

# CS-UY 4563: Lecture 19

## Convolutional Neural Networks

---

NYU Tandon School of Engineering, Prof. Christopher Musco

- I will be reviewing project proposals over the next few days. If you need feedback sooner, please email separately to set up a meeting, or stop by office hours.
- Written Homework 4 due **next Monday**.
- Work through `demo_convolution.ipynb`, now on course website.

# CONVOLUTIONAL FEATURE EXTRACTION

Last lecture: **Convolution** in 1, 2, and 3D is a powerful generic tool for designing natural **feature transformations** for time-series data, audio data, or images.

**X**

1	4	3	-1	2	-4	1	0	2	-1
---	---	---	----	---	----	---	---	---	----

**W**

1	2	1
---	---	---

 $\longrightarrow$

**X**

1	4	3	-1	2	-4	1	0	2	-1
---	---	---	----	---	----	---	---	---	----

**W**

1	2	1
---	---	---

 $\longrightarrow$

**X**

1	4	3	-1	2	-4	1	0	2	-1
---	---	---	----	---	----	---	---	---	----

**W**

1	2	1
---	---	---

 $\longrightarrow$

12	9	3	-1	-5	-1	3	3
----	---	---	----	----	----	---	---

$$\mathbf{u} = \mathbf{x} * \mathbf{w}$$

# 2D CONVOLUTION

$$w = \begin{bmatrix} 0 & 1 & 2 \\ 2 & 2 & 0 \\ 0 & 1 & 2 \end{bmatrix}$$

3 <sub>0</sub>	3 <sub>1</sub>	2 <sub>2</sub>	1	0
0 <sub>2</sub>	0 <sub>2</sub>	1 <sub>0</sub>	3	1
3 <sub>0</sub>	1 <sub>1</sub>	2 <sub>2</sub>	2	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 <sub>0</sub>	2 <sub>1</sub>	1 <sub>2</sub>	0
0	0 <sub>2</sub>	1 <sub>2</sub>	3 <sub>0</sub>	1
3	1 <sub>0</sub>	2 <sub>1</sub>	2 <sub>2</sub>	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2 <sub>0</sub>	1 <sub>1</sub>	0 <sub>2</sub>
0	0	1 <sub>2</sub>	3 <sub>2</sub>	1 <sub>0</sub>
3	1	2 <sub>0</sub>	2 <sub>1</sub>	3 <sub>2</sub>
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0 <sub>0</sub>	0 <sub>1</sub>	1 <sub>2</sub>	3	1
3 <sub>2</sub>	1 <sub>2</sub>	2 <sub>0</sub>	2	3
2 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0 <sub>0</sub>	1 <sub>1</sub>	3 <sub>2</sub>	1
3	1 <sub>2</sub>	2 <sub>2</sub>	2 <sub>0</sub>	3
2	0 <sub>0</sub>	0 <sub>1</sub>	2 <sub>2</sub>	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0	1 <sub>0</sub>	3 <sub>1</sub>	1 <sub>2</sub>
3	1	2 <sub>2</sub>	2 <sub>2</sub>	3 <sub>0</sub>
2	0	0 <sub>0</sub>	2 <sub>1</sub>	2 <sub>2</sub>
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0	1	3	1
3 <sub>0</sub>	1 <sub>1</sub>	2 <sub>2</sub>	2	3
2 <sub>2</sub>	0 <sub>2</sub>	0 <sub>0</sub>	2	2
2 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0	1	3	1
3	1 <sub>0</sub>	2 <sub>1</sub>	2 <sub>2</sub>	3
2	0 <sub>2</sub>	0 <sub>2</sub>	2 <sub>0</sub>	2
2	0 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

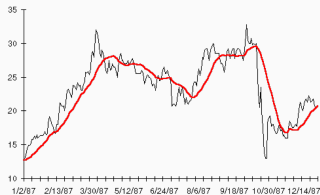
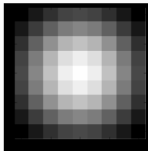
3	3	2	1	0
0	0	1	3	1
3	1	2 <sub>0</sub>	2 <sub>1</sub>	3 <sub>2</sub>
2	0	0 <sub>2</sub>	2 <sub>2</sub>	2 <sub>0</sub>
2	0	0 <sub>0</sub>	0 <sub>1</sub>	1 <sub>2</sub>

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0



# APPLICATION 1: SMOOTHING

A uniform or Gaussian filter can be used to smooth input data:

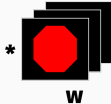


## APPLICATION 2: PATTERN MATCHING

Convolution can be used to find local patterns in images:



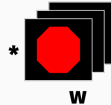
$x_1$



$W$



$x_2$



$W$



Good question from Piazza: Won't this filter yield lots of false positives?



