# CS 1674: Intro to Computer Vision Convolutional Neural Networks

Prof. Adriana Kovashka University of Pittsburgh November 1, 2018

## Plan for the next few lectures

#### Why (convolutional) neural networks?

#### Neural network basics

- Architecture
- Biological inspiration
- Loss functions
- Optimization / gradient descent
- Training with backpropagation

#### Convolutional neural networks (CNNs)

- Special operations
- Common architectures

#### **Practical matters**

- Tips and tricks for training
- Transfer learning
- Software packages

#### **Understanding CNNs**

- Visualization
- Synthesis / style transfer
- Breaking CNNs

# Neural network basics

#### Why (convolutional) neural networks?

#### State of the art performance on many problems Most papers in recent vision conferences use deep neural networks



## ImageNet Challenge 2012



[Deng et al. CVPR 2009]

- ~14 million labeled images, 20k classes
- Images gathered from Internet
- Human labels via Amazon Turk
- Challenge: 1.2 million training images, 1000 classes

A. Krizhevsky, I. Sutskever, and G. Hinton, <u>ImageNet Classification with Deep</u> <u>Convolutional Neural Networks</u>, NIPS 2012

Lana Lazebnik

# ImageNet Challenge 2012

- AlexNet: Similar framework to LeCun'98 but:
  - Bigger model (7 hidden layers, 650,000 units, 60,000,000 params)
  - More data (10<sup>6</sup> vs. 10<sup>3</sup> images)
  - GPU implementation (50x speedup over CPU)
    - Trained on two GPUs for a week
  - Better regularization for training (DropOut)



A. Krizhevsky, I. Sutskever, and G. Hinton, <u>ImageNet Classification with Deep</u> <u>Convolutional Neural Networks</u>, NIPS 2012

Adapted from Lana Lazebnik



#### ImageNet Challenge 2012

#### Krizhevsky et al. -- **16.4% error** (top-5) Next best (non-convnet) – **26.2% error**



### Example: CNN features for detection

#### **R-CNN:** Regions with CNN features



Object detection system overview. Our system (1) takes an input image, (2) extracts around 2000 bottom-up region proposals, (3) computes features for each proposal using a large convolutional neural network (CNN), and then (4) classifies each region using class-specific linear SVMs. R-CNN achieves a mean average precision (mAP) of 53.7% on PASCAL VOC 2010. For comparison, Uijlings et al. (2013) report 35.1% mAP using the same region proposals, but with a spatial pyramid and bag-of-visual-words approach. The popular deformable part models perform at 33.4%.

R. Girshick, J. Donahue, T. Darrell, and J. Malik, <u>Rich Feature Hierarchies for Accurate</u> <u>Object Detection and Semantic Segmentation</u>, CVPR 2014.

Lana Lazebnik

## What are CNNs?

- Convolutional neural networks are a type of neural network with layers that perform special operations
- Used in vision but also in NLP, biomedical etc.
- Often they are deep



## **Traditional Recognition Approach**



- Features are key to recent progress in recognition, but research shows they're flawed...
- Where next?

## What about learning the features?

- Learn a *feature hierarchy* all the way from pixels to classifier
- Each layer extracts features from the output of previous layer
- Train all layers jointly

#### "Shallow" vs. "deep" architectures

Traditional recognition: "Shallow" architecture



Deep learning: "Deep" architecture Image/ Video Pixels
Layer 1 ... Layer N Simple classifier Object Class

Lana Lazebnik

#### Neural network definition Figure 5.1 Network diagram for the twohidden units layer neural network corre $z_M$ sponding to (5.7). The input, $w_{MD}^{(1)}$ $w_{KM}^{(2)}$ hidden, and output variables are represented by nodes, and $x_D$ the weight parameters are rep $y_K$ resented by links between the nodes, in which the bias painputs outputs rameters are denoted by links coming from additional input $y_1$ and hidden variables $x_0$ and $z_0$ . Arrows denote the direc $x_1$ tion of information flow through the network during forward $w_{10}^{(2)}$ $z_1$ propagation. $x_0$ $z_0$ DRecall SVM: $a_j = \sum w_{ji}^{(1)} x_i + w_{j0}^{(1)}$ Activations: $w^T x + b$

hidden laver 1 hidden laver 2 hidden laver 3

• Nonlinear activation function h (e.g. sigmoid, RELU):  $z_i = h(a_i)$ 

i=1

Figure from Christopher Bishop



• Layer 3 (final)

