

CNN Architectures

I2DL: Prof. Dai



Convolution Layer



Convolution Layer



4

Convolution Layers: Dimensions

Input width of N

				ght		
Ζ				hei		
nt of				llter F F	-	
eigh	Filte	r wid	th	.0 19)	
ut he	of F					
Inpu						

Input: $N \times N$ Filter: $F \times F$ Stride: SOutput: $(\frac{N-F}{S} + 1) \times (\frac{N-F}{S} + 1)$

$$N = 7, F = 3, S = 1: \quad \frac{7-3}{1} + 1 = 5$$

$$N = 7, F = 3, S = 2: \quad \frac{7-3}{2} + 1 = 3$$

$$N = 7, F = 3, S = 3: \quad \frac{7-3}{3} + 1 = 2.3333$$

Fractions are illegal

Convolution Layers: Padding

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Types of convolutions:

• Valid convolution: using no padding

• Same convolution: output=input size

Set padding to $P = \frac{F-1}{2}$

mage 7x7 + zero padding

CNN Learned Filters



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013] L: Prof. Dai

CNN Prototype





Classic Architectures





Input: 32×32 grayscale images This one: Labeled as class "7"





- Valid convolution: size shrinks
- How many conv filters are there in the first layer?





• At that time average pooling was used, now max pooling is much more common





• Again valid convolutions, how many filters?





• Use of tanh/sigmoid activations \rightarrow not common now!





• Conv -> Pool -> Conv -> Pool -> Conv -> FC

60k parameters



- Conv -> Pool -> Conv -> Pool -> Conv -> FC
- As we go deeper: Width, Height Vumber of Filters

Deep Neural Clustering



(a) k-Means

(b) Autoencoder + k-Means

(c) Proposed

- Today MNIST is considered an easy classification problem.
- It can be solved (almost perfectly) without supervision as a clustering problem, where clustering is performed not in input space (a), but in the latent space of a deep network (b,c).

Aljalbout, Golkov, Siddiqui, Strobel, Cremers, "Clustering with Deep Learning: Taxonomy and New Methods", arxiv 2018. Haeusser, Golkov, Aljalbout, Cremers, "Associative Deep Clustering: Training a Classification Network with no Labels", GCPR 2018.

Test Benchmarks

ImageNet Dataset:
 ImageNet Large Scale Visual Recognition Competition (ILSVRC)





[Russakovsky et al., IJCV'15] "ImageNet Large Scale Visual Recognition Challenge."

Common Performance Metrics

- Top-1 score: check if a sample's top class (i.e. the one with highest probability) is the same as its target label
- Top-5 score: check if your label is in your 5 first predictions (i.e. predictions with 5 highest probabilities)
- → Top-5 error: percentage of test samples for which the correct class was not in the top 5 predicted classes



• Cut ImageNet error down in half







 $227 \times 227 \times 3$



 $227 \times 227 \times 3$



- Use of same convolutions
- As with LeNet: Width, Height Number of Filters





• Softmax for 1000 classes



• Similar to LeNet but much bigger (~1000 times)

• Use of ReLU instead of tanh/sigmoid





• Striving for simplicity

• CONV = 3x3 filters with stride 1, same convolutions

• MAXPOOL = 2x2 filters with stride 2



Conv=3x3,s=1,same Maxpool=2x2,s=2



[Simonyan and Zisserman ICLR'15] VGGNet

I2DL: Prof. Dai



Conv=3x3,s=1,same Maxpool=2x2,s=2



[Simonyan and Zisserman ICLR'15] VGGNet



Conv=3x3,s=1,same Maxpool=2x2,s=2



[Simonyan and Zisserman ICLR'15] VGGNet

I2DL: Prof. Dai

VGGNet





VGGNet





VGGNet

- Conv -> Pool -> Conv -> Pool -> Conv -> FC
- As we go deeper: Width, Height Vumber of Filters

