

Αλληλεπίδραση Ανθρώπου-Υπολογιστή

B7. Υπολογιστική όραση με συνελικτικά νευρωνικά δίκτυα (CNNs)

(2024-25)

Ίων Ανδρουτσόπουλος http://www.aueb.gr/users/ion/

Contents

- Convolutional neural networks (CNNs) and applications in image classification and object detection.
- Image to text generation with CNN encoders and RNN decoders.

Averaging each pixel with its neighboring values blurs an image:



0000

From the blog post "Understanding Convolutional Neural Networks for NLP" of Denny Britz, 2015. <u>http://www.wildml.com/</u> <u>2015/11/understanding-</u> <u>convolutional-neural-</u> <u>networks-for-nlp/</u>

Feature Map

		Input					Ke	mel (Fil	ter)
-1	1	-1	-1	-1	-1	1	-1	1	-1
1	1	1	-1	-1	-1		1	1	1
-1	1	-1	-1	-1	-1		-1	1	-1
-1	-1	-1	1	-1	1				
-1	-1	-1	-1	1	-1				
-1	-1	-1	1	-1	1				

- **Input: black/white image** with pixel values -1 or +1.
- Check if the input contains any crosses and report where.

Input

Kernel (Filter)

Feature Map

-1	1	-1	-1	-1	-1
1	1	1	-1	-1	-1
-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1
-1	-1	-1	-1	1	-1
-1	-1	-1	1	-1	1

-1	1	-1
1	1	1
-1	1	-1

-1	1	-1	-1	-1	-1
1	1	1	-1	-1	-1
-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1
-1	-1	-1	-1	1	-1
-1	-1	-1	1	-1	1

-1	1	-1
1	1	1
-1	1	-1



Input

Kernel (Filter)

Feature Map

-1	1	-1	-1	-1	-1
1	1	1	-1	-1	-1
-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1
-1	-1	-1	-1	1	-1
-1	-1	-1	1	-1	1

-1	1	-1
1	1	1
-1	1	-1

-1	1	-1	-1	-1	-1
1	1	1	-1	-1	-1
-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1
-1	-1	-1	-1	1	-1
-1	-1	-1	1	-1	1

-1	1	-1
1	1	1
-1	1	-1



-1	1	-1	-1	-1	-1
1	1	1	-1	-1	-1
-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1
-1	-1	-1	-1	1	-1
-1	-1	-1	1	-1	1





Input						ĸ	ernel (F	ilte	r)		Featu	ıre Map		
-1	1	-1	-1	-1	-1	-1	1		-1	1				
1	1	1	-1	-1	-1	1	1		1		9	-1	1	
-1	1	-1	-1	-1	-1	-1	1		-1					
-1	-1	-1	1	-1	1					•				
-1	-1	-1	-1	1	-1									
-1	-1	-1	1	-1	1									

- Let X be the part of the input where we apply the kernel (filter).
- Let *W* be the kernel.
- The resulting **feature** of the feature map is: $\sum_{i=1}^{3} \sum_{j=1}^{3} W_{i,j} X_{i,j}$
- In practice, we would also use an **activation function** and **bias term**: $f(\sum_{i=1}^{3} \sum_{j=1}^{3} W_{i,j}X_{i,j} + b)$

		Inpu	ut			Kernel (Filter)	Feature Map
-1 1 -1 -1 -1 -1	1 1 -1 -1	-1 1 -1 -1 -1	-1 -1 -1 1 -1 1	-1 -1 -1 -1 1 -1	-1 -1 -1 1 -1 1	-1 1 -1 1 1 1 -1 1 -1	9 -1 1
-1 1 -1 -1 -1	1 1 -1 -1	-1 1 -1 -1 -1	-1 -1 -1 1 -1	-1 -1 -1 -1 1	-1 -1 -1 1 -1	-11-1111-11-1	9 -1 1 -1 -1 -1 -1 -5 1 -1 -1 5 -1 -5 5 -7

• We can think of the resulting **feature map as a new "image"** that indicates the **position(s) of the cross(es)** in the original image.

No need to have the crosses at particular parts of the image.

• The new "image" is **4x4 instead of 6x6**, because the **kernel could not slide outside the boundaries** of the original image.



- We can **pad** the surrounding of the image with zeros, to allow the kernel to slide outside the image boundaries.
- We can now obtain a **feature map** with the **same resolution as the input** image (6x6).

