# ECE 449/590 – OOP and Machine Learning
## Lecture 17 Convolutional Networks

Professor Jia Wang
Department of Electrical and Computer Engineering
Illinois Institute of Technology

October 24, 2022

# Outline

Convolution

Pooling

# Reading Assignment

- ▶ This lecture: Deep Learning 9
- ▶ Next lecture: Deep Learning 5, 6

# Outline

Convolution

Pooling

## Convolutional Networks

- ▶ A.k.a. convolutional neural networks, or CNNs.
- ▶ For data that has a known grid-like topology, e.g.
  - ▶ 1-D: time-series
  - ▶ 2-D: images
- ▶ Make use of convolution in at least one of the neural network layers.
  - ▶ Specialized kind of linear operation.

## The Convolution Operation

▶ Convolution: average <u>input</u> $x$ with <u>kernel</u> $w$.

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

  ▶ Typically written as $s(t) = (x * w)(t)$
  ▶ From aspects of signals and systems, $w$ is the impulse response of the associated linear and time-invariant (LTI) system.

▶ Most neural network libearies implement convolution as cross-correlation.

  ▶ E.g. 2-D: $S(i,j) = \sum_m \sum_n I(i+m, j+n)K(m,n)$
  ▶ Note that $S(i,j)$, $I(i,j)$ could be vectors and $K(m,n)$ are therefore matrices, requiring $S$, $I$, $K$ to be tensors themselves.

# 2D Convolution



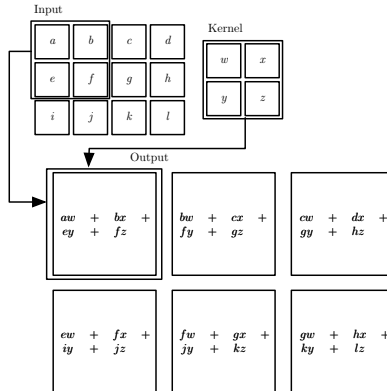Figure 9.1

(Goodfellow 2016)

▶ Only keep an output cell when all input cells are defined.
  ▶ Neural network libearies usually use <u>padding</u> to expand the input so that output could remain the same size.
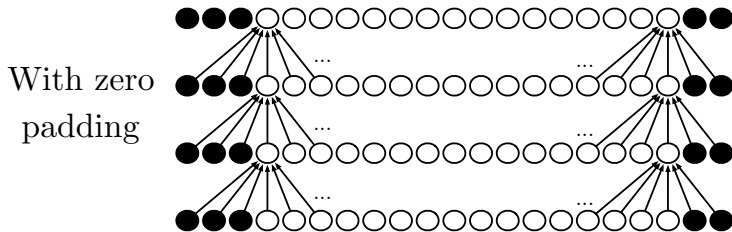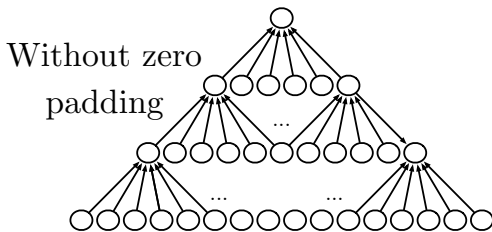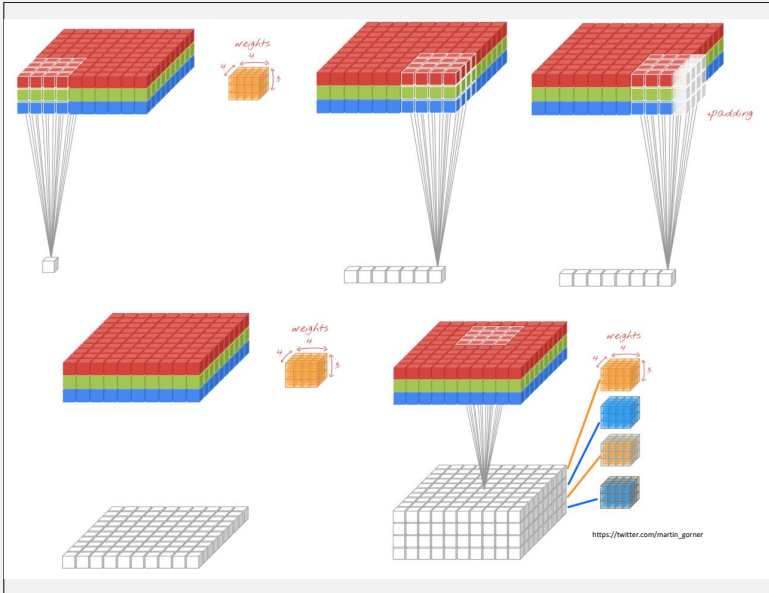
# Zero Padding Controls Size



Without zero padding

With zero padding

Figure 9.13

(Goodfellow 2016)

# Convolution with Input and Output Channels

▶ In practice, convolutions are usually applied to a group of input channels to generate a group of output channels.

