# **Convolutional Networks**

Hakan Bilen

Machine Learning Practical — MLP Lecture 7 29 October / 5 November 2019

# Recap: Fully-connected network for MNIST

On MNIST, we can get about 2% error (or even better) using these kind of networks



(image from: Michael Nielsen, Neural Networks and Deep Learning,

http://neuralnetworksanddeeplearning.com/chap6.html)

#### How about more complex images?



image credit: COCO dataset

### Generic object classification



- Large variations in appearance and pose
- Objects are not always centered
- Objects are in different scales
- High resolution images
- Background and clutter
- Occlusions

#### Fully-connected network in high dimension



For an input size of  $200 \times 200$ , a deep network with one hidden layer (1000 hidden units) has

- # input units is 40,000
- # connections is 40,000,000
- # parameters is 40,000,000

image credit: Lecun & Ranzato

### Fully-connected network in high dimension



image credit: Lecun & Ranzato

For an input size of 200  $\times$  200, a deep network with one hidden layer (1000 hidden units) has

- # input units is 40,000
- # connections is 40,000,000
- # parameters is 40,000,000

Observations:

- Too many parameters to learn!
- Spatial (2-D) structure of input image is ignored
- No easy way of learning the same feature at different coordinates

#### **Convolution operation**

• Convolution is an operation on two functions of a real- valued argument

$$s(t) = \int x(a)w(t-a)da$$

• Denoted with asterisk (\*)

$$s(t) = (x * w)(t)$$

• Discrete convolution

$$s(t) = \sum_{a=-\infty}^{a=\infty} x(a)w(t-a)da$$

• Convolution in multiple dimensions

$$s(i,j) = (\mathbf{x} * \mathbf{w})(i,j) = \sum_{k=0}^{m-1} \sum_{l=0}^{n-1} x_{k,l} w_{i-k,j-l}$$

#### **Convolution operation**

• Multiple dimensions

$$s(i,j) = (\mathbf{x} * \mathbf{w})(i,j) = \sum_{k=0}^{m-1} \sum_{l=0}^{n-1} x_{k,l} w_{i-k,j-l}$$

• Convolution is commutative

$$s(i,j) = (\mathbf{w} * \mathbf{x})(i,j) = \sum_{k=0}^{m-1} \sum_{l=0}^{n-1} x_{i-k,j-l} w_{k,l}$$

### **Convolution operation**

• Multiple dimensions

$$s(i,j) = (\mathbf{x} * \mathbf{w})(i,j) = \sum_{k=0}^{m-1} \sum_{l=0}^{n-1} x_{k,l} w_{i-k,j-l}$$

• Convolution is commutative

$$s(i,j) = (\mathbf{w} * \mathbf{x})(i,j) = \sum_{k=0}^{m-1} \sum_{l=0}^{n-1} x_{i-k,j-l} w_{k,l}$$

• Convolutions in deep learning libraries are actually cross-correlation ( $\otimes$ )

$$s(i,j) = (\mathbf{w} \otimes \mathbf{x})(i,j) = \sum_{k=0}^{m-1} \sum_{l=0}^{n-1} x_{i+k,j+l} w_{k,l}$$

If we "flip" (reflect horizontally and vertically) w (in cross-correlation) then we obtain w (in convolution)

- In machine learning the network learns the kernel appropriate to its orientation so if convolution is implemented with a flipped kernel, it will learn that it is a flipped implementation
- So it is OK to use an efficient (flipped) implementation of convolution for convolutional layers
- In convolutional networks, we also use a bias term:

$$s(i,j) = (\mathbf{w} \otimes \mathbf{x})(i,j) + b = \sum_{k=0}^{m-1} \sum_{l=0}^{n-1} x_{i+k,j+l} w_{k,l} + b$$

- $\mathbf{h} = f(\underbrace{\mathbf{x} \otimes \mathbf{w}}_{\mathbf{z}})$  where f is an activation function
- h is called feature map







$$\begin{bmatrix} w_{00} & w_{01} & w_{02} \\ w_{10} & w_{11} & w_{12} \\ w_{20} & w_{21} & w_{22} \end{bmatrix}$$