Possible fix that doesn't require image normalization:

red channel



blue channel



green channel



 = -1

 = 0

 = 1

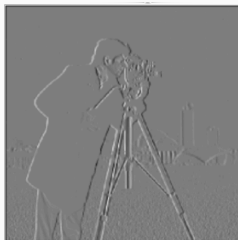
# APPLICATIONS OF CONVOLUTION

## Application 3: Edge detection.

Consider a 2D edge detection filter:

$$W_1 = \begin{bmatrix} 1 & -1 \end{bmatrix}$$

$$W_2 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$



$x * ?$



$x * ?$

# APPLICATIONS OF CONVOLUTION

Sobel filter is more commonly used:

$$W_1 = \begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$$

$$W_2 = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$$



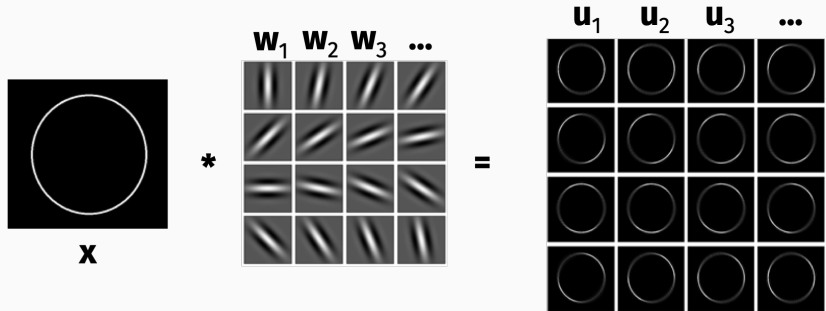
$x * ?$



$x * ?$

## DIRECTIONAL EDGE DETECTION

Can define edge detection filters for any orientation.



# EDGE DETECTION

How would edge detection as a feature extractor help you classify images of city-scapes vs. images of landscapes?





# EDGE DETECTION



$I_C$



$E_C$



$I_L$



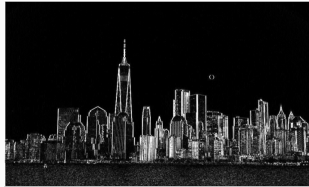
$E_L$

$$\text{mean}(I_C) = .108 \quad \text{vs.} \quad \text{mean}(I_L) = .123$$

The image with highest vertical edge response isn't the city-scape.

## EDGE DETECTION + PATTERN MATCHING

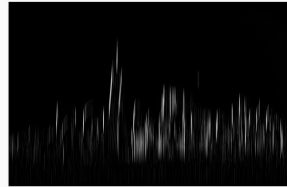
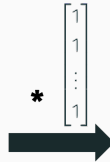
Feed edge detection result into pattern matcher that looks for long vertical lines.



$E_C$



$E_L$

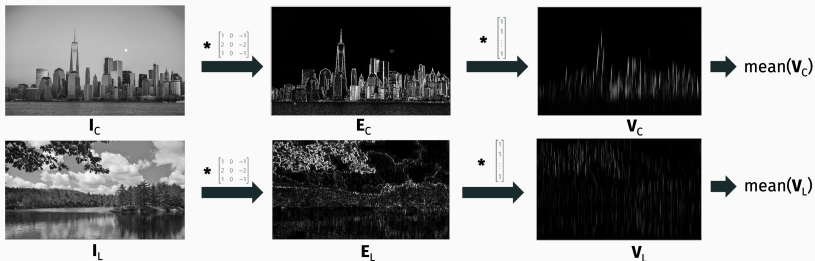


$V_C$



$V_L$

# HIERARCHICAL CONVOLUTIONAL FEATURES



$$\text{mean}(V_C) = .062 \quad \text{vs.} \quad \text{mean}(V_L) = .054$$

$$\text{mean}(V_C \cdot V_C) = .042 \quad \text{vs.} \quad \text{mean}(V_L \cdot V_L) = .018$$

The image with highest average response to (edge detector) + (vertical pattern) is the city scape.

$\text{mean}(V)$  is an extracted scalar feature which could be used for classifying cityscapes from landscapes using a linear classifier.

Hierarchical combinations of simple convolution filters are very powerful for understanding images.

In particular, edge detection seems like a critical first step.

Lots of evidence from biology.

# VISUAL SYSTEM

Light comes into the eye through the lens and is detected by an array of photosensitive cells in the **retina**.

