

Convolutional Neural Networks

I2DL: Prof. Dai

Fully Connected Neural Network















Why not simply more FC Layers?

We cannot make networks arbitrarily complex

- Why not just go deeper and get better?
 - No structure!!
 - It is just brute force!
 - Optimization becomes hard
 - Performance plateaus / drops!

Better Way than FC?

- We want to restrict the degrees of freedom
 - We want a layer with structure
 - Weight sharing → using the same weights for different parts of the image

Using CNNs in Computer Vision

Classification

Classification + Localization

Object Detection

Instance Segmentation



l2DL: Prof. Da



Convolutions



Convolution of two box functions

Convolution of two Gaussians

Application of a filter to a function — The 'smaller' one is typically called the filter kernel

Discrete case: box filter



'Slide' **filter kernel** from left to right; at each position, compute a single value in the output data

















Discrete case: box filter



What to do at boundaries?

Discrete case: box filter



What to do at boundaries?

Option 1: Shrink

3 0 0 1 10/3 4 4 16/3	
-----------------------	--



7/3	3	0	0	1	10/3	4	4	16/3	11/3
-----	---	---	---	---	------	---	---	------	------



















Image Filters

• Each kernel gives us a different image filter




Convolutions on RGB Images

$32 \times 32 \times 3$ image (pixels **X**)

3

32



 Z_i

 1 number at a time:
 equal to dot product between filter weights w and x_i - th chunk of the image. Here: 5 · 5 · 3 = 75-dim dot product + bias

 $z_i = \boldsymbol{w}^T \boldsymbol{x}_i + \boldsymbol{b}$

$$(5 \times 5 \times 3) \times 1$$
 $(5 \times 5 \times 3) \times 1$





Convolution Layer

Convolution Layer





Convolution Layer

- A basic layer is defined by
 - Filter width and height (depth is implicitly given)
 - Number of different filter banks (#weight sets)

• Each filter captures a different image characteristic

Different Filters



- Each filter captures different image characteristics:
 - Horizontal edges
 - Vertical edges
 - Circles
 - Squares

....

[Zeiler & Fergus, ECCV'14] Visualizing and Understanding Convolutional Networks I2DL: Prof. Dai



Dimensions of a Convolution Layer



Input:	7×7
Filter:	3 × 3
Output:	5 × 5



Input:	7×7
Filter:	3×3
Output:	5 × 5



Input:	7×7
Filter:	3 × 3
Output:	5 × 5



Input:	7×7
Filter:	3 × 3
Output:	5×5

			5	
<				
5				
-				

Input:	7×7
Filter:	3 × 3
Output:	5 × 5



Input:	7×7
Filter:	3×3
Stride:	1
Output:	5×5

Stride of *S*: apply filter every *S*-th spatial location; i.e. subsample the image



Input:	7×7
Filter:	3×3
Stride:	2
Output:	3×3



Input:	7×7
Filter:	3×3
Stride:	2
Output:	3×3



Input:	7×7
Filter:	3×3
Stride:	2
Output:	3×3



Input:	7×7
Filter:	3 × 3
Stride:	3
Output:	?×?



Input:	7×7
Filter:	3×3
Stride:	3
Output:	?×?



	input width of N								
				of F					
N				ght d					
t of				r hei					
lgh	Filter \	vidth c	of F	Filte					
It he									
ndu									
_									

1 1 1 11 1

Input:	N imes N
Filter	$F \times F$
Stride	S
Output:	$\left(\frac{N-F}{S}+1\right) \times \left(\frac{N-F}{S}+1\right)$

 $N = 7, F = 3, S = 1; \frac{7-3}{1} + 1 = 5$ $N = 7, F = 3, S = 2; \frac{7-3}{2} + 1 = 3$ $N = 7, F = 3, S = 3; \frac{7-3}{3} + 1 = 2.\overline{3}$ Fractions are illegal



Shrinking down so quickly $(32 \rightarrow 28 \rightarrow 24 \rightarrow 20)$ is typically not a good idea...

Why padding?

- Sizes get small too quickly
- Corner pixel is only used once

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

Why padding?