$$\boxed{z_{01} = \sum_{k=0}^{2} \sum_{l=0}^{2} x_{k,1+l} w_{kl} + b}$$

$$\begin{bmatrix} z_{00} & z_{01} \cdot z_{02} & z_{03} & z_{04} \\ z_{10} & z_{11} & z_{12} & z_{13} & z_{14} \\ z_{20} & z_{21} & z_{22} & z_{23} & z_{24} \\ z_{30} & z_{31} & z_{32} & z_{33} & z_{34} \\ z_{40} & z_{41} & z_{42} & z_{43} & z_{44} \end{bmatrix} = \begin{bmatrix} w_{00} & w_{01} & w_{02} \\ w_{10} & w_{11} & w_{12} \\ x_{10} & x_{01} & x_{02} & x_{03} \\ x_{10} & x_{11} & x_{12} & x_{13} \\ x_{10} & x_{11} & x_{12} & x_{13} \\ x_{20} & x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} \\ x_{30} & x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} \\ x_{40} & x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} \\ x_{50} & x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} \\ x_{60} & x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} \end{bmatrix}$$

$$\begin{pmatrix} w_{00} & w_{01} & w_{02} \\ w_{10} & w_{11} & w_{12} \\ w_{20} & w_{21} & w_{22} \end{pmatrix}$$

$$= \sum_{k=0}^{2} \sum_{l=0}^{2} x_{k,4+l} w_{kl} + b$$

$$\begin{pmatrix} z_{00} & z_{01} & z_{02} & z_{03} & z_{04} \\ z_{10} & z_{11} & z_{12} & z_{13} & z_{14} \\ z_{20} & z_{21} & z_{22} & z_{23} & z_{24} \\ z_{30} & z_{31} & z_{32} & z_{33} & z_{34} \\ z_{40} & z_{41} & z_{42} & z_{43} & z_{44} \end{pmatrix} = \begin{bmatrix} w_{00} & w_{01} & w_{02} \\ w_{10} & w_{11} & w_{12} \\ x_{00} & x_{01} & x_{02} & x_{03} & x_{04} & x_{05} & x_{06} \\ x_{10} & x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} \\ x_{20} & x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} \\ x_{30} & x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} \\ x_{40} & x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} \\ x_{50} & x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} \\ x_{60} & x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} \end{bmatrix}$$

$$\begin{bmatrix} w_{00} & w_{01} & w_{02} \\ w_{10} & w_{11} & w_{12} \\ w_{20} & w_{21} & w_{22} \end{bmatrix}$$

$$\begin{bmatrix} z_{10} = \sum_{k=0}^{2} \sum_{l=0}^{2} x_{1+k,l} w_{kl} + b \\ \hline \\ z_{10} \neq z_{11} & z_{12} & z_{13} & z_{14} \\ z_{20} & z_{21} & z_{22} & z_{23} & z_{24} \\ z_{30} & z_{31} & z_{32} & z_{33} & z_{34} \\ z_{40} & z_{41} & z_{42} & z_{43} & z_{44} \end{bmatrix} = \begin{bmatrix} w_{00} & x_{01} & x_{02} & x_{03} & x_{04} & x_{05} & x_{06} \\ \hline \\ x_{00} & x_{01} & x_{02} & x_{03} & x_{04} & x_{05} & x_{06} \\ \hline \\ x_{10} & x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} \\ \hline \\ x_{20} & x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} \\ \hline \\ x_{30} & x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} \\ \hline \\ x_{40} & x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} \\ \hline \\ x_{50} & x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} \\ \hline \\ x_{60} & x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} \end{bmatrix}$$

#### Number of ...

- input units are  $7 \times 7$
- weights is  $3 \times 3 + 1$  (9 for kernel + 1 for bias)
- hidden units is  $5\times 5$
- connections is  $5\times5\times3\times3$

#### Number of ...

- input units are  $7 \times 7$
- weights is  $3 \times 3 + 1$  (9 for kernel + 1 for bias)
- hidden units is  $5\times 5$
- connections is  $5\times5\times3\times3$

#### Properties

- Weights (kernel) are shared across all hidden units
- Spatial correspondence between pixels and hidden units ("2D matrix of hidden units" = "feature map")
- Translation invariance: extract same features irrespective of where an image patch is located in the input

# If **x** is $7 \times 7$ and **w** is $3 \times 3$ , then the feature map **h** is $5 \times 5$ . If **x** is $7 \times 7$ and **w** is $5 \times 5$ , what will be the dimensions of **h** ? (a) $3 \times 3$ , (b) $5 \times 5$ , (c) $7 \times 7$









If **x** is  $7 \times 7$  and **w** is  $5 \times 5$ , what will be the dimensions of **h** ?