Input

Kernel (Filter)

Feature Map

0	0	0	0	0	0	0	0
0	-1	1	-1	-1	-1	-1	0
0	1	1	1	-1	-1	-1	0
0	-1	1	-1	-1	-1	-1	0
0	-1	-1	-1	1	-1	1	0
0	-1	-1	-1	-1	1	-1	0
0	-1	-1	-1	1	-1	1	0
0	0	0	0	0	0	0	0

Г	-1	1	-1
	1	1	1
	-1	1	-1

	0	0	0	0	0	0	0	0
I	0	-1	1	-1	-1	-1	-1	0
L	0	1	1	1	-1	-1	-1	0
	0	-1	1	-1	-1	-1	-1	0
	0	-1	-1	-1	1	-1	1	0
	0	-1	-1	-1	-1	1	-1	0
	0	-1	-1	-1	1	-1	1	0
	0	0	0	0	0	0	0	0
					_			
	0	0	0	0	0	0	0	0
	0 0	0 -1	0	0 -1	0 -1	0 -1	0 -1	0 0
	0 0 0	0 -1 1	0 1 1	0 -1 1	0 -1 -1	0 -1 -1	0 -1 -1	0 0 0
	0 0 0	0 -1 1 -1	0 1 1 1	0 -1 1 -1	0 -1 -1 -1	0 -1 -1 -1	0 -1 -1 -1	0 0 0
	0 0 0 0	0 -1 1 -1 -1	0 1 1 1 -1	0 -1 1 -1 -1	0 -1 -1 -1 1	0 -1 -1 -1 -1	0 -1 -1 -1 1	0 0 0 0
	0 0 0 0 0	0 -1 1 -1 -1 -1	0 1 1 -1 -1	0 -1 1 -1 -1 -1	0 -1 -1 -1 1 -1	0 -1 -1 -1 -1 -1 1	0 -1 -1 -1 1 -1	0 0 0 0 0

-1	1	-1
1	1	1
-1	1	-1

-1

-1

-1

-1



0	-2		

Input

Kernel (Filter)

Feature Map

0	0	0	0	0	0	0	0
0	-1	1	-1	-1	-1	-1	0
0	1	1	1	-1	-1	-1	0
0	-1	1	-1	-1	-1	-1	0
0	-1	-1	-1	1	-1	1	0
0	-1	-1	-1	-1	1	-1	0
0	-1	-1	-1	1	-1	1	0
0	0	0	0	0	0	0	0

....

0

-1

1

-1

-1

-1

-1

0

0

-1

1

-1

-1

-1

-1

0

0

-1

-1

-1

1

-1

1

0

0

-1

-1

-1

1

-1

1

0

0

-1

-1

-1

-1

1

-1

0

0

-1

-1

-1

-1

1

-1

0

0

-1

-1

-1

1

-1

1

0

0

-1

-1

-1

1

-1

1

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

-1

1

-1

-1

-1

-1

0

0

-1

1

-1

-1

-1

-1

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

1

1

1

-1

-1

-1

0

0

1

1

1

-1

-1

-1

-1	1	-1
1	1	1
-1	1	-1

-2	0		
	-2	-2 0	-2 0

-1	1	-1
1	1	1
-1	1	-1

-2	0	-4	-2	-2
	-2	-2 0	-2 0 -4	-2 0 -4 -2

-1	1	-1
1	1	1
-1	1	-1

0	-2	0	-4	-2	-2
-2	9				

Input

Kernel (Filter)

Feature Map

0	0	0	0	0	0	0	0
0	-1	1	-1	-1	-1	-1	0
0	1	1	1	-1	-1	-1	0
0	-1	1	-1	-1	-1	-1	0
0	-1	-1	-1	1	-1	1	0
0	-1	-1	-1	-1	1	-1	0
0	-1	-1	-1	1	-1	1	0
0	0	0	0	0	0	0	0

-1	1	-1
1	1	1
-1	1	-1

0	-2	0	-4	-2	-2
-2	9	-1			

0	0	0	0	0	0	0	
-1	1	-1	-1	-1	-1	0	
1	1	1	-1	-1	-1	0	
-1	1	-1	-1	-1	-1	0	
-1	-1	-1	1	-1	1	0	_
-1	-1	-1	-1	1	-1	0	
-1	-1	-1	1	-1	1	0	
0	0	0	0	0	0	0	

-1	1	-1
1	1	1
-1	1	-1

0	-2	0	-4	-2	-2
-2	9	-1	1	-1	-2
0	-1	-1	-1	-5	0
-4	1	-1	-1	5	-2
-2	-1	-5	5	-7	4
-2	-2	0	-2	4	-2

- X: entire input image. F: feature map.
- W: kernel, but with rows and columns numbered -1, 0, 1.
- Feature map values: $F_{i,j} = \sum_{k=-1}^{1} \sum_{l=-1}^{1} W_{k,l} X_{i+k,j+l}$
- In practice: $F_{i,j} = f(\sum_{k=-1}^{1} \sum_{l=-1}^{1} W_{k,l} X_{i+k,j+l} + b)$

Convolution or cross-correlation?

