
- In neural networks, normalization may happen in later or earlier layers
- This interpretation is not necessary to understand neural networks
- Nice to have a mental model, though
- Many neurons wired together can approximate any function we want: A neural network is a *universal function approximator*

-

・ロト ・ 一 ト ・ ヨ ト ・ ヨ ト

#### Layers and Networks

• A layer is a set of neurons that share the same input



- A neural network is a cascade of layers
- A neural network is *deep* if it has many layers
- *Two* layers can make a universal approximator
- If neurons did not have nonlinearities, any cascade of layers would collapse to a single layer

## **Convolutional Layers**

- A layer with input **x** ∈ ℝ<sup>d</sup> and output **y** ∈ ℝ<sup>e</sup> has *e* neurons, each with *d* gains and one bias
- Total of (d + 1)e weights to be trained in a single layer
- For images, *d*, *e* are in the order of hundreds of thousands or even millions
- Too many parameters
- Convolutional layers are layers restricted in a special way
- Many fewer parameters to train
- Also good justification in terms of basic principles

・ 同 ト ・ ヨ ト ・ ヨ ト

#### Hierarchy, Locality, Reuse

- To find a person, look for a face, a torso, limbs,...
- To find a face, look for eyes, nose, ears, mouth, hair,...
- To find an eye look for a circle, some corners, some curved edges,...
- A *hierarchical* image model is less sensitive to viewpoint, body configuration, ...
- Hierarchy leads to a *cascade* of layers
- Low-level features are *local*: A neuron doesn't need to see the entire image
- Circles are circles, regardless of where they show up: A single neuron can be *reused* to look for circles anywhere in the image

э

・ロッ ・ 一 ・ ・ ー ・ ・ ・ ・ ・ ・

# Correlation, Locality, and Reuse

• Does the drumbeat on the left show up in the clip on the right?



- Drumbeat **g** has 25 samples, clip **x** has 100
- Make 100 25 + 1 = 76 neurons that look for g in every possible position

• 
$$y_i = \rho(\mathbf{v}_i^T \mathbf{x} + b_i)$$
 where  $\mathbf{v}_i^T = [\underbrace{0, \dots, 0}_{i-1}, \underbrace{g_0, \dots, g_{24}}_{\mathbf{g}}, \underbrace{0, \dots, 0}_{76-i}]$   
• Layer gain matrix  $V = \begin{bmatrix} g_0 & \dots & g_{24} & 0 & 0 & \dots & 0 \\ 0 & g_0 & \dots & g_{24} & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \dots & 0 & g_0 & \dots & g_{24} \end{bmatrix}$ 

## **Compact Computation**

- Gain matrix  $V = \begin{bmatrix} g_0 & \cdots & g_{24} & 0 & 0 & \cdots & 0 \\ 0 & g_0 & \cdots & g_{24} & 0 & \cdots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & \cdots & \cdots & 0 & g_0 & \cdots & g_{24} \end{bmatrix}$
- Rows of Vx:  $z_i = \mathbf{v}_i^T \mathbf{x} = \sum_{a=0}^{24} g_a x_{i+a}$  for i = 0, ..., 75
- In general,

$$z_i = \sum_{a=0}^{k-1} g_a x_{i+a}$$
 for  $i = 0, \dots, e-1 = 0, \dots, d-k$ 

- (One-dimensional) correlation
- g is the kernel

#### A Small Example

$$z_i = \sum_{a=0}^2 g_a x_{i+a} \quad \text{for} \quad i = 0, \dots, 5$$





э

## **Correlation and Convolution**

- The correlation of **x** with  $\mathbf{g} = [g_0, \dots, g_{k-1}]$  is the *convolution* of **x** with  $\mathbf{r} = [r_0, \dots, r_{k-1}] = [g_{k-1}, \dots, g_0]$
- There are deep reasons why mathematicians prefer convolution to correlation
- We do not need to get into these, but see notes
- A layer whose gain matrix V is a correlation matrix is called a *convolutional layer*
- Also includes biases b

4 周 5 4 3 5 4 3 5 5

### Input Padding

- If the input has d entries and the kernel has k, then the output has e = d k + 1 entries
- This shrinkage is inconvenient when cascading several layers
- Pad input with *k* 1 zeros to make the output have *d* entries
- Padding is typically asymmetric when index is time, symmetric when index is position in space



Shape-preserving or 'same' correlation or convolution

э

#### **2D Correlation**

• Generalize in a straightforward way for 2D images:



(日)

#### Stride

- Images often vary slowly over space
- Output  $z_{ij}$  is often similar to  $z_{i,j+1}$  and  $z_{i+1,j}$
- Reduce the redundancy in the output by computing correlations with a *stride s*<sub>m</sub> greater than one
- Only compute every s<sub>m</sub> output values in dimension m ∈ {1,2}
- Output size shrinks from  $d_1 \times d_2$  to about  $d_1/s_1 \times d_2/s_2$

A (10) A (10)

# Max Pooling

- Another way to reduce output resolution is max pooling
- This is a layer of its own, separate from correlation
- Consider *k* × *k* windows with stride *s*
- Often s = k (adjacent, non-overlapping windows)
- · For each window, output the maximum value
- Output is about  $d_1/s \times d_2/s$
- Returns highest response in window, rather than the response in a fixed position

-

< 同 > < 回 > < 回 > -

#### The Input Layer of AlexNet

- AlexNet *circa* 2012, classifies color images into one of 1000 categories
- Trained on ImageNet, a large database with millions of labeled images



#### A more Compact Drawing



э

#### AlexNet



224x224x3

æ

## Output

• The last layer of a neural net used for classification is a soft-max layer

$$\mathbf{p} = \sigma(\mathbf{y}) = \frac{\exp(\mathbf{y})}{\mathbf{1}^T \exp(\mathbf{y})}$$

- The function from **x** to **p** is (nearly) differentiable
- Use cross-entropy loss on **p** to train
- After training, replace loss function with arg max p

#### AlexNet Numbers

- Input is  $224 \times 224 \times 3$  (color image)
- First layer has 96 feature maps of size  $55 \times 55$
- A fully-connected layer would have about  $224 \times 224 \times 3 \times 55 \times 55 \times 96 \approx 4.4 \times 10^{10}$  weights
- With convolutional kernels of size 11  $\times$  11, there are only 96  $\times$  11^2 = 11,616 weights
- That's a big deal! Locality and reuse
- Most of the complexity is in the last few, fully-connected layers, which still have millions of parameters
- More recent neural nets have much lighter final layers, but many more layers

-

< 同 > < 回 > < 回 > -