**Q**. If **x** is  $q \times r$  and **w** is  $m \times n$  dimensional, what is the output dimensionality? **A**.  $(q - m + 1) \times (r - n + 1)$  Padding

#### Q. What can we do to get $7 \times 7$ feature map size?



 $Z_{00}$  $Z_{01}$  $Z_{02}$ Z03 Z04  $Z_{10}$  $Z_{11}$  $z_{12}$   $z_{13}$   $z_{14}$ Z23 Z20  $Z_{21}$  $Z_{22}$ Z30 Z31 Z32 Z33 Z41 Z42 Z43  $Z_{40}$  $Z_{51}$ Z52 Z53  $Z_{50}$ Z60 Z61 Z62 Z63

Padding

 $Z_{00}$ 

 $Z_{10}$  $Z_{11}$ 

 $Z_{20}$ 

Z30

 $Z_{40}$ 

 $Z_{50}$ 

 $Z_{60}$ Z61

Z01

 $Z_{21}$ 

Z31

Z41

 $Z_{51}$ 

 $Z_{02}$ 

 $Z_{22}$ 

Z32

Z42

 $Z_{52}$ 

Z62

#### Q. What can we do to get $7 \times 7$ feature map size?



Padding





![](_page_29_Figure_1.jpeg)

![](_page_30_Figure_1.jpeg)

![](_page_31_Figure_1.jpeg)

![](_page_32_Figure_1.jpeg)

![](_page_33_Figure_1.jpeg)

# What is the feature map size when **x** is $q \times r$ and **w** is $m \times n$ dimensional, for padding p and stride s?

### **Receptive field**

- Biology: The **receptive field** of an individual sensory neuron is the particular region of the sensory space,
- Convolutional networks: The region in the input space that a hidden unit is looking at.

$h_{00}$	$h_{01}$	$h_{02}$	h <sub>03</sub>	h <sub>04</sub>
$h_{10}$	$h_{11}$	$h_{12}$	$h_{13}$	$h_{14}$
$h_{20}$	$h_{21}$	h <sub>22</sub>	h <sub>23</sub>	h <sub>24</sub>
h <sub>30</sub>	$h_{31}$	h <sub>32</sub>	h <sub>33</sub>	h <sub>34</sub>
h <sub>40</sub>	$h_{41}$	h <sub>42</sub>	h <sub>43</sub>	h <sub>44</sub>

<i>x</i> <sub>00</sub>	<i>x</i> <sub>01</sub>	<i>x</i> <sub>02</sub>	<i>x</i> 03	<i>x</i> <sub>04</sub>	<i>x</i> 05	<i>x</i> 06
<i>x</i> <sub>10</sub>	$x_{11}$	<i>x</i> <sub>12</sub>	<i>x</i> <sub>13</sub>	<i>x</i> <sub>14</sub>	$x_{15}$	<i>x</i> <sub>16</sub>
<i>x</i> <sub>20</sub>	<i>x</i> <sub>21</sub>	<i>x</i> <sub>22</sub>	<i>x</i> <sub>23</sub>	<i>x</i> <sub>24</sub>	<i>x</i> <sub>25</sub>	<i>x</i> <sub>26</sub>
<i>x</i> <sub>30</sub>	<i>x</i> <sub>31</sub>	<i>x</i> <sub>32</sub>	<i>X</i> 33	<i>X</i> 34	<i>X</i> 35	<i>x</i> 36
<i>x</i> <sub>40</sub>	<i>x</i> <sub>41</sub>	<i>x</i> <sub>42</sub>	<i>x</i> <sub>43</sub>	<i>x</i> <sub>44</sub>	<i>x</i> 45	<i>x</i> <sub>46</sub>
<i>x</i> 50	<i>x</i> 51	<i>x</i> <sub>52</sub>	<i>X</i> 53	<i>X</i> 54	<i>X</i> 55	<i>X</i> 56
<i>x</i> 60	<i>x</i> <sub>61</sub>	<i>x</i> <sub>62</sub>	<i>x</i> 63	<i>x</i> 64	<i>x</i> 65	<i>x</i> 66

### **Receptive field**

- Biology: The **receptive field** of an individual sensory neuron is the particular region of the sensory space,
- Convolutional networks: The region in the input space that a hidden unit is looking at.

