# Lecture 19: Recognition with CNNs

COMP 590/776: Computer Vision Instructor: Soumyadip (Roni) Sengupta TA: Mykhailo (Misha) Shvets



Course Website: Scan Me!

# Image Classification: a core task in computer vision

• Assume given set of discrete labels, e.g.

{cat, dog, cow, apple, tomato, truck, ... }



#### Linear Classifier

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



# Linear Classifier

 Parameters define a hyperplane for each class:

$$f(x_i, W, b) = Wx_i + b$$

 We can think of each class score as defining a distribution that is proportional to distance from the corresponding hyperplane



# Simpler example: binary classification

- Two classes (e.g., "cat" and "not cat")
  - AKA "positive" and "negative" classes









not cat

## Simpler example: binary classification

- Find linear function (*hyperplane*) to separate positive and negative examples
  - $\mathbf{x}_i$  positive:  $\mathbf{x}_i \cdot \mathbf{w} + b \ge 0$  $\mathbf{x}_i$  negative:  $\mathbf{x}_i \cdot \mathbf{w} + b < 0$

Which hyperplane is best? We need a **loss function** to decide



#### Linear classification



airplane	-3.45	-0.51	3.42
automobile	-8.87	6.04	4.64
bird	0.09	5.31	2.65
cat	2.9	-4.22	5.1
deer	4.48	-4.19	2.64
dog	8.02	3.58	5.55
frog	3.78	4.49	-4.34
horse	1.06	-4.37	-1.5
ship	-0.36	-2.09	-4.79
truck	-0.72	-2.93	6.14

Cat image by Nikita is licensed under CC-BY 2.0; Car image is CC0 1.0 public domain; Frog image is in the public domain

#### Output scores

#### TODO:

- 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)

## Loss functions

3.2

5.1

-1.7

cat

car

frog

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



1.3

4.9

2.0

2.2

2.5

-3.1

A **loss function** tells how good our current classifier is

Given a dataset of examples  $\{(x_i, y_i)\}_{i=1}^N$ 

Where  $oldsymbol{x_i}$  is image and  $oldsymbol{y_i}$  is (integer) label

Loss over the dataset is a sum of loss over examples:

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

## Softmax classifier

 Interpret Scores as unnormalized log probabilities of classes

$$f(x_i, W) = Wx_i$$
 (score function)



softmax function

Squashes values into *probabilities* ranging from 0 to 1

$$P(y_i \mid x_i; W)$$

Example with three classes:

 $[1,-2,0] \rightarrow [e^1,e^{-2},e^0] = [2.71,0.14,1] \rightarrow [0.7,0.04,0.26]$ 

#### Softmax classifier

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



#### Cross-entropy loss

 $f(x_i, W) = Wx_i$  (score function)

#### Cross-entropy loss

 $f(x_i, W) = W x_i$  (score function)



#### Cross-entropy loss

 $f(x_i, W) = W x_i$  (score function)



# Summary

- Have score function and loss function
  - Currently, score function is based on linear classifier
  - Next, will generalize to convolutional neural networks
- Find W and b to minimize loss



#### (Deep) Neural Networks

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

(**Before**) Linear score function:

(Now) 2-layer Neural Network

$$egin{aligned} f &= Wx \ f &= W_2 \max(0, W_1 x) \ \end{array}$$

Non-linear Activation Function (many other choices exist)



(Before) Linear score function: f = Wx(Now) 2-layer Neural Network  $f = W_2 \max(0, W_1 x)$  $x W_1 h W_2 s$ 3072 10

Total number of weights to learn:
 3,072 x 100 + 100 x 10 = 308,200

(Before) Linear score function: f = Wx(Now) 2-layer Neural Network  $f = W_2 \max(0, W_1 x)$ or 3-layer Neural Network  $f = W_3 \max(0, W_2 \max(0, W_1 x))$ 

> also called "Multi-Layer Perceptrons" (MLPs)

