

Warm-up as You Log In

Suppose we have a function that takes in a vector and squares each element individually, returning another vector, $\mathbf{y} = f(\mathbf{x})$.

$$f\left(\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}\right) \rightarrow \begin{bmatrix} x_1^2 \\ x_2^2 \\ x_3^2 \end{bmatrix}$$

$$f\left(\begin{bmatrix} 7 \\ 3 \\ 5 \end{bmatrix}\right) \rightarrow \begin{bmatrix} 49 \\ 9 \\ 25 \end{bmatrix}$$

What is $\partial \mathbf{y} / \partial \mathbf{x}$?

Announcements

Assignments

- HW5
 - Due Mon, 10/26, 11:59 pm
 - Start early

Recitation

- No recitation this Friday

An abstract graphic on the left side of the slide. It features a sphere-like shape composed of a dense grid of lines. The lines are primarily red and green, with some blue lines interspersed. The grid is curved, following the shape of the sphere, and the lines are of varying thicknesses, creating a complex, woven appearance. The sphere is set against a dark gray background.

Introduction to Machine Learning

Neural Networks

Instructor: Pat Virtue

Plan

Last Time

- Neural Networks
 - Building blocks
 - Optimization

Today

- Neural Networks
 - Wrap-up calculus
 - Universal Approximation Theorem
 - Convolutional neural networks

Backpropagation (so-far)

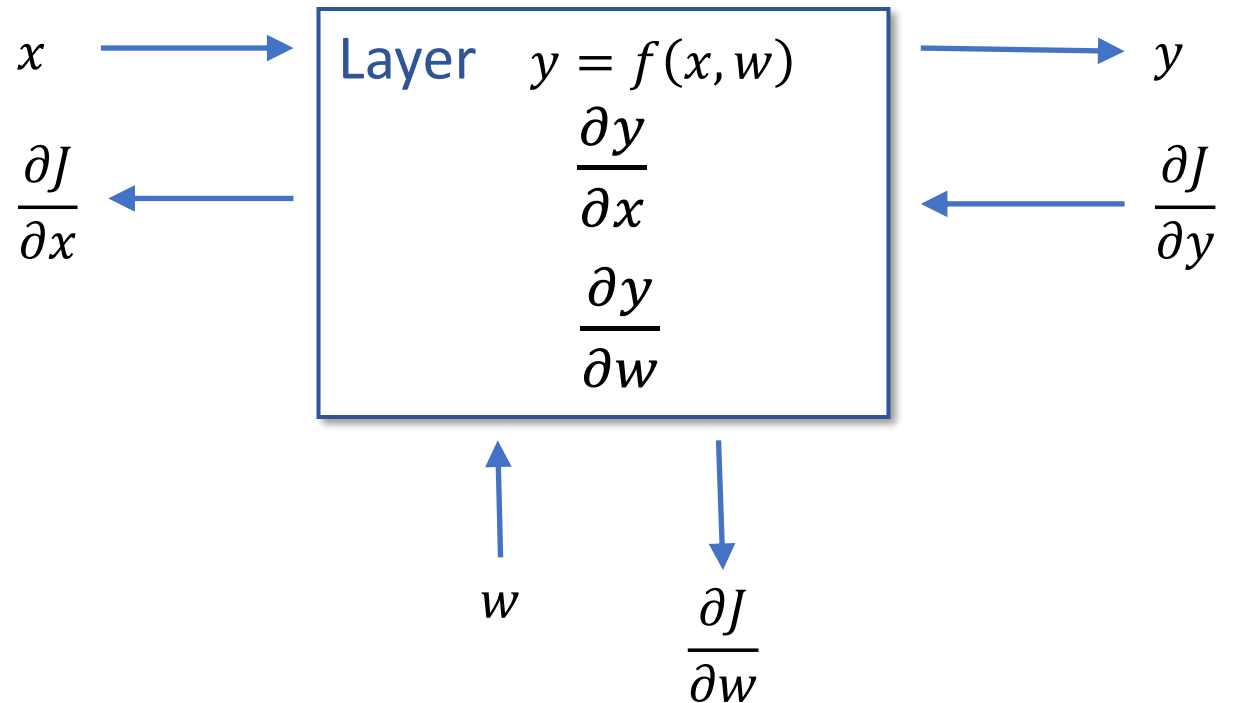
Compute derivatives per layer, utilizing previous derivatives

Objective: $J(\mathbf{w})$

Arbitrary layer: $y = f(x, w)$

Need:

- $\frac{\partial J}{\partial x} = \frac{\partial J}{\partial y} \frac{\partial y}{\partial x}$
- $\frac{\partial J}{\partial w} = \frac{\partial J}{\partial y} \frac{\partial y}{\partial w}$