								( x <sub>00</sub>	<i>x</i> 01	<i>x</i> <sub>02</sub>	<i>x</i> 03	<i>x</i> 04	X05	<i>x</i> 06
/	h <sub>00</sub>	$h_{01}$	h <sub>02</sub>	h <sub>03</sub>	h <sub>04</sub>			<i>x</i> <sub>10</sub>	<i>x</i> <sub>11</sub>	<i>x</i> <sub>12</sub>	<i>x</i> <sub>13</sub>	<i>x</i> <sub>14</sub>	<i>x</i> <sub>15</sub>	<i>x</i> <sub>16</sub>
	$h_{10}$	h <sub>11</sub>	$h_{12}$	$h_{13}$	$h_{14}$			<i>x</i> <sub>20</sub>	<i>x</i> <sub>21</sub>	<i>x</i> <sub>22</sub>	<i>x</i> 23	<i>x</i> 24	<i>x</i> 25	<i>x</i> <sub>26</sub>
	h <sub>20</sub>	$h_{21}$	h <sub>22</sub>	h <sub>23</sub>	h <sub>24</sub>			X30	<i>x</i> 31	<i>x</i> <sub>32</sub>	<i>x</i> 33	<i>x</i> 34	X35	<i>x</i> 36
	h <sub>30</sub>	$h_{31}$	h <sub>32</sub>	h <sub>33</sub>	h <sub>34</sub>			x <sub>40</sub>	<i>x</i> <sub>41</sub>	<i>x</i> <sub>42</sub>	<i>x</i> 43	<i>x</i> 44	<i>x</i> 45	x <sub>46</sub>
	h <sub>40</sub>	$h_{41}$	h <sub>42</sub>	h <sub>43</sub>	h <sub>44</sub>			<i>x</i> 50	<i>x</i> 51	<i>X</i> 52	<i>X</i> 53	<i>X</i> 54	<i>X</i> 55	<i>X</i> 56
						/	(	×60	<i>x</i> <sub>61</sub>	<i>x</i> <sub>62</sub>	<i>x</i> 63	<i>x</i> <sub>64</sub>	<i>x</i> 65	<i>x</i> 66

$$\mathbf{h}^{(2)} = f(\mathbf{w}^2 \otimes f(\mathbf{w}^1 \otimes \mathbf{x}))$$

**Q**. What size is the receptive field of a hidden unit in second convolutional layer?

![](_page_37_Figure_5.jpeg)

$$\mathbf{h}^{(2)} = f(\mathbf{w}^2 \otimes f(\mathbf{w}^1 \otimes \mathbf{x}))$$

**Q**. What size is the receptive field of a hidden unit in second convolutional layer?

![](_page_38_Figure_5.jpeg)

$$\mathbf{h}^{(2)} = f(\mathbf{w}^2 \otimes f(\mathbf{w}^1 \otimes \mathbf{x}))$$

 ${f Q}.$  What size is the receptive field of a hidden unit in second convolutional layer?

![](_page_39_Figure_5.jpeg)

$$\mathbf{h}^{(2)} = f(\mathbf{w}^2 \otimes f(\mathbf{w}^1 \otimes \mathbf{x}))$$

 ${\bf Q}.$  What size is the receptive field of a hidden unit in second convolutional layer?

![](_page_40_Figure_5.jpeg)

$$\mathbf{h}^{(2)} = f(\mathbf{w}^2 \otimes f(\mathbf{w}^1 \otimes \mathbf{x}))$$

 ${f Q}.$  What size is the receptive field of a hidden unit in second convolutional layer?

![](_page_41_Figure_5.jpeg)

- Constrain each hidden unit *h<sub>ij</sub>* to extract the same feature by **sharing weights** across the receptive fields
- Local receptive fields with shared weights result in a **feature map** which is a map showing where the feature corresponding to the shared weight matrix (kernel) occurs in the image
- Feature map encodes **translation invariance**: extract the same features irrespective of where an image is located in the input

#### Multiple output feature maps

 $\left(\begin{array}{ccccc} h_{100} & h_{101} & \dots & h_{104} \\ h_{110} & h_{111} & \dots & h_{114} \\ \vdots & \vdots & \ddots & \vdots \\ h_{140} & h_{141} & \dots & h_{144} \end{array}\right) = \left(\begin{array}{cccc} x_{00} & x_{01} & \dots & x_{06} \\ x_{10} & x_{11} & \dots & x_{16} \\ \vdots & \vdots & \ddots & \vdots \\ x_{60} & x_{61} & \dots & x_{66} \end{array}\right) \\ \hline \otimes \left(\begin{array}{cccc} w_{100} & w_{101} & w_{102} \\ w_{110} & w_{111} & w_{112} \\ w_{120} & w_{121} & w_{122} \end{array}\right)$ 