- Very coarse generalization of neural networks:
  - Linear functions chained together and separated by non-linearities (*activation functions*), e.g. "max"

$$f=W_3\max(0,W_2\max(0,W_1x))$$

- Why separate linear functions with non-linear functions?
- *Very roughly* inspired by real neurons



#### Activation functions



# Leaky ReLU $\max(0.1x, x)$



Maxout  $\max(w_1^T x + b_1, w_2^T x + b_2)$ 



#### Neural network architecture

• Computation graph for a 2-layer neural network



#### Neural networks: Architectures

![](_page_23_Figure_1.jpeg)

**Deep** networks typically have many layers and potentially millions of parameters

#### Example feed-forward computation of a neural network

![](_page_24_Picture_1.jpeg)

# 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)

# Optimizing parameters with gradient descent

- How do we find the best **W** and **b** parameters?
- In general: *gradient descent* 
  - 1. Start with a guess of a good **W** and **b** (or randomly initialize them)
  - 2. Compute the loss function for this initial guess and the *gradient* of the loss function
  - 3. Step some distance in the negative gradient direction (direction of steepest descent)
  - 4. Repeat steps 2 & 3
- Note: efficiently performing step 2 for deep networks is called *backpropagation*

#### **The Learning Cycle: Forward Propagation**

![](_page_26_Figure_1.jpeg)

![](_page_26_Picture_2.jpeg)

#### **The Learning Cycle: Backward Propagation**

![](_page_27_Picture_1.jpeg)

![](_page_27_Picture_2.jpeg)

How do we change our weights?

![](_page_28_Figure_0.jpeg)

negative gradient direction

Gradient descent: walk in the direction opposite gradient

- **Q**: How far?
- A: Step size: *learning rate*
- Too big: will miss the minimum
- Too small: slow convergence

# 2D example of gradient descent

![](_page_29_Figure_1.jpeg)

- In reality, in deep learning we are optimizing a highly complex loss function with millions of variables (or more)
- More on this later...

#### Convolutional Neural Networks (CNNs)

#### Convolutional neural networks

A bit of history: Gradient-based learning applied to document recognition [LeCun, Bottou, Bengio, Haffner 1998]

![](_page_31_Figure_2.jpeg)

## Fast-forward to today: ConvNets are everywhere

#### Classification

Retrieval

![](_page_32_Picture_3.jpeg)

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

#### Image features vs ConvNets

![](_page_33_Figure_1.jpeg)

training

## Last layer of most CNNs is a linear classifier

![](_page_34_Figure_1.jpeg)

Input Perform everything with a big neural Pixels network, trained end-to-end

**Key:** perform enough processing so that by the time you get to the end of the network, the classes are linearly separable

#### Visualizing AlexNet in 2D with t-SNE

![](_page_35_Figure_1.jpeg)

(2D visualization using t-SNE)

[Donahue, "DeCAF: DeCAF: A Deep Convolutional ...", arXiv 2013]
# Convolutional neural networks

- Layer types:
  - Convolutional layer
  - Pooling layer
  - Fully-connected layer

#### **AlexNet: An Early Example**





https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf

3



32x32x3 image



5x5x3 filter

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

# Convolution (Recap)

• Same as cross-correlation, except that the kernel is "flipped" (horizontally and vertically)

$$\begin{split} G[i,j] &= \sum_{u=-k}^{k} \sum_{v=-k}^{k} H[u,v] F[i-u,j-v] \\ G[i,j] &= \sum_{u=-k}^{k} \sum_{v=-k}^{k} H[u,v] F[i+u,j+v] \quad \text{Cross-correlation} \end{split}$$

• Convolution is **commutative** and **associative** 



32x32x3 image 5x5x3 filter 32 **Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products" 32 Number of weights: 5 x 5 x 3 + 1 = **76** (vs. 3072 for a fully-connected layer) (+1 for bias)





activation map



consider a second, green filter

activation maps

28

28

#### Convolution Layer



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!