• Called VGG-16: 16 layers that have weights

138M parameters

• Large but simplicity makes it appealing

[Simonyan and Zisserman ICLR'15] VGGNet

VGG<u>Net</u>

		ConvNet C	onfiguration					
Α	A-LRN	В	С	D D	E			
11 weight	11 weight	13 weight	16 weight	16 weight	19 weight			
layers	layers	layers	layers	layers	layers			
	i	nput ($224 imes 22$	24 RGB image					
conv3-64	conv3-64	conv3-64	conv3-64	conv3-64	conv3-64			
	LRN	conv3-64	conv3-64	conv3-64	conv3-64			
maxpool								
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	conv3-256	conv3-256	conv3-256	conv3-256	conv3-256			
			conv1-256	conv3-256	conv3-256			
					conv3-256			
		max	pool	l i				
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512			
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512			
			conv1-512	conv3-512	conv3-512			
					conv3-512			
		max	pool					
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512			
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512			
			conv1-512	conv3-512	conv3-512			
					conv3-512			
	maxpool							
	FC-4096							
FC-4096								
	FC-1000							
	soft-max							

Table 2: Number of parameters (in millions).

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

• A lot of architectures were analyzed

[Simonyan and Zisserman 2014]

I2DL: Prof. Dai



Skip Connections

The Problem of Depth

• As we add more and more layers, training becomes harder

• Vanishing and exploding gradients

• How can we train very deep nets?
• Two layers $x^{L-1} \longrightarrow x^{L+1} \longrightarrow x^{L+1}$



$$\longrightarrow x^{L+1} = f(W^{L+1}x^L + b^{L+1})$$



Main path



• Two layers



- Usually use a same convolution since we need same dimensions
- Otherwise we need to convert the dimensions with a matrix of learned weights or zero padding

ResNet Block







- Xavier/2 initialization
- SGD + Momentum (0.9)

ResNet-152: 60M parameters

- Learning rate 0.1, divided by 10 when plateau
- Mini-batch size 256
- Weight decay of 1e-5
- No dropout



• Without ResNet, if we make the network deeper, at some point performance starts to degrade:



ResNet

With Residual Blocks, performance gets better with deeper network:





• How is this block really affecting me?





• We kept the same values and added a non-linearity

$$x^{L+1} = f(x^{L-1})$$



- The identity is easy for the residual block to learn
- Guaranteed it will not hurt performance, can only improve



Recall: Convolutions on Images





What is the output size?



-10		



-10	6	4	-10	6	
8	6	4	2	-6	
2	0	6	6	10	
-4	0	2	8	8	
10	12	14	18	-2	

	-5	3	2	-5	3		-10	6	4	-10	6
Image 5x5	4	3	2	1	-3		8	6	4	2	-6
	1	0	3	3	5		2	0	6	6	10
	-2	0	1	4	4		-4	0	2	8	8
	5	6	7	9	-1		10	12	14	18	-2

• 1x1 kernel: keeps the dimensions and scales input



• Same as having a 3 neuron fully connected layer



• As always we use more convolutional filters

Using 1x1 Convolutions

- Use it to shrink the number of channels
- Further adds a non-linearity → one can learn more complex functions







• 3x3 max pooling is with stride 1



Not sustainable!

Inception Layer: Computational Cost



Multiplications: 5x5x200 x 32x32x92 ~ 470 million

1 value of the output volume

Inception Layer: Computational Cost



Multiplications: 1x1x200x32x32x16

5x5x16x32x32x92

~ 40 million

Reduction of multiplications by 1/10



(a) Inception module, naïve version

(b) Inception module with dimensionality reduction

[Szegedy et al CVPR'15] GoogLeNet

Inception Layer: Dimensions



GoogLeNet: Using the Inception Layer



reduce dimensionality

[Szegedy et al CVPR'15] GoogLeNet

Xception Net

- "Extreme version of Inception": applying (modified)
 Depthwise Separable Convolutions instead of normal convolutions
- 36 conv layers, structured into several modules with skip connections
- outperforms Inception Net V3



Normal convolutions act on all channels.



Filters are applied only at certain depths of the features. Normal convolutions have groups set to 1, the convolutions used in this image have groups set to 3.

classtorch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, groups=3)

classtorch.nn.ConvTranspose2d(in_channels, out_channels, kernel_size, stride=1, padding=0, groups=3)



But the depth size is always the same!

classtorch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, groups=3)

classtorch.nn.ConvTranspose2d(in_channels, out_channels, kernel_size, stride=1, padding=0, groups=3)



I2DL: Prof. Dai

But why?



Original convolution 256 kernels of size 5x5x3

Multiplications: 256x5x5x3 x (8x8 locations) = 1.228.800

But why?