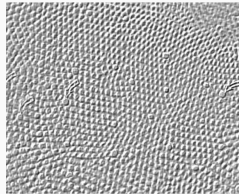
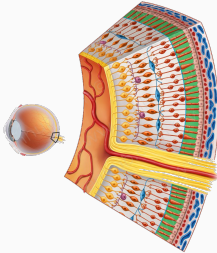
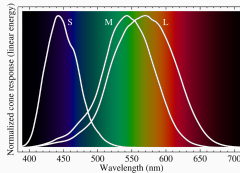


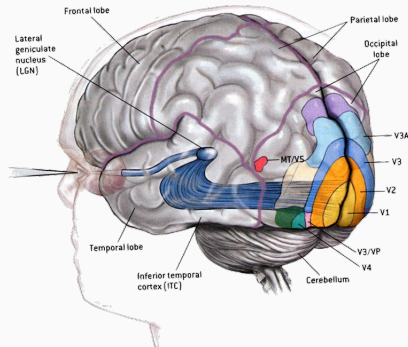
Fig. 13. Tangential section through the human fovea. Larger cones (arrows) are blue cones. From Ahnelt et al. 1987.

**Rod** cells are sensitive to all light, larger **cone** cells are sensitive to specific colors. We have three types of cones:



# VISUAL SYSTEM

Signal passes from the retina to the primary (V1) visual cortex, which has neurons that connect to higher level parts of the brain.

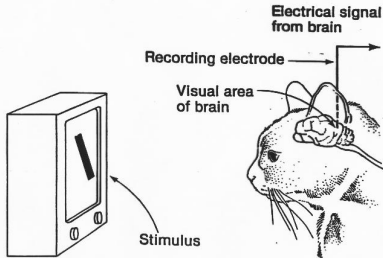


What sort of processing happens in the primary cortex?

**Lots of edge detection!**

## EDGE DETECTORS IN CATS

Huber + Wiesel, 1959: "Receptive fields of single neurones in the cat's striate cortex." Won Nobel prize in 1981.

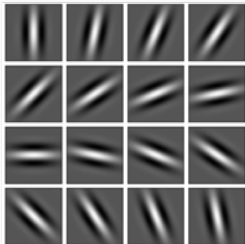


Different neurons fire when the cat is presented with stimuli at different angles. Cool video at <https://www.youtube.com/watch?v=0GxVfKJqX5E>.

"What the Frog's Eye Tells the Frog's Brain", Lettvin et al. 1959. Found explicit edge detection circuits in a frogs visual cortex.

State of the art until  $\sim 10$  years ago:

- Convolve image with edge detection filters at many different angles.
- Hand engineer features based on the responses.
- **SIFT** and **HOG** features were especially popular.

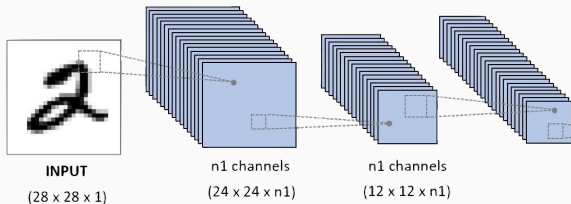




# CONVOLUTIONAL NEURAL NETWORKS

**Neural network approach:** Learn the parameters of the convolution filters based on training data.

## Convolutional Layer



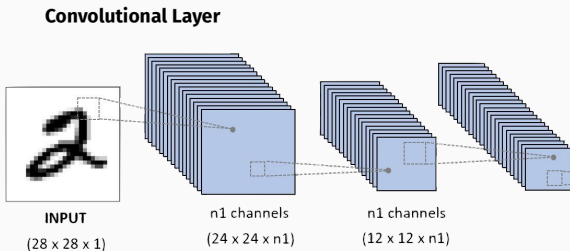
First convolutional layer involves  $n1$  convolution filters  $\mathbf{W}_1, \dots, \mathbf{W}_{n1}$ . Each is small, e.g.  $5 \times 5$ . Every entry in  $\mathbf{W}_i$  is a free parameter:  
 $\sim 25 \cdot n1$  parameters to learn.

Produces  $n1$  matrices of hidden variables: i.e. a tensor with depth  $n1$ .

# CONVOLUTIONAL NEURAL NETWORKS

A fully connected layer that extracts the same feature would require  $(28 \cdot 28 \cdot 24 \cdot 24) \cdot n1 = 451,584 \cdot n1$  parameters.