 $a_k =$ 

Outputs (e.g. sigmoid/softmax)

(binary) 
$$y_k = \sigma(a_k) = \frac{1}{1 + \exp(-a_k)}$$
 (multiclass)  $y_k = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$   
• Finally:

(binary)  
$$y_k(\mathbf{x}, \mathbf{w}) = \sigma \left( \sum_{j=1}^M w_{kj}^{(2)} h \left( \sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right)$$

### Activation functions



### A multi-layer neural network



- Nonlinear classifier
- Can approximate any continuous function to arbitrary accuracy given sufficiently many hidden units

### Inspiration: Neuron cells

- Neurons
  - accept information from multiple inputs,
  - transmit information to other neurons.
- Multiply inputs by weights along edges
- Apply some function to the set of inputs at each node
- If output of function over threshold, neuron "fires"



#### **Biological analog**



### **Biological analog**



### Multilayer networks

- Cascade neurons together
- Output from one layer is the input to the next
- Each layer has its own sets of weights















#### Deep neural networks

- Lots of hidden layers
- Depth = power (usually)



#### How do we train them?

- The goal is to iteratively find such a set of weights that allow the activations/outputs to match the desired output
- We want to *minimize a loss function*
- The loss function is a function of the weights in the network
- For now let's simplify and assume there's a single layer of weights in the network

#### **Classification goal**

airplane	the second second	-	X	*	+	2			-
automobile				-	The state			-	*
bird		12	X		-	1		- Se	4
cat			50		1		Å,	1 and a start	
deer	1 20	X	R		Y	Ŷ	N.	-	
dog	17 A.	X		1		-	C?	1	The second
frog				2			S.		5
horse		A	2	P	H TAB	-3	No.	6	Y
ship	2	diri:	-	144	-	2	12	10-1	
truck							1		ALL.

Example dataset: CIFAR-10 10 labels 50,000 training images each image is 32x32x3 10,000 test images.

#### **Classification scores**

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



**10** numbers, indicating class scores

[32x32x3] array of numbers 0...1 (3072 numbers total)

#### Linear classifier



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



Andrej Karpathy

#### Linear classifier

#### Going forward: Loss function/Optimization



TODO:

- 1. Define a **loss function** that quantifies our unhappiness with the scores across the training data.
- 2. Come up with a way of efficiently finding the parameters that minimize the loss function.
  (optimization)

#### Linear classifier

Suppose: 3 training examples, 3 classes. With some W the scores f(x, W) = Wx are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

## Linear classifier: Hinge loss

Suppose: 3 training examples, 3 classes. With some W the scores f(x, W) = Wx are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

#### Hinge loss:

Given an example  $(x_i, y_i)$ where  $x_i$  is the image and where  $y_i$  is the (integer) label,

and using the shorthand for the scores vector:  $s = f(x_i, W)$ 

the loss has the form:

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

Want:  $s_{y_i} \ge s_j + 1$ i.e.  $s_j - s_{y_i} + 1 \le 0$ 

If true, loss is 0 If false, loss is magnitude of violation

### Linear classifier: Hinge loss

Suppose: 3 training examples, 3 classes. With some W the scores f(x, W) = Wx are:



#### Hinge loss:

Given an example  $(x_i, y_i)$ where  $x_i$  is the image and where  $y_i$  is the (integer) label,

and using the shorthand for the scores vector:  $s = f(x_i, W)$ 

the loss has the form:

 $\begin{aligned} L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\ &= \max(0, 5.1 - 3.2 + 1) \\ &+ \max(0, -1.7 - 3.2 + 1) \\ &= \max(0, 2.9) + \max(0, -3.9) \\ &= 2.9 + 0 \\ &= 2.9 \end{aligned}$
Suppose: 3 training examples, 3 classes. With some W the scores f(x, W) = Wx are:



#### Hinge loss:

Given an example  $(x_i, y_i)$ where  $x_i$  is the image and where  $y_i$  is the (integer) label,

and using the shorthand for the scores vector:  $s = f(x_i, W)$ 

the loss has the form:

 $\begin{aligned} L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\ &= \max(0, 1.3 - 4.9 + 1) \\ &+ \max(0, 2.0 - 4.9 + 1) \\ &= \max(0, -2.6) + \max(0, -1.9) \\ &= 0 + 0 \\ &= 0 \end{aligned}$ 