(total number of parameters:  $6 \times (75 + 1) = 456$ )

Activation functions



Leaky ReLU  $\max(0.1x, x)$ 



 $\begin{array}{l} \textbf{Maxout} \\ \max(w_1^T x + b_1, w_2^T x + b_2) \end{array}$ 



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



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



#### **AlexNet: An Early Example**





https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf

3

## **Convolutional Networks**

Learnable 3x3 Convolutional Kernels



Slide Credits: Gedas Bertasius

### How Else to Shrink the Model Size?

Pooling Layer:

- Max Pooling
- Other pooling options like average pooling are also used







## **Convolutional Networks**

Learnable 3x3 Convolutional Kernels



Slide Credits: Gedas Bertasius



AlexNet: An Early Example

https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf

## Fully Connected Layer



## **Convolutional Networks**

Learnable FC Layer Weight Matrix



 $W \in \mathbb{R}^{d \times C}$ 

- d feature dimensionality (4800 in this example)
- ${\cal C}$  number of classes



Slide Credits: Gedas Bertasius

## **Convolutional Networks**

Learnable FC Layer Weight Matrix



 $W \in \mathbb{R}^{d \times C}$ 

- d feature dimensionality (4800 in this example)
- ${\cal C}$  number of classes



#### Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



### **Fully Connected Layer**

32x32x3 image -> stretch to 3072 x 1



#### Same as a linear classifer!

#### Where Models Learn Features of an Image







# AlexNet (2012)

#### 6M parameters in total



Elgendy, Deep Learning for Vision Systems, https://livebook.manning.com/book/grokking-deep-learning-for-computer-vision/chapter-5/v-3/







[He et al. CVPR 2016]

[Krizhevsky et al. NIPS 2012]

2] [Szegedy et al. CVPR 2015]

[Simonyan & Zisserman, ICLR 2015]

# Big picture

- A convolutional neural network can be thought of as a function from images to class scores
  - With millions of adjustable weights...
  - ... leading to a very non-linear mapping from images to features / class scores.
  - We will set these weights based on classification accuracy on training data...
  - ... and hopefully our network will generalize to new images at test time

## Data is key—enter ImageNet

- ImageNet (and the ImageNet Large-Scale Visual Recognition Challege, aka **ILSVRC**) has been key to training deep learning methods
  - J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li and L. Fei-Fei, ImageNet: A Large-Scale Hierarchical Image Database. *CVPR*, 2009.
- **ILSVRC**: 1,000 object categories, each with ~700-1300 training images. Test set has 100 images per categories (100,000 total).
- Standard ILSVRC error metric: top-5 error
  - if the correct answer for a given test image is in the top 5 categories, your answer is judged to be correct

# Performance improvements on ILSVRC

- ImageNet Large-Scale Visual Recognition Challenge
- Held from 2011-2017
- 1000 categories, 1000 training images per category
- Test performance on held-out test set of images





Image credit: Zaid Alyafeai, Lahouari Ghouti

## Closer look at Convolution

A closer look at spatial dimensions:



A closer look at spatial dimensions:



7x7 input (spatially) assume 3x3 filter

A closer look at spatial dimensions:



7x7 input (spatially) assume 3x3 filter


7x7 input (spatially) assume 3x3 filter



7x7 input (spatially) assume 3x3 filter



7x7 input (spatially) assume 3x3 filter

=> 5x5 output



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



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



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



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



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

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



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

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

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

0	0	0	0	0	0		
0							
0							
0							
0							

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!

Output filter size: (N + 2\*pad - F)/S + 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 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 => zero pad with 1

F = 5 => zero pad with 2

F = 7 = 2 zero pad with 3

# Padding & Stride in CNN

CLASS torch.nn.Conv2d(*in\_channels*, *out\_channels*, *kernel\_size*, *stride=1*, *padding=0*, *dilation=1*, *groups=1*, *bias=True*, *padding\_mode='zeros'*, *device=None*, *dtype=None*) [SOURCE]

Applies a 2D convolution over an input signal composed of several input planes.