Original convolution 256 kernels of size 5x5x3

Multiplications: 256x5x5x3 x (8x8 locations) = 1.228.800

Depth-wise convolution 3 kernels of size 5x5x1

256 kernels of size 1x1x3

1x1 convolution

Multiplications: 5x5x3 x (8x8 locations) = 4800

Less computation!

Multiplications: 256x1x1x3x (8x8 locations) = 49152
ImageNet Benchmark

ImageNet Classification top-5-error (%)





Fully Convolutional Network

74





FCN: Becoming Fully Convolutional

convolution



Convert fully connected layers to convolutional layers!

FCN: Becoming Fully Convolutional

convolution



FCN: Upsampling Output

convolution



Semantic Segmentation (FCN)



[Long and Shelhamer. 15] FCN

• 1. Interpolation





• 1. Interpolation

Original image 🌆 🗴 10





Nearest neighbor interpolation Bilinear interpolation

Bicubic interpolation

Image: Michael Guerzhoy

• 1. Interpolation

✓ Few artifacts



[Dosovitskiy, Springenberg, Brox, "Learning to Generate Chairs with Convolutional Networks", CVPR 2015] Prof. Dai

- 2. Transposed convolution
 - Unpooling
 - Convolution filter (learned)
 - Also called up-convolution
 (never deconvolution)



Refined Outputs

• If one does a cascade of unpooling + conv operations, we get to the encoder-decoder architecture

• Even more refined: Autoencoders with skip connections (aka U-Net)



U-Net architecture: Each blue box is a multichannel feature map. Number of channels denoted at the top of the box, dimensions at the top. White boxes are copied feature maps.

[Ronneberger, Fischer, Brox, "U-net: Convolutional networks for biomedical image segmentation", MICCAI'15]

U-Net: Encoder

Left side: Contraction Path (Encoder)

- Captures context of the image
- Follows typical architecture of a CNN:
 - Repeated application of 2 unpadded 3x3 convolutions
 - Each followed by ReLU activation
 - 2x2 maxpooling operation with stride 2 for downsampling
 - At each downsampling step, # of channels is doubled
- \rightarrow as before: Height, Width \downarrow , Depth: \uparrow

[Ronneberger, Fischer, Brox, "U-net: Convolutional networks for biomedical image segmentation", MICCAI'15] Prof. Dai

U-Net: Decoder

Right Side: Expansion Path (Decoder):

- Upsampling to recover spatial locations for assigning class labels to each pixel
 - 2x2 up-convolution that halves number of input channels
 - Skip Connections: outputs of up-convolutions are concatenated with feature maps from encoder
 - Followed by 2 ordinary 3x3 convs
 - final layer: 1x1 conv to map 64 channels to # classes
- Height, Width:
 ↑ Depth:

[Ronneberger, Fischer, Brox, "U-net: Convolutional networks for biomedical image segmentation", MICCAI'15] : Prof. Dai



Object Detection

What is Object Detection?

Classification

Classification + Localization

Object Detection

Instance Segmentation



What is object detection?

- Goal: localize + classify objects in an image
- Input: image
- Output: list of bounding boxes, class labels (+ confidences)

Object Detection



CAT, DOG, DUCK

Challenges in Detection

- Convolutions are translationally equivariant
 - Convolutions share weights: same filters slide across all positions
 - Shifting input image -> output feature maps will shift by same amount



https://chriswolfvision.medium.com/what-is-translation-equivariance-and-why-do-we-use-convolutions-to-get-it-6f18139d4c59

Challenges in Detection

- Convolutions are translationally equivariant
- Good for identifying "what" (e.g., classification), makes "where" challenging

• Detection: need to detect multiple objects, likely different sizes/places

Region-Based CNNs (R-CNNs)



- Use selective search to generate candidate regions
- Warp/rescale each region to fixed image size, pass through a CNN
- CNN extracts features for predicting class + box offsets
- Slow: runs CNN separately for each region

Fast R-CNN

- Instead of running CNN once per region, run CNN once for the entire image
- Use region of interest (ROI) pooling to extract fixedsize feature maps for each region



I2DL: Prof. Dai

- Used to handle variable-sized object proposals efficiently
- Proposed regions have different sizes and aspect ratios
- Fully-connected layers expect fixed-size inputs
- Convert each variable-sized region of interest (ROI) from the feature map and rescale it into a fixed-sized output using max pooling