Suppose: 3 training examples, 3 classes. With some W the scores f(x, W) = Wx are:



#### Hinge loss:

Given an example  $(x_i, y_i)$ where  $x_i$  is the image and where  $y_i$  is the (integer) label,

and using the shorthand for the scores vector:  $s = f(x_i, W)$ 

the loss has the form:

$$\begin{split} L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\ &= \max(0, 2.2 - (-3.1) + 1) \\ &+ \max(0, 2.5 - (-3.1) + 1) \\ &= \max(0, 5.3 + 1) \\ &+ \max(0, 5.6 + 1) \\ &= 6.3 + 6.6 \\ &= 12.9 \end{split}$$

Suppose: 3 training examples, 3 classes. With some W the scores f(x, W) = Wx are:



cat	3.2	1.3	2.2	
car	5.1	4.9	2.5	
frog	-1.7	2.0	-3.1	
Losses:	2.9	0	12.9	-

#### Hinge loss:

Given an example  $(x_i, y_i)$ where  $x_i$  is the image and where  $y_i$  is the (integer) label,

and using the shorthand for the scores vector:  $s = f(x_i, W)$ 

the loss has the form:

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

and the full training loss is the mean over all examples in the training data:

$$L = rac{1}{N} \sum_{i=1}^N L_i$$

L = (2.9 + 0 + 12.9)/3 = 15.8 / 3 = **5.3** 

f(x,W) = Wx

 $L = rac{1}{N} \sum_{i=1}^{N} \sum_{j 
eq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1)$ 

#### Weight Regularization

 $\lambda$  = regularization strength (hyperparameter)

$$L=rac{1}{N}\sum_{i=1}^N\sum_{j
eq y_i} \max(0,f(x_i;W)_j-f(x_i;W)_{y_i}+1)+\lambda R(W)$$

#### In common use: L2 regularization L1 regularization Dropout (will see later)

$$egin{aligned} R(W) &= \sum_k \sum_l W_{k,l}^2 \ R(W) &= \sum_k \sum_l |W_{k,l}| \end{aligned}$$

# Another loss: Softmax (cross-entropy)



3.2

5.1

-1.7

scores = unnormalized log probabilities of the classes.

$$P(Y=k|X=x_i)=rac{e^{s_k}}{\sum_j e^{s_j}}$$
 where  $oldsymbol{s}=f(x_i;W)$ 

Want to maximize the log likelihood, or (for a loss function) to minimize the negative log likelihood of the correct class:

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

Andrej Karpathy

cat

car

frog

# Another loss: Softmax (cross-entropy)



# Other losses

• Triplet loss (Schroff, FaceNet)

$$\sum_{i=1}^{N} \left[ \|f(x_{i}^{a}) - f(x_{i}^{p})\|_{2}^{2} - \|f(x_{i}^{a}) - f(x_{i}^{n})\|_{2}^{2} + \alpha \right]_{+}$$

a denotes anchor p denotes positive n denotes negative



Figure 3. The **Triplet Loss** minimizes the distance between an *an-chor* and a *positive*, both of which have the same identity, and maximizes the distance between the *anchor* and a *negative* of a different identity.

Anything you want!

## How to minimize the loss function?



## How to minimize the loss function?

In 1-dimension, the derivative of a function:

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

In multiple dimensions, the gradient is the vector of (partial derivatives).

current W:	
[0.34, -1.11, 0.78, 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,]	
loss 1.25347	

gradient dW:



current W:	W + h (first dim):	gradient dW:
[0.34,	[0.34 <b>+ 0.0001</b> ,	[?,
-1.11,	-1.11,	?,
0.78,	0.78,	?.
0.12,	0.12,	?
0.55,	0.55,	?
2.81,	2.81,	?,
-3.1,	-3.1,	?
-1.5,	-1.5,	?
0.33,]	0.33,]	?]
loss 1.25347	loss 1.25322	

current W:	
[0.34,	
-1.11,	
0.78,	
0.12,	
0.55,	
2.81,	
-3.1,	
-1.5,	
0.33,]	
loss 1.25347	