- **Cross-correlation**: $F_{i,j} = \sum_{k=-\infty}^{+\infty} \sum_{l=-\infty}^{+\infty} W_{k,l} X_{i+k,j+l}$ Optional study
- Convolution: $F_{i,j} = \sum_{k=-\infty}^{+\infty} \sum_{l=-\infty}^{+\infty} W_{k,l} X_{i-k,j-l} = W * X$
- We are actually computing cross-correlations, not convolutions.
 - The **cross-correlations** we compute are **equal to convolutions with the kernel (or the image) flipped** both vertically and horizontally.
 - Convolution is like cross-correlation, but flips one of the two signals. We don't flip the kernel inside the cross-correlation, which is equivalent to giving the kernel already flipped to the convolution; the convolution will flip the kernel once more, ending up using the kernel without flipping.
 - So we actually compute **convolutions with flipped kernels** or **cross-correlations with the original kernels**.
 - The example kernels were symmetric, so no difference.
 - In CNNs (Convolutional Neural Networks), the kernels are learned, so we don't care if they are flipped in the "convolutions" we compute.
 - So we usually say **CNNs "compute convolutions"**, though we actually use the formulae of cross-correlations.

Two kernels

			Input					Two Kernels	Feature Map of Kernel 1 ("+")	Feature Map of Kernel 2 ("X")
0 0 0 0 0 0 0	0 -1 -1 -1 -1 -1 -1 0	0 1 1 -1 -1 -1 0	0 -1 -1 -1 -1 -1 -1 0	0 -1 -1 -1 1 -1 1 0	0 -1 -1 -1 1 -1 0	0 -1 -1 -1 1 -1 1 0	0 0 0 0 0 0	-1 1 -1 1 1 1 -1 1 -1 1 -1 1 -1 1 -1 1 -1 1 -1 1 -1 1 -1 1		
0 0 0 0 0 0	0 -1 1 -1 -1 -1 -1 -1 0	0 1 1 -1 -1 -1 0	0 -1 -1 -1 -1 -1 -1 -1	0 -1 -1 -1 1 -1 1 0	0 -1 -1 -1 -1 1 -1	0 -1 -1 -1 -1 -1 1 0	0 0 0 0 0	-1 1 -1 1 1 1 -1 1 -1 1 -1 1 -1 1 -1 1 -1 1 1 -1 1 1 1 1	0	-2

- We now want to check the input image for crosses and "X"s.
- We use **two kernels**, one for crosses, one for "X"s.

Two kernels

			Input					Tw	0
0	0	0	0	0	0	0	0	-1	
0	-1	1	-1	-1	-1	-1	0	1	
0	1	1	1	-1	-1	-1	0	-1	
0	-1	1	-1	-1	-1	-1	0		
0	-1	-1	-1	1	-1	1	0	1	
0	-1	-1	-1	-1	1	-1	0	-1	
0	-1	-1	-1	1	-1	1	0	1	
0	0	0	0	0	0	0	0		
_									
0	0	0	0	0	0	0	0	-1	
0	-1	1	-1	-1	-1	-1	0	1	
0	1	1	1	-1	-1	-1	0	-1	
0	-1	1	-1	-1	-1	-1	0		
0	-1	-1	-1	1	-1	1	0	1	ſ

-1

1

0

0

-1

-1

-1

1

-1

1

0

0

-1

-1

-1

1

-1

1

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

-1

-1

0

0

-1

1

-1

-1

-1

-1

0

0

-1

1

-1

-1

-1

-1

0

-1

-1

0

0

1

1

1

-1

-1

-1

0

0

1

1

1

-1

-1

-1

0

-1

-1

0

0

-1

1

-1

-1

-1

-1

0

0

-1

1

-1

-1

-1

-1

0

1

-1

0

0

-1

-1

-1

-1

1

-1

0

0

-1

-1

-1

-1

1

-1

0

-1

1

0

0

-1

-1

-1

1

-1

1

0

0

-1

-1

-1

1

-1

1

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

-1	1	-1
1	1	1
-1	1	-1
1	-1	1
-1	1	-1

-1

1

-1

1

-1

1

1

-1

1

-1 1 1 1 0 -2 -1 1 -1 1 -1 -1 1 -1 1 -1 -1 1 1 -2 -2 1 1 0 4 -1 1 -1 1 -1 1 1 -1 -1 -1 1 1 -1 -1 1 1 1 -2 0 -2 1 0 -2 4 -1 -1

Feature Map of Kernel 1 ("+")

Feature Map of Kernel 2 ("X")

Two kernels

We can think of the two feature maps as two "channels" of the new image, one for "+" info, one for "X" info.

Input

_	0	0	0	0	0	0	0
ſ	0	-1	1	-1	-1	-1	-1
I	0	1	1	1	-1	-1	-1
L	0	-1	1	-1	-1	-1	-1
	0	-1	-1	-1	1	-1	1
	0	-1	-1	-1	-1	1	-1
	0	-1	-1	-1	1	-1	1
	0	0	0	0	0	0	0

0

-1

1

-1

-1

-1

-1

0

0

-1

-1

-1

1

-1

1

0

0

-1

-1

-1

-1

1

-1

0

0

1

1

1

-1

-1

-1

0

0

0

0

0

0

0

0

0

0

0 0

0

0

0

0

0

-1

-1

0

-1

-1

-1

0

-1

1

-1

-1

-1

-1

0

-1	1	-1
1	1	1
-1	1	-1
1	-1	1
-1	1	-1
-	-	-

Two Kernels

0	-2	0	-4	-2	-2
-2					

Feature Map of Kernel 1 ("+")