- Sizes get small too quickly
- Corner pixel is only used once

+ Zero padding

 7×7

Image

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

Input $(N \times N)$: 7×7 Filter $(F \times F)$: 3×3 Padding (P):1Stride (S):1Output 7×7

Most common is 'zero' padding

Output Size:

$$\left(\left\lfloor\frac{N+2\cdot P-F}{S}\right\rfloor+1\right)\times \left(\left\lfloor\frac{N+2\cdot P-F}{S}\right\rfloor+1\right)$$

[] denotes the floor operator (as in practice an integer division is performed)

+ zero padding

∽ ×

Image

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 $P = \frac{F-1}{2}$

+ Zero padding

∽ ×

~

Image



 $i \in 32 \times 32 \times 10$

2DL: Prof. Dai





Number of parameters (weights): Each filter has $5 \times 5 \times 3 + 1 = 76$ params (+1 for bias) -> $76 \cdot 10 = 760$ parameters in layer



- You are given a convolutional layer with **4** filters, kernel size **5**, stride **1**, and no padding that operates on an RGB image.
- Q1: What are the dimensions and the shape of its weight tensor?

A1: (3,4,5,5) A2: (4,5,5) A3: depends on the width and height of the image



- You are given a convolutional layer with **4** filters, kernel size **5**, stride **1**, and no padding that operates on an RGB image.
- Q1: What are the dimensions and the shape of its weight tensor?





Convolutional Neural Network (CNN)

CNN Prototype

ConvNet is concatenation of Conv Layers and activations



CNN Learned Filters



[Zeiler & Fergus, ECCV'14] Visualizing and Understanding Convolutional Networks 12DL: Prof. Dai



Pooling



[Li et al., CS231n Course Slides] Lecture 5: Convolutional Neural Networks 2DL: Prof. Dai
Pooling Layer: Max Pooling

Single depth slice of input





Pooling Layer

Conv Layer = 'Feature Extraction'
Computes a feature in a given region

- Pooling Layer = 'Feature Selection'
 - Picks the strongest activation in a region

Pooling Layer

- Input is a volume of size $W_{in} \times H_{in} \times D_{in}$
- Two hyperparameters

 - Stride S

Spatial filter extent F
Stride S

• Output volume is of size $W_{out} \times H_{out} \times D_{out}$

$$-W_{out} = \frac{W_{in}-F}{S} + 1$$

$$-H_{out} = \frac{H_{in} - F}{S} + 1$$

$$- D_{out} = D_{in}$$

Does not contain parameters; e.g. it's fixed function

Pooling Layer

- Input is a volume of size $W_{in} \times H_{in} \times D_{in}$
- Two hyperparameters
 - Spatial filter extent **F**
 - Stride S
- Output volume is of size $W_{out} \times H_{out} \times D_{out}$

$$-W_{out} = \frac{W_{in} - F}{S} + 1$$

$$-H_{out} = \frac{H_{in}-F}{S} + 1$$

$$- D_{out} = D_{in}$$

• Does not contain parameters; e.g. it's fixed function

Common settings: *F* = 2, S = 2 *F* = 3, *S* = 2

Pooling Layer: Average Pooling

Single depth slice of input



Average pool with 2 × 2 filters and stride 2

'Pooled' output

2.5	6
1.75	3

• Typically used deeper in the network

CNN Prototype



Final Fully-Connected Layer

- Same as what we had in 'ordinary' neural networks
 - Make the final decision with the extracted features from the convolutions
 - One or two FC layers typically

Convolutions vs Fully-Connected

- In contrast to fully-connected layers, we want to restrict the degrees of freedom
 - FC is somewhat brute force
 - Convolutions are **structured**
- Sliding window to with the same filter parameters to extract image features
 - Concept of weight sharing
 - Extract same features independent of location



• Spatial extent of the connectivity of a convolutional filter





• Spatial extent of the connectivity of a convolutional filter







3x3 output



3x3 receptive field = 1 output pixel is connected to 9 input pixels

• Spatial extent of the connectivity of a convolutional filter





• Spatial extent of the connectivity of a convolutional



3x3 receptive field = 1 output pixel is connected to 9 input pixels

7x7 input

• Spatial extent of the connectivity of a convolutional



3x3 receptive field = 1 output pixel is connected to 9 input pixels

7x7 input

Spatial extent of the connectivity of a convolutional filter



5x5 receptive field on the original input: one output value is connected to 25 input pixels

7x7 input



See you next time!

92

References

Goodfellow et al. "Deep Learning" (2016),
– Chapter 9: Convolutional Networks

• <u>http://cs231n.github.io/convolutional-networks/</u>

 Useful info on convolutions in image processing: <u>https://visionbook.mit.edu/linear_image_filtering.html</u>