W + h (first dim): [0.34 + **0.0001**, -1.11, 0.78, 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,...] loss 1.25322



current W:	<b>W</b> + h (
[0.34,	[0.34,
-1.11,	-1.11 +
0.78,	0.78,
0.12,	0.12,
0.55,	0.55,
2.81,	2.81,
-3.1,	-3.1,
-1.5,	-1.5,
0.33,]	0.33,]
loss 1.25347	loss 1.2

second dim): 0.0001, 25353

gradient dW:

[-2.5, ?, ?, ?, ?, ?, ?, ?, ?, ?,

current W:	W + h
[0.34,	[0.34,
-1.11,	-1.11
0.78,	0.78,
0.12,	0.12,
0.55,	0.55,
2.81,	2.81,
-3.1,	-3.1,
-1.5,	-1.5,
0.33,]	0.33,.
loss 1.25347	loss '

n (second dim): + 0.0001, . . 1.25353



current W:	W + h (third dim):
[0.34,	[0.34,
-1.11,	-1.11,
0.78,	0.78 + <b>0.0001</b> ,
0.12,	0.12,
0.55,	0.55,
2.81,	2.81,
-3.1,	-3.1,
-1.5,	-1.5,
0.33,]	0.33,]
loss 1.25347	loss 1.25347

gradient dW:

[-2.5, 0.6, ?, ?, ?, ?, ?, ?, ?, ?,

#### This is silly. The loss is just a function of W:

$$egin{aligned} L &= rac{1}{N} \sum_{i=1}^N L_i + \sum_k W_k^2 \ L_i &= \sum_{j 
eq y_i} \max(0, s_j - s_{y_i} + 1) \ s &= f(x; W) = Wx \end{aligned}$$

want  $\nabla_W L$ 

This is silly. The loss is just a function of W:  

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

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

$$s = f(x; W) = Wx$$
want  $\nabla_W L$ 
Calculus
$$\nabla_W L = ...$$

.

.

Andrej Karpathy

. .

- - 

current W:		gradient dW:
[0.34, -1.11, 0.78, 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,]	dW = (some function data and W)	[-2.5, 0.6, 0, 0.2, 0.7, -0.5, 1.1, 1.3, -2.1,]

# Loss gradients

- Denoted as (diff notations):  $\frac{\partial E}{\partial w_{ji}^{(1)}} = \nabla_W L$
- i.e. how does the loss change as a function of the weights
- We want to change the weights in such a way that makes the loss decrease as fast as possible <u>1</u>



# Gradient descent

- We'll update weights
- Move in direction opposite to gradient:

$$\mathbf{w}^{(\tau+1)}_{\uparrow} = \mathbf{w}^{(\tau)} - \eta \nabla E(\mathbf{w}^{(\tau)})$$
  
Time Learning rate



# Gradient descent

- Iteratively subtract the gradient with respect to the model parameters (w)
- I.e. we're moving in a direction opposite to the gradient of the loss
- I.e. we're moving towards *smaller* loss

# Mini-batch gradient descent

- In classic gradient descent, we compute the gradient from the loss for all training examples
- Could also only use some of the data for each gradient update
- We cycle through all the training examples multiple times
- Each time we've cycled through all of them once is called an 'epoch'
- Allows faster training (e.g. on GPUs), parallelization

### Learning rate selection



# Gradient descent in multi-layer nets

- We'll update weights
- Move in direction opposite to gradient:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E(\mathbf{w}^{(\tau)})$$

- How to update the weights at all layers?
- Answer: backpropagation of error from higher layers to lower layers

# Backpropagation: Graphic example

First calculate error of output units and use this to change the top layer of weights.



# Backpropagation: Graphic example

# Next calculate error for hidden units based on errors on the output units it feeds into.



# Backpropagation: Graphic example

Finally update bottom layer of weights based on errors calculated for hidden units.



$$f(x, y, z) = (x + y)z$$
  
e.g. x = -2, y = 5, z = -4



$$f(x, y, z) = (x + y)z$$
  
e.g. x = -2, y = 5, z = -4  
 $q = x + y$   $\frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$   
 $f = qz$   $\frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$ 