-2 4	4	-2	2	0	0

Feature Map of Kernel 2 ("X")

0	0	-1	1	
-1	0	1	1	
-1	0	-1	1	ſ
-1	0			
1	0	1	-1	
-1	0	-1	1	
1	0	- 1	- 1	

0 0

0

0

0

0

0

0

-1	1	-1	
1	1	1	
-1	1	-1	
1	-1	1	
1 -1	-1 1	1 -1	

0	2	0	4	2	2
0	-2	0	-4	-2	-2
-2	9				

-2	4	-2	2	0	0
4	-7				

)	0	0	0	0	0	0
)	-1	1	-1	-1	-1	-1
)	1	1	1	-1	-1	-1
)	-1	1	-1	-1	-1	-1
)	-1	-1	-1	1	-1	1
)	-1	-1	-1	-1	1	-1
)	-1	-1	-1	1	-1	1
· ·	0	0	0	0	0	0

-1	1	-1
1	1	1
-1	1	-1
1	-1	1
-1	1	-1

_						
Г	0	-2	0	-4	-2	-2
I	-2	9	-1			
I						
I						
I						
I						

Г	-2	4	-2	2	0	0
L	4	-7	3			
L						
Т						

0	0	0	0	0	0
-1	1	-1	-1	-1	-1
1	1	1	-1	-1	-1
-1	1	-1	-1	-1	-1

-1

-1

-1

0

1

-1

1

0

-1

1

-1

0

1

-1

1

0

1 -1	1 1	1 -1
1	-1	1
-1	1	-1
1	-1	1

-1 1

-1

0	-2	0	-4	-2	-2
-2	9	-1	1	-1	-2
0	-1	-1	-1	-5	0
-4	1	-1	-1	5	-2
-2	-1	-5	5	-7	4
-2	-2	0	-2	4	-2

-2	4	-2	2	0	0
4	-7	3	-3	-1	0
-2	3	-1	-1	3	-2
2	-3	-1	3	-7	4
0	-1	3	-7	9	-6
0	0	-2	4	-6	4
	-2 4 -2 2 0 0	-2 4 4 -7 -2 3 2 -3 0 -1 0 0	-2 4 -2 4 -7 3 -2 3 -1 2 -3 -1 0 -1 3 0 0 -2	-2 4 -2 2 4 -7 3 -3 -2 3 -1 -1 2 -3 -1 3 0 -1 3 -7 0 0 -2 4	-2 4 -2 2 0 4 -7 3 -3 -1 -2 3 -1 -1 3 2 -3 -1 3 -7 0 -1 3 -7 9 0 0 -2 4 -6

Two input channels too

		Input C	hanne	11							Input	t Chann	el 2					т	wo Two	o-chan	nel Kerne	els		F	ature Map of Kernel 1 ("+")	Feature Map of Kernel 2 ("X")
	0 -1 -1 -1 -1 -1 -1 0	0 1 1 -1 -1 -1 0	0 -1 -1 -1 -1 -1 -1 0	0 -1 -1 0,9 -1 0,9 0,9	0 -1 -1 -1 -1 0,9 -1 0	0 -1 -1 -1 0,9 -1 0,9 0	0 0 0 0 0 0 0		0 0 0 0 0 0	0 -1 0,9 -1 -1 -1 -1 0	0 0,9 0,9 -1 -1 -1 -1 0	0 -1 0,9 -1 -1 -1 -1 -1 0	0 -1 -1 -1 1 -1 1 0	0 -1 -1 -1 -1 1 -1 0	0 -1 -1 -1 1 -1 1 0	0 0 0 0 0 0 0	-1 1 -1 1 -1 1	1 1 -1 1 -1	-1 1 -1 1 -1 1	&	-1 1 -1 1 -1 1	1 1 -1 1 -1	-1 1 -1 1 -1 1			
0 0 0 0 0 0	0 -1 -1 -1 -1 -1 -1	0 1 1 -1 -1 -1 -1	0 -1 1 -1 -1 -1 -1	0 -1 -1 0,9 -1 0,9	0 -1 -1 -1 -1 0,9 -1	0 -1 -1 -1 0,9 -1 0,9	0 0 0 0 0	[0 0 0 0 0 0	0 -1 0,9 -1 -1 -1 -1	0 0,9 0,9 -1 -1 -1	0 -1 0,9 -1 -1 -1 -1	0 -1 -1 -1 1 -1 1 1	0 -1 -1 -1 -1 1 -1	0 -1 -1 -1 1 -1 1 1	0 0 0 0 0 0	-1 1 -1 1 -1 1	1 1 1 -1 1 -1	-1 1 -1 1 -1 1	&	-1 1 -1 1 -1 1	1 1 1 -1 1 -1	-1 1 -1 1 -1 1	-0,	ı	-3,9

• The **input image** now also has **two channels** (e.g., from grayscale and depth cameras). **Each kernel** now operates on **both input channels**.