▶ For a 2D image,

▶ With color channels, the input to the first layer is a 3D tensor (W, H, input channels).

▶ A kernel as a 3D tensor will only generate a 2D output.

▶ Instead, we use a kernel that is a 4D tensor so the output is a 3D tensor (W, H, output channels)

▶ Layers using 4D kernels can now be easily composed together.
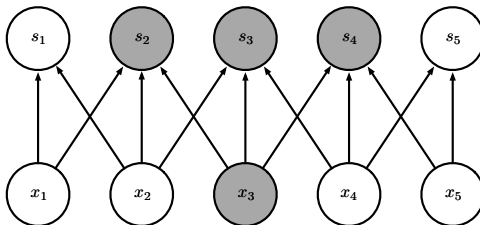
https://twitter.com/martin_gorner

# Why convolution?

▶ Sparse interactions: by using a small kernel, interaction between inputs and outputs are limited.

▶ Parameter sharing: same kernel is applied many times.

▶ Equivariant representations: help to learn representations that are invariant to input locations.

▶ Compared to the (fully-connected) linear layer,
  ▶ Able to work with inputs with variable size
  ▶ Much less capacity and storage $O(k)$ (vs $O(n^2)$).
  ▶ Less computation $O(kn)$ (vs $O(n^2)$).

# Sparse Connectivity

Sparse connections due to small convolution kernel
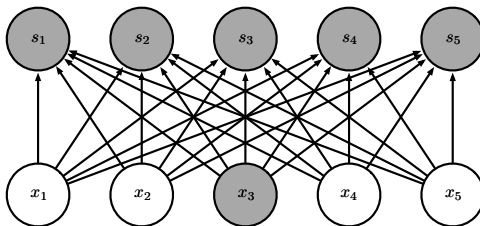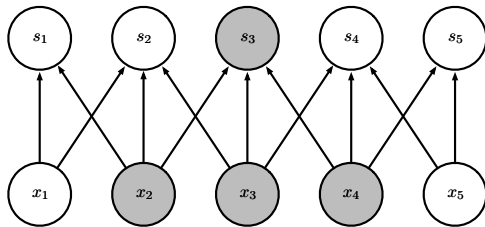


Dense connections

Figure 9.2

ECE 449/590 – Object-Oriented Programming and Machine Learning, Dept. of ECE, IIT

(Goodfellow 2016)

# Sparse Connectivity



Sparse connections due to small convolution kernel

Dense connections
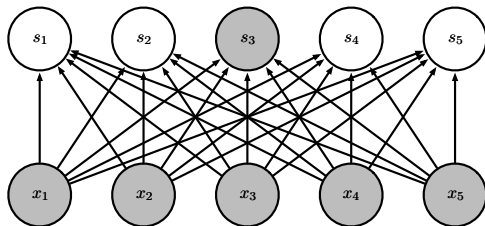
Figure 9.3

(Goodfellow 2016)

# Growing Receptive Fields

Figure 9.4

(Goodfellow 2016)
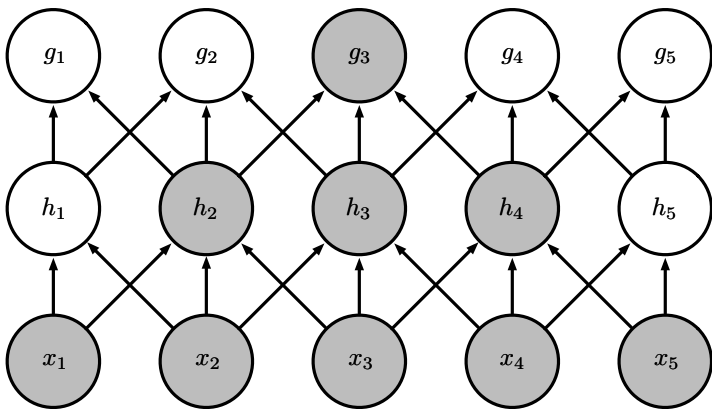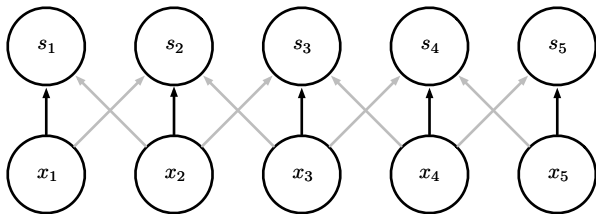
# Parameter Sharing



Convolution shares the same parameters across all spatial locations

Traditional matrix multiplication does not share any parameters

Figure 9.5

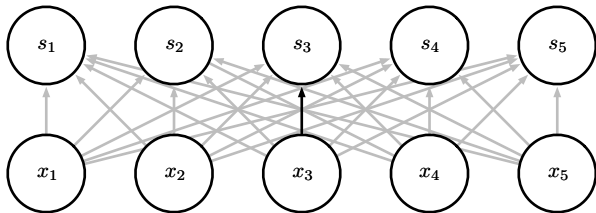(Goodfellow 2016)