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

$$f(x, y, z) = (x + y)z$$
  
e.g.  $x = -2, y = 5, z = -4$   

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

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

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

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

$$f(x, y, z) = (x + y)z$$
  
e.g. x = -2, y = 5, z = -4  

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

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

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

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

$$f(x, y, z) = (x + y)z$$
  
e.g.  $x = -2, y = 5, z = -4$   

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

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

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

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

$$f(x, y, z) = (x + y)z$$
  
e.g.  $x = -2, y = 5, z = -4$   

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

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

$$x = -2$$
  

$$y = 5$$
  

$$z = -4$$
  

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

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

$$f(x, y, z) = (x + y)z$$
  
e.g. x = -2, y = 5, z = -4  

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

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

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

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

$$f(x, y, z) = (x + y)z$$
  
e.g.  $x = -2, y = 5, z = -4$   

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

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

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

Want: 
$$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$
$$f(x, y, z) = (x + y)z$$
  
e.g.  $x = -2, y = 5, z = -4$   

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

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

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

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

$$f(x, y, z) = (x + y)z$$
  
e.g.  $x = -2, y = 5, z = -4$   

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

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

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

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











# **Convolutional neural networks**

#### Convolutional Neural Networks (CNN)

- Neural network with specialized connectivity structure
- Stack multiple stages of feature extractors
- Higher stages compute more global, more invariant, *more abstract* features
- Classification layer at the end





Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, <u>Gradient-based learning applied to document</u> recognition, Proceedings of the IEEE 86(11): 2278–2324, 1998.

#### Convolutional Neural Networks (CNN)

- Feed-forward feature extraction:
  - 1. Convolve input with learned filters
  - 2. Apply non-linearity
  - 3. Spatial pooling (downsample)
- Supervised training of convolutional filters by back-propagating classification error



#### 1. Convolution

- Apply learned filter weights
- One feature map per filter
- Stride can be greater than
   1 (faster, less memory)





#### Feature Map

Adapted from Rob Fergus

#### 2. Non-Linearity

- Per-element (independent)
- Some options:
  - Tanh
  - Sigmoid: 1/(1+exp(-x))
  - Rectified linear unit (ReLU)
    - Avoids saturation issues



#### 3. Spatial Pooling

 Sum or max over non-overlapping / overlapping regions



## 3. Spatial Pooling

- Sum or max over non-overlapping / overlapping regions
- Role of pooling:
  - Invariance to small transformations
  - Larger receptive fields (neurons see more of input)







Sum





32x32x3 image



5x5x3 filter

**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"

#### **Convolution Layer**



#### **Convolution Layer**



activation map



**Convolution Layer** 

consider a second, green filter

32x32x3 image 5x5x3 filter convolve (slide) over all spatial locations

28

activation maps

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a "new image" of size 28x28x6!

**Preview:** ConvNet is a sequence of Convolution Layers, interspersed with activation functions



**Preview:** ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



#### **Preview**

[From recent Yann LeCun slides]



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]



A closer look at spatial dimensions:



A closer look at spatial dimensions:



A closer look at spatial dimensions:



A closer look at spatial dimensions:



A closer look at spatial dimensions:



A closer look at spatial dimensions:



7x7 input (spatially) assume 3x3 filter => 5x5 output

A closer look at spatial dimensions:



7x7 input (spatially) assume 3x3 filter applied **with stride 2** 

A closer look at spatial dimensions:



7x7 input (spatially) assume 3x3 filter applied **with stride 2** 

A closer look at spatial dimensions:



7x7 input (spatially) assume 3x3 filter applied with stride 2 => 3x3 output!

A closer look at spatial dimensions:



7x7 input (spatially) assume 3x3 filter applied **with stride 3?**
A closer look at spatial dimensions:



7x7 input (spatially) assume 3x3 filter applied **with stride 3?** 

doesn't fit! cannot apply 3x3 filter on 7x7 input with stride 3.

Andrej Karpathy

Ν



Output size: (N - F) / stride + 1

e.g. N = 7, F = 3: stride 1 => (7 - 3)/1 + 1 = 5stride 2 => (7 - 3)/2 + 1 = 3stride 3 => (7 - 3)/3 + 1 = 2.33 :\

#### In practice: Common to zero pad the border



e.g. input 7x7 **3x3** filter, applied with stride 1
pad with 1 pixel border => what is the output?

(recall:) (N - F) / stride + 1

#### In practice: Common to zero pad the border



e.g. input 7x7 **3x3** filter, applied with **stride 1 pad with 1 pixel** border => what is the output?

7x7 output!

#### In practice: Common to zero pad the border



e.g. input 7x7 **3x3** filter, applied with stride 1 **pad with 1 pixel** border => what is the output?