 $h_{ijk}$ : *i*-th feature-map, *j*-th row, *k*-th column

- # input units is  $7\times7$
- # feature maps is  $F_{out}$
- # hidden units is  $\textit{F}_{out} \times (5 \times 5)$
- # of parameters is  $F_{\text{out}} \times (3 \times 3 + 1)$
- # of connections is  $\textit{F}_{out} \times (5 \times 5 \times 3 \times 3)$

#### Multiple input feature maps (or input channels)

![](_page_45_Figure_1.jpeg)

- # input channels is  $F_{in}$
- # input units is  $F_{in} \times 7 \times 7$
- # hidden units is  $5\times 5$
- # parameters is  $\textit{F}_{in} \times 3 \times 3 + 1$  for bias
- # connections is  $F_{in} \times 5 \times 5 \times 3 \times 3$
- Typically we do not tie weights across feature maps
- Local receptive fields across multiple input images

- # input channels is  $F_{in}$
- # feature maps is  $F_{out}$
- # input units is  $F_{in} \times 7 \times 7$
- # hidden units is  $\textit{F}_{out} \times 5 \times 5$
- # parameters is  $F_{\text{in}} \times F_{\text{out}} \times 3 \times 3 + F_{\text{out}}$  for bias
- # connections is  $\textit{F}_{in} \times \textit{F}_{out} \times 5 \times 5 \times 3 \times 3$

# Pooling (subsampling)

![](_page_48_Figure_1.jpeg)

# **Pooling (subsampling)**

![](_page_49_Figure_1.jpeg)

- Similar to convolution, slides over input pixels but no learnable parameters
- Has local receptive field too
- Typical stride s > 1

# Pooling

- Pooling or subsampling takes a feature map and reduces it in size e.g. by transforming a set of 2x2 regions to a single unit
- Reduces computation time and memory use
- Pooling functions
  - Max-pooling takes the maximum value of the units in the region
  - $L_p$ -pooling take the  $L_p$  norm of the units in the region:

$$h' = \left(\sum_{i \in \text{region}} h_i^p\right)^{1/p}$$

- Average- / Sum-pooling takes the average / sum value of the pool
- Information reduction pooling removes precise location information for a feature
- Apply pooling to each feature map separately

# Learning image kernels

![](_page_51_Figure_1.jpeg)

https://en.wikipedia.org/wiki/ Kernel\_(image\_processing)

- Image kernels have been hand designed and used for feature extraction in image processing (e.g. edge detection)
- Pros: No need for data and training
- Cons 1: Learning filters can be more optimal (minimising network cost function)
- Cons 2: Difficult to design filters for complex tasks (e.g. recognising a cat)
- Automating feature engineering

![](_page_52_Figure_1.jpeg)

# MNIST Results (1997)

![](_page_53_Figure_1.jpeg)

#### ImageNet Classification ("AlexNet")

Krizhevsky, Sutskever and Hinton, "ImageNet Classification with Deep Convolutional Neural Networks", NIPS'12.

 $\verb+http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf$ 

![](_page_54_Figure_3.jpeg)

Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264– 4096–4096–1000.

Model	Top-1	Top-5
Sparse coding [2]	47.1%	28.2%
SIFT + FVs [24]	45.7%	25.7%
CNN	37.5%	17.0%

#### **Hierarchical Representations**

$$\mathsf{Pixel} \to \mathsf{edge} \to \mathsf{texton} \to \mathsf{motif} \to \mathsf{part} \to \mathsf{object}$$

![](_page_55_Figure_2.jpeg)

Zeiler & Fergus, "Visualizing and Understanding Convolutional Networks", ECCV'14.

https://cs.nyu.edu/~fergus/papers/zeilerECCV2014.pdf

Slide credits: Lecun & Ranzato

- Train convolutional networks with a straightforward but careful application of backprop / SGD
- Exercise: prior to the next lecture, write down the gradients for the weights and biases of the feature maps in a convolutional network. Remember to take account of weight sharing.
- Next lecture: implementing convolutional networks: how to deal with local receptive fields and tied weights, computing the required gradients...

Convolutional networks include local receptive fields, weight sharing, and pooling leading to:

- Modelling the spatial structure
- Translation invariance
- Local feature detection

#### Reading:

Michael Nielsen, *Neural Networks and Deep Learning* (ch 6) http://neuralnetworksanddeeplearning.com/chap6.html

Yann LeCun et al, "Gradient-Based Learning Applied to Document Recognition", Proc

*IEEE*, 1998. http://dx.doi.org/10.1109/5.726791

Ian Goodfellow, Yoshua Bengio & Aaron Courville, Deep Learning (ch 9)

http://www.deeplearningbook.org/contents/convnets.html