• It has **two slices**, one per input channel (c = 1, c = 2).

- We have **two kernels**, so the **output** also has **two channels**.
- At the output feature map of kernel $W^{(m)}$, the value at cell (i, j) is:

$$F_{i,j,m} = \sum_{k=-1}^{1} \sum_{l=-1}^{1} \sum_{c=1}^{2} W_{k,l,c}^{(m)} X_{i+k,j+l,c}$$

• In practice, we would also have an activation function and bias term.

Two input channels too

Input Channel 1	Input Channel 2	Two Two-channel Kernels	Feature Map of Kernel 1 ("+")	Feature Map of Kernel 2 ("X")
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	-0,1 -4	-3,9 7,8
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	-0,1 -4 -0,1	-3,9 7,8 -3,9
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	-0,1 -4 -0,1 -7,9 -4 -4 -4	-3,9 7,8 -3,9 3,9 0 0 7,8
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	-0,1 -4 -0,1 -7,9 -4 -4 -4 17,5	-3,9 7,8 -3,9 3,9 0 0 7,8

Two input channels too

		Inp	ut Char	nel 1							Input C	annel 2					Т	wo Two	-chann	nel Kerr	nels			Fe	ature M	Nap of	(ernel 1	("+")		Feat	ure Map	o of Ker	mel 2 ("	X")	
0 0	0 -1	0	0 -1	0 -1	0 -1	0 -1	0 0	0	0	L	0 0 0,9 -1	0	0	0 -1	0	-1 1	1 1	-1 1	&	-1 1	1	-1 1	1	-0,1	4	-0,1	-7,9	-4	-4	 -3,9	7,8	-3,9	3,9	0	0
0	1	1	1	-1	-1	-1	0	0	0,9	9 0	0,9 O,	9 -1	-1	-1	0	-1	1	-1		-1	1	-1		-4	17,5	-2				7,8	-13,7	5,8			
0	-1	-1	-1	0,9	-1	0,9	0	0	-1		-1 -1	. 1	-1	1	0	1	-1	1		1	-1	1	1												
0	-1	-1	-1	0,9	-1	0,9	0	0	-1		-1 -1	1	-1	1	0	1	-1	1	¢.	-1	-1	1													
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	-1	1	-1		-1	1	-1]	_											
0	-1	1	-1	-1	-1	-1	0	0	-1	L (0,9 -1 no o	-1 • -1	-1	-1	0	1	1	1	&	1	1	1		-0,1	-4	-0,1	-7,9	-4	-4	-3,9 7.8	7,8	-3,9 5.8	3,9	0	0
0	-1	1	-1	-1	-1	-1	o	0	-1	L C	0,9 -1 0,9 -1	-1	-1	-1	o	1	1	1		1	1		3	-0,1	-2	-2	-2	-9,8	-0,1	-3,9	5,8	-2	-2	5,8	-3,9
0	-1	-1	-1	0,9	-1	0,9	0	. 0	-1	L -	-1 -1	1	-1	1	0	1	-1	1		1	-1	1	1	-7,9	1,9	-2	-2	9,7	-4	3,9	-5,9	-2	5,8	-13,7	7,8
0	-1	-1	-1	-1	0,9	-1	0	0	-1	L ·	-1 -1	-1	1	-1	0	-1	1	-1	&	-1	1	-1		-4	-2	-9,8	9,7	-13,7	7,7	0	-2	5,8	-13,7	17,5	-11,7
<u> </u>	-1	-1	-1	0,9	-1	0,9	0	0	-1		-1 -1	. 1	-1	1		1	-1	1		1	-1	1		-4	-4	-0,1	-4	1,1	-4	 U	U	-5,9	7,8	-11,7	7,8

- We now have a mechanism, a "convolutional layer", that maps an input image of any number of channels to a new output "image" of any number of channels (feature maps).
 - The kernels will have as many slices as the input channels.
 - The number of kernels will be equal to the number of output channels.
- We can stack multiple convolutional layers.
 - Each one will operate on the "image" produced by the previous layer.
 - All kernels will be randomly initialized and learned via backpropagation.

Max-pooling

Feature Map of Kernel 1 ("+")

Feature Map of Kernel 2 ("X")

Max-Pooling (2,2) with Stride (2,2)