#### 7x7 output!

in general, common to see CONV layers with stride 1, filters of size FxF, and zero-padding with (F-1)/2. (will preserve size spatially)

e.g.  $F = 3 \Rightarrow zero pad$  with 1

- $F = 5 \Rightarrow zero pad with 2$
- F = 7 => zero pad with 3

(N + 2\*padding - F) / stride + 1

Examples time:

Input volume: **32x32x3** 10 5x5 filters with stride 1, pad 2

Output volume size: ?



Examples time:

Input volume: 32x32x3 10 5x5 filters with stride 1, pad 2



Output volume size: (32+2\*2-5)/1+1 = 32 spatially, so 32x32x10

Examples time:

Input volume: **32x32x3** 10 5x5 filters with stride 1, pad 2



Number of parameters in this layer?

Examples time:

Input volume: 32x32x3 10 5x5 filters with stride 1, pad 2



Number of parameters in this layer? each filter has 5\*5\*3 + 1 = 76 params (+1 for bias) => 76\*10 = 760

# Putting it all together



## A Common Architecture: AlexNet



# Case Study: VGGNet

		ConvNet C	onfiguration			
Α	A-LRN	B	С	D	E 19 weight layers	
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers		
	i	nput ( $224 \times 2$	24 RGB imag	:)		
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	
		max	pool			
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	
		max	pool			
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-250 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256	
	20 N	max	pool	0		
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512	
		max	pool			
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512	
		max	pool			
		FC-	4096			
		FC-	4096			
		FC-	1000			
		soft	-max			

[Simonyan and Zisserman, 2014]

Only 3x3 CONV stride 1, pad 1 and 2x2 MAX POOL stride 2

#### best model

11.2% top 5 error in ILSVRC 2013 -> 7.3% top 5 error

#### Table 2: Number of parameters (in millions).

Network	A,A-LRN	В	С	D	E
Number of parameters	133	133	134	138	144

# Case Study: GoogLeNet



Andrej Karpathy

## Case Study: ResNet

[He et al., 2015]

ILSVRC 2015 winner (3.6% top 5 error)



Slide from Kaiming He's recent presentation <u>https://www.youtube.com/watch?v=1PGLj-uKT1w</u>

Case Study: ResNet



(slide from Kaiming He's recent presentation)

## Case Study: ResNet





(slide from Kaiming He's recent presentation)

# **Practical matters**

# Comments on training algorithm

- Not guaranteed to converge to zero training error, may converge to local optima or oscillate indefinitely.
- However, in practice, does converge to low error for many large networks on real data, with good choice of hyperparameters (e.g. learning rate).
- Thousands of epochs (epoch = network sees all training data once) may be required, hours or days to train.
- To avoid local-minima problems, run several trials starting with different random weights (*random restarts*), and take results of trial with lowest training set error.
- May be hard to set learning rate and to select number of hidden units and layers.
- Neural networks had fallen out of fashion in 90s, early 2000s; back with a new name and improved performance (deep networks trained with dropout and lots of data).

# **Over-training prevention**

• Running too many epochs can result in over-fitting.



 Keep a hold-out validation set and test accuracy on it after every epoch. Stop training when additional epochs actually increase validation error.

#### **Training: Best practices**

- Use regularization
- To initialize weights, use "Xavier initialization"
- Use RELU or leaky RELU or ELU, don't use sigmoid
- Use data augmentation
- Use mini-batch
- Center (subtract mean from) your data
- Use cross-validation for your parameters
- Learning rate: too high? Too low?

## **Regularization: Dropout**



- Randomly turn off some neurons
- Allows individual neurons to independently be responsible for performance

Dropout: A simple way to prevent neural networks from overfitting [Srivastava JMLR 2014]

# Data Augmentation (Jittering)

#### Create virtual training samples

- Horizontal flip
- Random crop
- Color casting
- Geometric distortion





Andrej Karpathy

#### **Transfer Learning with CNNs**

- The more weights you need to learn, the more data you need
- That's why with a deeper network, you need more data for training than for a shallower network
- One possible solution:



Set these to the already learned Learn these on your own task weights from another network

# Transfer Learning with CNNs



Another option: use network as feature extractor, train SVM on extracted features for target task