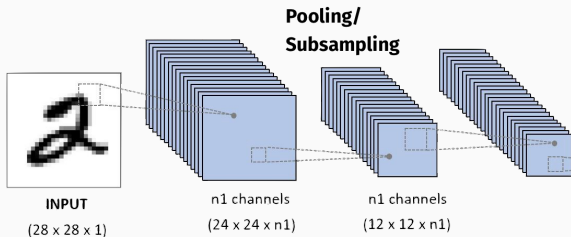
By “baking in” knowledge about what type of features matter, we greatly simplify the network.



Each of the  $n1$  outputs is typically processed with a **non-linearity**. Most commonly a Rectified Linear Unity (ReLU):  $x = \max(\bar{x}, 0)$ .

## POOLING AND DOWNSAMPLING

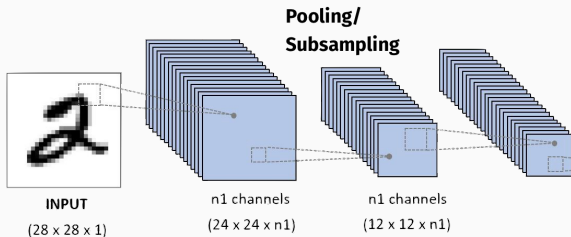
Convolution + non-linearity are typically followed by a layer which performs **pooling + down-sampling**.



Most common approach is **max-pooling**.

## POOLING AND DOWNSAMPLING

Convolution + non-linearity are typically followed by a layer which performs **pooling + down-sampling**.



Most common approach is **max-pooling**.

# POOLING AND DOWNSAMPLING

Max Pooling

29	15	28	184
0	100	70	38
12	12	7	2
12	12	45	6

2 x 2  
pool size

100	184
12	45

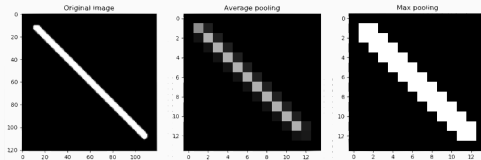
Average Pooling

31	15	28	184
0	100	70	38
12	12	7	2
12	12	45	6

2 x 2  
pool size

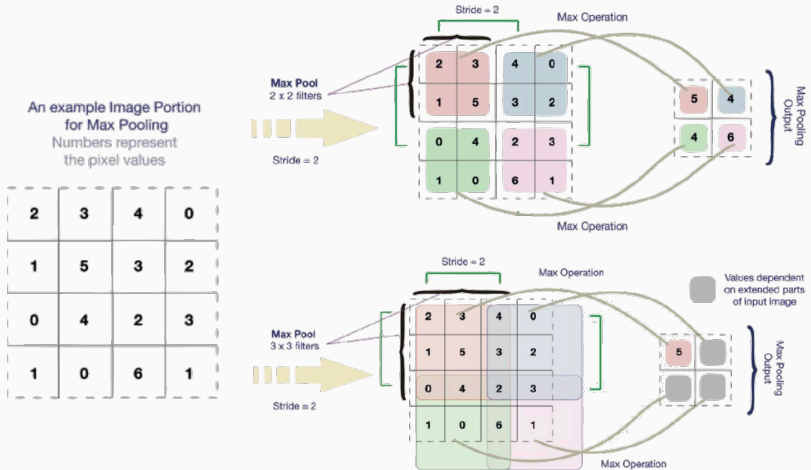
36	80
12	15

- Reduces number of variables, helps prevent over-fitting and speed up training.
- Helps “smooth” result of convolutional filters.
- Improves shift-invariance.

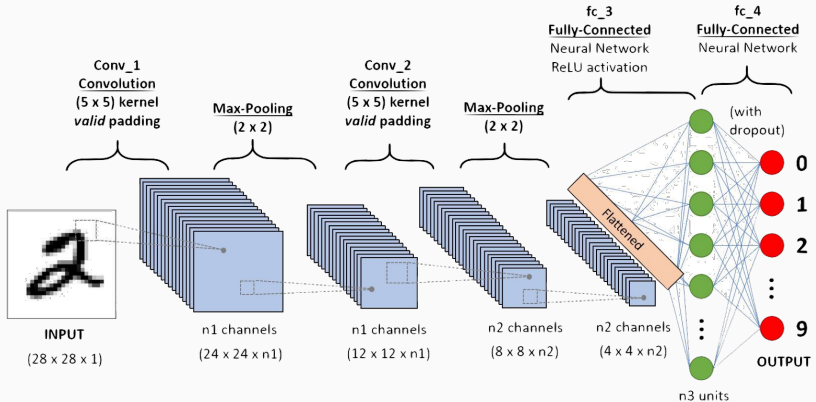


# POOLING AND DOWNSAMPLING

Many possible variations on standard 2x2 max-pooling.