-0,1	-4	-0,1	-7,9	-4	-4	1	-3,9	7,8	-3,9	3,9	0	0		17,5			1	7,8		
-4	17,5	-2	1,9	-2	-4		7,8	-13,7	5,8	-5,9	-2	0								
-0,1	-2	-2	-2	-9,8	-0,1		-3,9	5,8	-2	-2	5,8	-3,9					1			
-7,9	1,9	-2	-2	9,7	-4		3,9	-5,9	-2	5,8	-13,7	7,8					-			
-4	-2	-9,8	9,7	-13,7	7,7		0	-2	5,8	-13,7	17,5	-11,7								
-4	-4	-0,1	-4	7,7	-4		0	0	-3,9	7,8	-11,7	7,8								
						-							I							
																	_			
-0,1	-4	-0,1	-7,9	-4	-4	1	-3,9	7,8	-3,9	3,9	0	0		17,5	1,9		1	7,8	5,8	
-4	17,5	-2	1,9	-2	-4		7,8	-13,7	5,8	-5,9	-2	0								
-0,1	-2	-2	-2	-9,8	-0,1		-3,9	5,8	-2	-2	5,8	-3,9								
-7,9	1,9	-2	-2	9,7	-4		3,9	-5,9	-2	5,8	-13,7	7,8		Ť						
-4	-2	-9,8	9,7	-13,7	7,7		0	-2	5,8	-13,7	17,5	-11,7								
-4	-4	-0,1	-4	7,7	-4		0	0	-3,9	7,8	-11,7	7,8								
													I							
-0,1	-4	-0,1	-7,9	-4	-4	1	-3,9	7,8	-3,9	3,9	0	0		17,5	1,9	-2	1	7,8	5,8	0
-4	17,5	-2	1,9	-2	-4		7,8	-13,7	5,8	-5,9	-2	0								
-0,1	-2	-2	-2	-9,8	-0,1	1	-3,9	5,8	-2	-2	5 <u>,</u> 8	-3,9								
-7,9	1,9	-2	-2	9,7	-4		3,9	-5,9	-2	5,8	-13,7	7,8		ř			-			
-4	-2	-9,8	9,7	-13,7	7,7		0	-2	5,8	-13,7	17,5	-11,7								
-4	-4	-01	-4	77	-4		0	0	-39	78	-11 7	78								

- We keep the **max value of each window**, separately from each channel.
- The stride determines how much the window shifts vertically & horizontally.

Max-pooling

Fea	ture Ma	ap of Ke	ernel 1	("+")			Feat	ure Ma	p of Ke	rnel 2 ("X")			I	Max-Po	ooling (2,2) wi	th Stride	e (2,2)	
-0,1	-4	-0,1	-7,9	-4	-4		-3,9	7,8	-3,9	3,9	0	0		17,5	1,9	-2		7,8	5,8	0
-4	17,5	-2	1,9	-2	-4		7,8	-13,7	5,8	-5,9	-2	0		1,9				5,8		
-0,1	-2	-2	-2	-9,8	-0,1		-3,9	5,8	-2	-2	5,8	-3,9								
-7,9	1,9	-2	-2	9,7	-4		3,9	-5,9	-2	5,8	-13,7	7,8								
-4	-2	-9,8	9,7	-13,7	7,7		0	-2	5,8	-13,7	17,5	-11,7								
-4	-4	-0,1	-4	7,7	-4		0	0	-3,9	7,8	-11,7	7,8								
						•							•							
						-														
-0,1	-4	-0,1	-7,9	-4	-4		-3,9	7,8	-3,9	3,9	0	0		17,5	1,9	-2		7,8	5,8	0
-4	17,5	-2	1,9	-2	-4		7,8	-13,7	5,8	-5,9	-2	0		1,9	-2	9,7		5,8	5,8	7,8
-0,1	-2	-2	-2	-9,8	-0,1		-3,9	5,8	-2	-2	5,8	-3,9		-2	9,7	7,7		0	7,8	17,5
-7,9	1,9	-2	-2	9,7	-4		3,9	-5,9	-2	5,8	-13,7	7,8					•			
-4	-2	-9,8	9,7	-13,7	7,7		0	-2	5,8	-13,7	17,5	-11,7								
-4	-4	-0,1	-4	7,7	-4		0	0	-3,9	7,8	-11,7	7,8								

• Max-pooling layers are usually placed between stacked convolutional layers.

Stacking convolution, pooling, dense layers



- Max-pooling gradually reduces the resolution at higher layers, allowing us to use more channels (for the same total number of trainable parameters).
- It also helps increase more quickly the receptive field.

		Inp	out Cha	nnel 1						Inpu	rt Chanr	nel 2					T	wo Two	-chann	el Kern	els		Fea	ture Ma	p of Ke	rnel 1 ("+")			Featu	re Map	of Ker	nel 2 ("	'X")			N	Max-Po	oling (2	1,2) with	1 Stride	2,2)	
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-1	1	-1	1	-1	1	-1																					
0	-1	1	-1	-1	-1	-1	0	0	-1	0,9	-1	-1	-1	-1	0	1	1	1	&	1	1	1	-0,1	-4	-0,1	-7,9	-4	-4	[-3,9	7,8	-3,9	3,9	0	0	- I <u>-</u>	17,5	1,9	-2	Г	7,8	5,8	0
0	1	1	1	-1	-1	-1	0	0	0,9	0,9	0,9	-1	-1	-1	0	-1	1	-1		-1	1	-1	-4	17,5	-2	1,9	-2	-4		7,8	-13,7	5,8	-5,9	-2	0		1,9	-2	9,7		5,8	5,8	7,8
D	-1	1	-1	-1	-1	-1	0	0	-1	0,9	-1	-1	-1	-1	0								-0,1	-2	-2	-2	-9,8	-0,1		-3,9	5,8	-2	-2	5,8	-3,9		-2	9,7	7,7		0	7,8	17,5
D	-1	-1	-1	0,9	-1	0,9	0	0	-1	-1	-1	1	-1	1	0	1	-1	1		1	-1	1	-7,9	1,9	-2	-2	9,7	-4		3,9	-5,9	-2	5,8	-13,7	7,8	_							
D	-1	-1	-1	-1	0,9	-1	0	0	-1	-1	-1	-1	1	-1	0	-1	1	-1	&	-1	1	-1	-4	-2	-9,8	9,7	-13,7	7,7		0	-2	5,8	-13,7	17,5	-11,7								
D	-1	-1	-1	0,9	-1	0,9	0	0	-1	-1	-1	1	-1	1	0	1	-1	1		1	-1	1	-4	-4	-0,1	-4	7,7	-4		0	0	-3,9	7,8	-11,7	7,8								
o '	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0								-						•				_		_								