- Input: 8x8 feature map
- Output size: 2x2

input								
0.88	0.44	0.14	0.16	0.37	0.77	0.96	0.27	
0.19	0.45	0.57	0.16	0.63	0.29	0.71	0.70	
0.66	0.26	0.82	0.64	0.54	0.73	0.59	0.26	
0.85	0.34	0.76	0.84	0.29	0.75	0.62	0.25	
0.32	0.74	0.21	0.39	0.34	0.03	0.33	0.48	
0.20	0.14	0.16	0.13	0.73	0.65	0.96	0.32	
0.19	0.69	0.09	0.86	0.88	0.07	0.01	0.48	
0.83	0.24	0.97	0.04	0.24	0.35	0.50	0.91	

Source: https://deepsense.ai/region-of-interest-pooling-explained/

- Input: 8x8 feature map
- Output size: 2x2

region proposal							
0.88	0.44	0.14	0.16	0.37	0.77	0.96	0.27
0.19	0.45	0.57	0.16	0.63	0.29	0.71	0.70
0.66	0.26	0.82	0.64	0.54	0.73	0.59	0.26
0.85	0.34	0.76	0.84	0.29	0.75	0.62	0.25
0.32	0.74	0.21	0.39	0.34	0.03	0.33	0.48
0.20	0.14	0.16	0.13	0.73	0.65	0.96	0.32
0.19	0.69	0.09	0.86	0.88	0.07	0.01	0.48
0.83	0.24	0.97	0.04	0.24	0.35	0.50	0.91

Source: https://deepsense.ai/region-of-interest-pooling-explained/

- Input: 8x8 feature map
- Output size: 2x2

pooling sections								
0.88	0.44	0.14	0.16	0.37	0.77	0.96	0.27	
0.19	0.45	0.57	0.16	0.63	0.29	0.71	0.70	
0.66	0.26	0.82	0.64	0.54	0.73	0.59	0.26	
0.85	0.34	0.76	0.84	0.29	0.75	0.62	0.25	
0.32	0.74	0.21	0.39	0.34	0.03	0.33	0.48	
0.20	0.14	0.16	0.13	0.73	0.65	0.96	0.32	
0.19	0.69	0.09	0.86	0.88	0.07	0.01	0.48	
0.83	0.24	0.97	0,04	0.24	0.35	0.50	0.91	

Source: https://deepsense.ai/region-of-interest-pooling-explained/

Faster R-CNN

- Use convolutions to get sliding window effect
- Replace selective search with a region proposal network (RPN)
 - Small CNN that slides over feature map
 - Produces region proposals based on anchors



[Ren et al. NeurIPS'15] Faster R-CNN

Anchors

• Potential bounding box candidates where an object can be detected



Instance Segmentation

 Predict list of object bounding boxes, classes, and pixel-wise masks

Classification

Classification + Localization

Object Detection

Instance Segmentation



12DL: Proi, Dai

Mask R-CNN

• Extend Faster R-CNN to also predict pixel-wise mask for each object



C x 14 x 14

[He et al. ICCV'17] Mask R-CNN

R-CNNs for Detection + Segmentation

- R-CNN, Fast R-CNN, Faster R-CNN, Mask-RCNN: two-stage detectors
- First generate region proposals, then refine boxes + predict class labels (+ masks)

• Today: end-to-end, unified approach with transformers

Next Time: Guest Lecture

July 14: Guest Lecture!
In person, regular lecture slot (recording released later)



Phillip Isola, MIT

- Leader in vision and generative models
- Created Pix2Pix, CycleGAN,
 - representation learning
- Foundations of Computer
 Vision textbook
- Live Q&A ask questions directly!





See you next time!

References

We highly recommend to read through these papers!

- <u>AlexNet</u> [Krizhevsky et al. 2012]
- <u>VGGNet</u> [Simonyan & Zisserman 2014]
- <u>ResNet</u> [He et al. 2015]
- <u>GoogLeNet</u> [Szegedy et al. 2014]
- <u>Xception</u> [Chollet 2016]
- Fast R-CNN [Girshick 2015]
- <u>U-Net</u> [Ronneberger et al. 2015]
- <u>EfficientNet</u> [Tan & Le 2019]