# Transfer Learning with CNNs



# Pre-training on ImageNet

- Have a source domain and target domain
- Train a network to classify ImageNet classes
  - Coarse classes and ones with fine distinctions (dog breeds)
- Remove last layers and train layers to replace them, that predict target classes



Oquab et al., "Learning and Transferring Mid-Level Image Representations...", CVPR 2014

#### Transfer learning with CNNs is pervasive...



 $y_t$ 

 $h_t$ 

 $x_t$ 

 $W_{oh}$ 

 $W_{hx}$ 

**Object Detection** Ren et al., "Faster R-CNN", NIPS 2015

Adapted from Andrej Karpathy

#### Semantic segmentation



Andrej Karpathy

#### Photographer identification

#### Who took this photograph?



- Deep net features achieve 74% accuracy
  - Chance is less than 3%, human performance is 47%
- Method learns what proto-objects + scenes authors shoot

# Analysis of pre-training on ImageNet

- Source:
  - distinguish 1000 ImageNet categories (incl. many dog breeds)
- Target tasks:
  - object detection and action recognition on PASCAL
  - scene recognition on SUN
- Pre-training with 500 images per class is about as good as pre-training with 1000
- Pre-training with 127 classes is about as good as pre-training with 1000
- Pre-training with (fewer classes, more images per class) > (more classes, fewer images)
- Small drop in if classes with fine-grained distinctions removed from pre-training set

Huh et al., "What makes ImageNet good for transfer learning?", arxiv 2016

#### Packages

**TensorFlow** 

Torch / PyTorch

**Keras** 

Caffe and Caffe Model Zoo

#### Some Learning Resources

http://deeplearning.net/

http://cs231n.stanford.edu

# **Understanding CNNs**

#### **Recall: Biological analog**



Layer 1





Visualizing and Understanding Convolutional Networks [Zeiler and Fergus, ECCV 2014]
#### Layer 2



 Activations projected down to pixel level via decovolution  Patches from validation images that give maximal activation of a given feature map

Visualizing and Understanding Convolutional Networks [Zeiler and Fergus, ECCV 2014]

### Layer 3



Visualizing and Understanding Convolutional Networks [Zeiler and Fergus, ECCV 2014]

### Layer 4 and 5



Visualizing and Understanding Convolutional Networks [Zeiler and Fergus, ECCV 2014]

# **Occlusion experiments**



(d) Classifier, probability of correct class

> (as a function of the position of the square of zeros in the original image)

[Zeiler & Fergus 2014]

## **Occlusion experiments**



(as a function of the position of the square of zeros in the original image)

[Zeiler & Fergus 2014]

# What image maximizes a class score?



#### Repeat:

- 1. Forward an image
- 2. Set activations in layer of interest to all zero, except for a 1.0 for a neuron of interest
- 3. Backprop to image
- 4. Do an "image update"

# What image maximizes a class score?



[Understanding Neural Networks Through Deep Visualization, Yosinski et al., 2015] http://yosinski.com/deepvis

### What image maximizes a class score?



# **Breaking CNNs**



Take a correctly classified image (left image in both columns), and add a tiny distortion (middle) to fool the ConvNet with the resulting image (right).

#### Intriguing properties of neural networks [Szegedy ICLR 2014]

Andrej Karpathy

# **Breaking CNNs**



Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images [Nguyen et al. CVPR 2015]

Jia-bin Huang

# Fooling a linear classifier



Fooled linear classifier: The starting image (left) is classified as a kit fox. That's incorrect, but then what can you expect from a linear classifier? However, if we add a small amount 'goldfish' weights to the image (top row, middle), suddenly the classifier is convinced that it's looking at one with high confidence. We can distort it with the school bus template instead if we wanted to.

To fool a linear classifier, add a small multiple of the weight vector to the training example:

 $\mathbf{x} \leftarrow \mathbf{x} + \alpha \mathbf{w}$ 

#### http://karpathy.github.io/2015/03/30/breaking-convnets/

# Summary

- We use deep neural networks because of their strong performance in practice
- Convolutional neural network (CNN)
  - Convolution, nonlinearity, max pooling
- Training deep neural nets
  - We need an objective function that measures and guides us towards good performance
  - We need a way to minimize the loss function: stochastic gradient descent
  - We need backpropagation to propagate error towards all layers and change weights at those layers
- Practices for preventing overfitting
  - Dropout; data augmentation; transfer learning