- Each feature of the max-pooled feature maps is derived from (is "looking at") 4 features of the pre-pooled feature maps, and 16 features of the input.
- By stacking convolution and pooling layers, we can get features that are increasingly aware of larger parts of the input (larger "receptive field").

ustrated

in Dive

Stacking convolution, pooling, dense layers



- The features of the top feature maps are concatenated to a single vector and passed to a dense (fully connected) layer or an MLP (with hidden layers).
 - To recognize the digit (0-9) in an image, the dense layer (or output layer of the MLP) would have 10 neurons with softmax, and we would use cross-entropy loss.
 - To output the **coordinates of the eyes** in images (or video frames) of faces, the **dense layer** (or output layer of the MLP) could have **4 neurons** (x1, y1, x2, y2) with no activation function, and we could use the **mean squared error** as loss. (But better, more advanced models can be used...)
 - The training examples would be digit or face images (or video frames) annotated with the correct responses (digits or coordinates of the eyes).
- In practice we would also include **dropout** layers and **residuals**.

nitecture

ustrated

in Dive

What do the layers learn?



- The kernels of lower layers tend to detect low-level features (e.g., edges of different directions). The kernels of higher layers tend to detect higher-level features (e.g., eyes, ears).
- Pre-trained kernels of lower levels can be useful in many different tasks.
 Figure from the recommended book "Deep Learning with Python" by F. Chollet, Manning Publications, 1st edition. Also covers Keras. Optionally consult Chapter 5 (Deep Learning for Computer Vision) for ways to visualize what CNN layers learn. <u>https://www.manning.com/books/deep-learning-with-python</u>

https://www.manning.com/books/deep-learning-with-python-second-edition

Re-using pretrained layers



- In practice, we start with a CNN pre-trained on a very large dataset.
 - o Often ImageNet, 1.4 million images, 1,000 classes (e.g., dogs, cats).
- We replace the top layers with a task-specific classification/regression layer.
 - We train the task-specific layer on task-specific data, keeping the pre-trained convolutional layers frozen (no weight updates in the frozen layers).
 - We may then **gradually unfreeze some of the convolutional layers too** (weight updates in both the task-specific layers and the unfrozen convolutional layers).

Figure from the recommended book "Deep Learning with Python" by F. Chollet, Manning Publications, 1st edition. Also covers Keras. <u>https://www.manning.com/books/deep-learning-</u> <u>with-python</u> <u>https://www.manning.com/books/deep-learning-with-python-second-edition</u>



Re-using pretrained layers

Figure from the recommended book "Deep Learning with Python" by F. Chollet, Manning Publications, 1st edition. Also covers Keras. <u>https://www.manning.com/books/deeplearning-with-python</u> <u>https://www.manning.com/books/deeplearning-with-python-second-edition</u>

Figure 5.19 Fine-tuning the last convolutional block of the VGG16 network

Data augmentation



Figure 5.11 Generation of cat pictures via random data augmentation

- We can **increase the number of task-specific training examples** by adding artificial training examples.
 - For example, we can **rotate**, **squeeze**, **flip** etc. the task-specific **training images**.
 - **Big improvements** usually.

Figure from the recommended book "Deep Learning with Python" by F. Chollet, Manning Publications, 1st edition. Also covers data augmentation in Keras. <u>https://www.manning.com/books/deep-learning-with-python</u> <u>https://www.manning.com/books/deep-learning-with-python-second-edition</u>

Object detection



- Find and classify **bounding boxes** of particular **object types**.
 - E.g., people, cars, horses, buses. We no longer classify entire images.

Images from "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", by Ren et al., NeurIPS 2015. <u>https://arxiv.org/abs/1506.01497</u>

Object detection with Faster R-CNN



- The **image** has **already** been turned into a **feature map** (possibly with several channels per cell in the left figure) using a **CNN image encoder**.
- **Region Proposal Network (RPN)**: We slide a convolutional window (e.g., 3×3 with 256 filters, shown in red) on the feature map.
- For each widow placement of RPN, we consider k (e.g., 9) anchor boxes (boxes of different fixed shapes/sizes, shown blue) of the *input image* whose centers correspond to the center of the window placement.