# Edge Detection by Convolution



Input

Output

Kernel

| 1 | -1 |
|---|----|

Figure 9.6

(Goodfellow 2016)

# Kinds of Connectivity



Local connection: like convolution, but no sharing

Convolution

Fully connected

Figure 9.14

# Efficiency of Convolution

Input size: 320 by 280

Kernel size: 2 by 1

Output size: 319 by 280

|  | Convolution | Dense matrix | Sparse matrix |
|---|---|---|---|
| **Stored floats** | 2 | 319*280*320*280 > 8e9 | 2*319*280 = 178,640 |
| **Float muls or adds** | 319*280*3 = 267,960 | > 16e9 | Same as convolution (267,960) |

(Goodfellow 2016)

# Gabor Functions



Figure 9.18

(Goodfellow 2016)

▶ Primary visual cortex (V1), where our brain begins to process visual input, can be described by Gabor functions.

# Gabor-like Learned Kernels



Figure 9.19

ECE 449/590 – Object-Oriented Programming and Machine Learning, Dept. of ECE, IIT

(Goodfellow 2016)

# Outline

ECE 449/590 – Object-Oriented Programming and Machine Learning, Dept. of ECE, IIT

# Convolutional Network Components



Figure 9.7

(Goodfellow 2016)

# Pooling

- ▶ Summarize statistic of the nearby outputs.
    - ▶ Max pooling
    - ▶ Average pooling
    - ▶ Help to learn representations that are invariant to small translations, e.g. noises on where an eye is located.
- ▶ Pooling may be applied beyond convolution to learn invariant from a group of inputs.
    - ▶ E.g. to learn representations invariant to rotation and scaling from different rotations and scalings of input.
- ▶ Usually reduce input sizes to speed up following layers.

# Max Pooling and Invariance to Translation



POOLING STAGE

... 1. 1. 1. 0.2 ...

... 0.1 1. 0.2 0.1 ...

DETECTOR STAGE

POOLING STAGE

... 0.3 1. 1. 1. ...

... 0.3 0.1 1. 0.2 ...

DETECTOR STAGE

Figure 9.8

# Cross-Channel Pooling and Invariance to Learned Transformations

(Goodfellow 2016)

# Pooling with Downsampling



ECE 449/590 – Object-Oriented Programming and Machine Learning, Dept. of ECE, IIT

Figure 9.10

(Goodfellow 2016)

# Example Classification Architectures



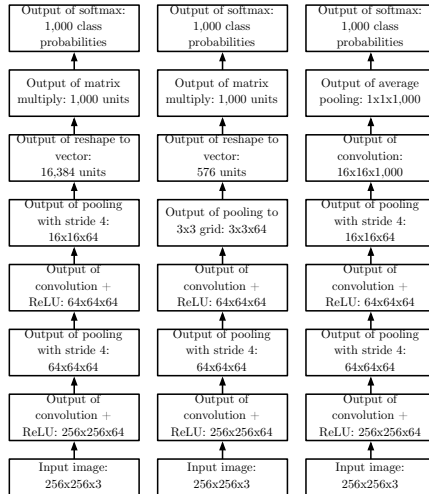| Output of softmax:<br>1,000 class<br>probabilities | Output of softmax:<br>1,000 class<br>probabilities | Output of softmax:<br>1,000 class<br>probabilities |
|---|---|---|
| Output of matrix<br>multiply: 1,000 units | Output of matrix<br>multiply: 1,000 units | Output of average<br>pooling: 1x1x1,000 |
| Output of reshape to<br>vector:<br>16,384 units | Output of reshape to<br>vector:<br>576 units | Output of<br>convolution:<br>16x16x1,000 |
| Output of pooling<br>with stride 4:<br>16x16x64 | Output of pooling to<br>3x3 grid: 3x3x64 | Output of pooling<br>with stride 4:<br>16x16x64 |
| Output of<br>convolution +<br>ReLU: 64x64x64 | Output of<br>convolution +<br>ReLU: 64x64x64 | Output of<br>convolution +<br>ReLU: 64x64x64 |
| Output of pooling<br>with stride 4:<br>64x64x64 | Output of pooling<br>with stride 4:<br>64x64x64 | Output of pooling<br>with stride 4:<br>64x64x64 |
| Output of<br>convolution +<br>ReLU: 256x256x64 | Output of<br>convolution +<br>ReLU: 256x256x64 | Output of<br>convolution +<br>ReLU: 256x256x64 |
| Input image:<br>256x256x3 | Input image:<br>256x256x3 | Input image:<br>256x256x3 |

Figure 9.11

(Goodfellow 2016)

# Summary

- CNNs apply the same kernel over many parts of the input.
- Pooling help to summarize statistic of the nearby outputs.