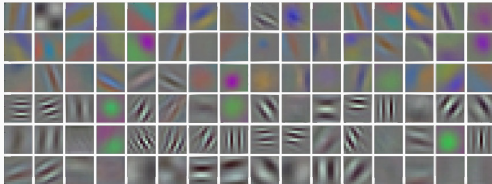
# OVERALL NETWORK ARCHITECTURE



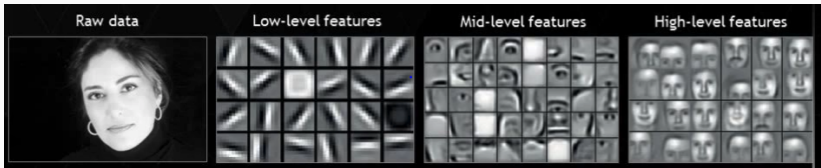
Each layer contains a 3D tensor of variables. Last few layers are standard fully connected layers.

## UNDERSTANDING LAYERS

What type of convolutional filters do we learn from gradient descent? **Lots of edge detectors in the first layer!**



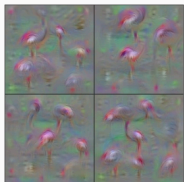
Other layers are harder to understand... but the hypothesis is that hidden variables later in the network encode for “higher level features”:



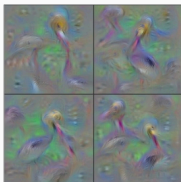


## UNDERSTANDING LAYERS

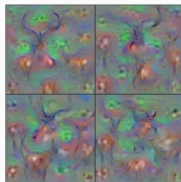
**Technique to probe later neurons:** Use optimization to find images that most strongly “activate” a given neuron deep in the network.



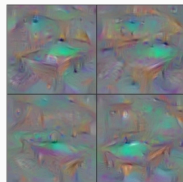
Flamingo



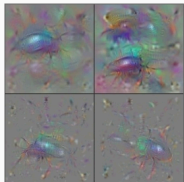
Pelican



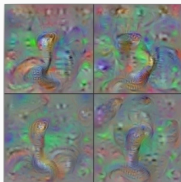
Hartebeest



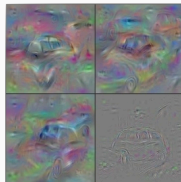
Billiard Table



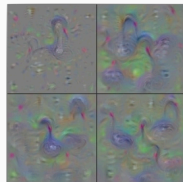
Ground Beetle



Indian Cobra



Station Wagon



Black Swan

“Understanding Neural Networks Through Deep Visualization”, Yosinski et al.

Beyond techniques discussed for general neural nets (back-prop, batch gradient descent, adaptive learning rates) training convolutional networks requires a lot of “tricks”.

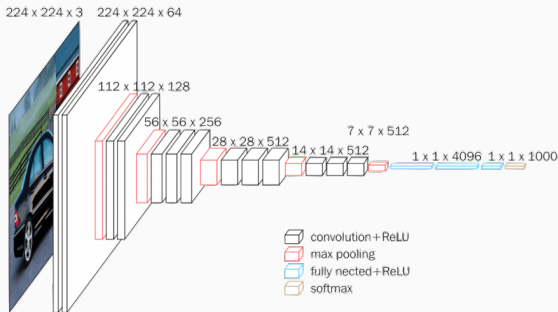
- Batch normalization (accelerate training).
- Dropout (prevent over-fitting)
- Residual connections (accelerate training, allow for more depth – 100s of layers).

And convolutional networks require **lots of training data**.

# TRANSFER LEARNING

What if you want to apply deep convolutional networks to a problem where you don't have a lot of data?

**Idea behind transfer learning:** features transformations learned when training a classifier on e.g. Imagenet are often useful in other problems, even with different inputs, classes, etc.

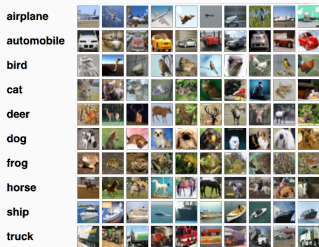


- Download state of the art pre-trained network (Alexnet, VGG, Inception, etc.)
- Chop off classification layer.
- Use first part of network as feature extractor.
- Solve classification problem using more scalable methods: kernel SVM, logistic regression, shallow fully connected net, etc.

**Very easy to do in Tensorflow/Keras.** Many pre-trained networks are made available.

## Two demos to be released shortly:

- Classification of CIFAR-10 dataset using 2 layer neural nets in Keras. You will likely want to use Google Collab to access a GPU.



- Transfer learning in Keras.