Images from "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", by Ren et al., NeurIPS 2015. <u>https://arxiv.org/abs/1506.01497</u>

Object detection with Faster R-CNN (cont'ed)



- The values (e.g., 256 channels) generated by each window placement of RPN (red) go through two different dense layers (or 1x1 convolutions).
 - The **first dense** layer ("cls", with sigmoids) predicts the "**objectness**" score of **each anchor box** (contains object or not, *k* **anchor boxes**, 2*k* scores).
 - The second dense layer produces the coordinates of the bounding box inside the anchor box (4k coordinates in total, see paper for details).
 - Among **overlapping predicted bounding boxes** (from different anchor boxes), we (roughly) keep the one with the **highest objectness score**.

Images from "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", by Ren et al., NeurIPS 2015. <u>https://arxiv.org/abs/1506.01497</u>

Object detection with Faster R-CNN (cont'ed)



- The resulting **bounding boxes** ("proposals", "Regions of Interest", ROIs) produced by the **RPN** are then **classified** (e.g., as person, bus, bicycle).
 - Using a **CNN-based classifier** and the **same convolution layers** as the RPN.
 - The **RPN's proposals** tell the classifier **where to look** in the original image.
 - The **sizes/shapes** of all proposals become **equal** by changing their resolutions (up/down-sampling, "ROI pooling").
- The **RPN** and **classifier** can be **trained jointly**, while also **fine-tuning** the **pre-trained convolutional layers**.

Images from "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", by Ren et al., NeurIPS 2015. <u>https://arxiv.org/abs/1506.01497</u>



(a) General



⁽b) Biomedical

Image captioning

Optional material

Possible applications:

- Image retrieval via captions.
- Eyesight problems.
- Drafting medical reports.

Figure 1: Example of a caption produced by the model of Vinyals et al. (2017) for a non-biomedical image (1a) and an example of a PEIR Radiology image with its associated caption (1b).

From I. Pavlopoulos, V. Kougia, I. Androutsopoulos, "A Survey on Biomedical Image Captioning". <u>https://www.aclweb.org/anthology/W19-1803/</u>

Biomedical image to text generation



Optional

material

- CNN converts each patch to a vector, producing "visual features".
- MLP ("MLC") predicts tags (classes) given the visual features.
- The word embeddings of the tags are "semantic features".
- Sentence-level LSTM produces sentence embeddings ("topics").
 - A stop control (classifier) decides when to stop producing sentences.
 - At each time-step, attention over visual and semantic features.
- For each sentence embedding, word-level LSTM produces words.
 B. Jing, P. Xie, E.P. Xing, "On the Automatic Generation of Medical Imaging Reports", ACL 2018 (<u>http://www.aclweb.org/anthology/P18-1240</u>).

NLP with CNNs and Transformers

- CNNs can also be applied to texts
 - Viewed as **1D images**. Each **"pixel" is a word**. The **channels of the input** 1D image are the **dimensions of the word embeddings**.
 - To be discussed in the **following lectures**.
 - Faster than RNNs, but usually worse results.
- Pre-trained layers recently led to big improvements in NLP.
 - Mostly using **Transformers**, a type of neural nets to be discussed in the **following lectures**. Used in **BERT**, **ChatGPT**, ...
- **Transformers** are starting to be used in **Computer Vision** too.



Figure from J. Alammar's "The Illustrated BERT, ELMo, and co." (<u>http://jalammar.github.io/illustrated-bert/</u>).

Recommended reading

- F. Chollet, *Deep Learning in Python*, Manning Publications, 1st edition, 2017, Chapter 5.
 - The 1st edition is freely available, suffices for this course: <u>https://www.manning.com/books/deep-learning-with-python</u>
 - 2nd edition also available, requires payment, recommended: <u>https://www.manning.com/books/deep-learning-with-python-second-edition</u>
- A. Zhang et al., *Dive into Deep Learning*, Chapter 6.
 Freely available at: <u>https://d21.ai/</u>
- Y. Goldberg, *Neural Network Models for Natural Language Processing*, Morgan & Claypool Publishers, 2017.
 - Chapter 13 discusses applying CNNs to text.
- See also the recommended reading and resources of Part B5 of this course.





Βιβλιογραφία – συνέχεια

- Αν έχετε από το μάθημα της ΤΝ το βιβλίο των Russel & Norvig «Τεχνητή Νοημοσύνη – Μια σύγχρονη προσέγγιση», 4^η έκδοση, Κλειδάριθμος, 2021, μπορείτε να συμβουλευτείτε τα κεφάλαια 21 και 25.
 - Κυρίως τις ενότητες 21.3, 25.4, 25.5, 25.7.2. Η περιγραφή του Faster R-CNN δεν είναι πολύ καλή όμως.