In the simplest case, the output value of the layer with input size  $(N, C_{\rm in}, H, W)$  and output  $(N, C_{\rm out}, H_{\rm out}, W_{\rm out})$  can be precisely described as:

$$\operatorname{out}(N_i, C_{\operatorname{out}_j}) = \operatorname{bias}(C_{\operatorname{out}_j}) + \sum_{k=0}^{C_{\operatorname{in}}-1} \operatorname{weight}(C_{\operatorname{out}_j}, k) \star \operatorname{input}(N_i, k)$$

where  $\star$  is the valid 2D cross-correlation operator, N is a batch size, C denotes a number of channels, H is a height of input planes in pixels, and W is width in pixels.

$$egin{aligned} H_{out} &= igg igg [rac{H_{in}+2 imes ext{padding}[0]- ext{dilation}[0] imes ( ext{kernel\_size}[0]-1)-1}{ ext{stride}[0]}+1 igg ] \ W_{out} &= igg [rac{W_{in}+2 imes ext{padding}[1]- ext{dilation}[1] imes ( ext{kernel\_size}[1]-1)-1}{ ext{stride}[1]}+1 igg ] \end{aligned}$$



padding=1, stride=2

#### Remember back to...

E.g. 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially! (32 -> 28 -> 24 ...). Shrinking too fast is not good, doesn't work well.



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

Output volume size: ?



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



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

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



Number of parameters in this layer?

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 (btw, 1x1 convolution layers make perfect sense)



# Convolutional layer—properties

- Small number of parameters to learn compared to a fully connected layer
- Preserves spatial structure—output of a convolutional layer is shaped like an image
- **Translation equivariant**: passing a translated image through a convolutional layer is (almost) equivalent to translating the convolution output (but be careful of image boundaries)

# Self-study

# [ConvNetJS demo: training on CIFAR-10]

#### ConvNetJS CIFAR-10 demo

#### Description

This demo trains a Convolutional Neural Network on the <u>CIFAR-10 dataset</u> in your browser, with nothing but Javascript. The state of the art on this dataset is about 90% accuracy and human performance is at about 94% (not perfect as the dataset can be a bit ambiguous). I used <u>this python script</u> to parse the <u>original files</u> (python version) into batches of images that can be easily loaded into page DOM with img tags.

This dataset is more difficult and it takes longer to train a network. Data augmentation includes random flipping and random image shifts by up to 2px horizontally and verically.

By default, in this demo we're using Adadelta which is one of per-parameter adaptive step size methods, so we don't have to worry about changing learning rates or momentum over time. However, I still included the text fields for changing these if you'd like to play around with SGD+Momentum trainer.

Report questions/bugs/suggestions to @karpathy.



https://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html

# Train a neural network for classification on CIFAR10 dataset in google colab

**Google Colab Page** 

# Where to Look for More Information

- Explore existing computer vision and machine learning frameworks
  - o <u>https://pytorch.org/</u>
  - o <u>https://www.tensorflow.org/</u>
  - o <u>https://keras.io/</u>
  - o <u>https://opencv.org/</u>
- Watch more in-depth lecture series
  - The Ancient Secrets of Computer Vision Joseph Redmon
  - O Deep Learning Specialization Andrew Ng
- Checkout other online courses and guides
  - o <u>https://ai.google/education/</u>
  - o <u>https://www.udacity.com/course/deep-learning-pytorch--ud188</u>



# Slide Credits

- <u>CS5670, Introduction to Computer Vision</u>, Cornell Tech, by Noah Snavely.
- CS 543 Computer Vision, by Stevlana Lazebnik, UIUC.
- EECS 442 <u>Computer Vision</u>, by Justin Johnson & David Fouhey, U Michigan.