Matrix Calculus

Jacobian: Vector in, vector out

Numerator-layout

$$\mathbf{y} = f(\mathbf{x}) \quad \mathbf{y} \in \mathbb{R}^N, \quad \mathbf{x} \in \mathbb{R}^M, \quad \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \in \mathbb{R}^{N \times M}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_M} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_N}{\partial x_1} & \cdots & \frac{\partial y_N}{\partial x_M} \end{bmatrix}$$

Matrix Calculus

Vector in, scalar out

Numerator-layout

$$y = f(\boldsymbol{x}) \quad y \in \mathbb{R}, \quad \boldsymbol{x} \in \mathbb{R}^M, \quad \frac{\partial y}{\partial \boldsymbol{x}} \in \mathbb{R}^{1 \times M}$$

$$\frac{\partial y}{\partial \boldsymbol{x}} = \left[\frac{\partial y}{\partial x_1} \quad \cdots \quad \frac{\partial y}{\partial x_M} \right]$$

Matrix Calculus

Scalar in, vector out

Numerator-layout

$$\mathbf{y} = f(x) \quad \mathbf{y} \in \mathbb{R}^N, \quad x \in \mathbb{R}, \quad \frac{\partial \mathbf{y}}{\partial x} \in \mathbb{R}^{N \times 1}$$

$$\frac{\partial \mathbf{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \vdots \\ \frac{\partial y_N}{\partial x} \end{bmatrix}$$

Matrix Calculus

Gradient: Vector in, scalar out

Transpose of numerator-layout

$$y = f(\boldsymbol{x}) \quad y \in \mathbb{R}, \quad \boldsymbol{x} \in \mathbb{R}^M, \quad \frac{\partial y}{\partial \boldsymbol{x}} \in \mathbb{R}^{1 \times M}, \quad \nabla_{\boldsymbol{x}} f \in \mathbb{R}^{M \times 1}$$

$$\frac{\partial y}{\partial \boldsymbol{x}} = (\nabla_{\boldsymbol{x}} f)^T$$

Matrix Calculus

Matrix in, scalar out

Keep same dimensions as matrix

$$y = f(\mathbf{X}) \quad y \in \mathbb{R}, \quad \mathbf{X} \in \mathbb{R}^{N \times M}, \quad \frac{\partial y}{\partial \mathbf{X}} \in \mathbb{R}^{N \times M}$$

$$\frac{\partial y}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial y}{\partial X_{1,1}} & \cdots & \frac{\partial y}{\partial X_{1,M}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial X_{N,1}} & \cdots & \frac{\partial y}{\partial X_{N,M}} \end{bmatrix}$$

Warm-up as You Log In

Suppose we have a function that takes in a vector and squares each element individually, returning another vector, $\mathbf{y} = f(\mathbf{x})$.

$$f\left(\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}\right) \rightarrow \begin{bmatrix} x_1^2 \\ x_2^2 \\ x_3^2 \end{bmatrix}$$

$$f\left(\begin{bmatrix} 7 \\ 3 \\ 5 \end{bmatrix}\right) \rightarrow \begin{bmatrix} 49 \\ 9 \\ 25 \end{bmatrix}$$

What is $\partial \mathbf{y} / \partial \mathbf{x}$?

Calculus Chain Rule

Scalar:

$$y = f(z)$$

$$z = g(x)$$

$$\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx}$$

Multivariate:

$$y = f(\mathbf{z})$$

$$\mathbf{z} = g(x)$$

$$\frac{dy}{dx} = \sum_j \frac{\partial y}{\partial z_j} \frac{\partial z_j}{\partial x}$$

Multivariate:

$$\mathbf{y} = f(\mathbf{z})$$

$$\mathbf{z} = g(\mathbf{x})$$

$$\frac{dy_i}{dx_k} = \sum_j \frac{\partial y_i}{\partial z_j} \frac{\partial z_j}{\partial x_k}$$

Piazza Poll 1

$$y = f(\mathbf{z}) \quad y \in \mathbb{R}, \mathbf{z} \in \mathbb{R}^N, x \in \mathbb{R}$$

$$\mathbf{z} = g(x)$$

Select all that apply

$$\frac{\partial y}{\partial x} = \dots$$

A. $\frac{\partial y}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial x}$

B. $\left(\frac{\partial y}{\partial \mathbf{z}}\right)^T \frac{\partial \mathbf{z}}{\partial x}$

C. $\frac{\partial y}{\partial \mathbf{z}} \left(\frac{\partial \mathbf{z}}{\partial x}\right)^T$

D. $\left(\frac{\partial y}{\partial \mathbf{z}}\right)^T \left(\frac{\partial \mathbf{z}}{\partial x}\right)^T$

E. $\left(\frac{\partial y}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial x}\right)^T$

F. None of the above

Piazza Poll 1

$$y = f(\mathbf{z}) \quad y \in \mathbb{R}, \mathbf{z} \in \mathbb{R}^N, x \in \mathbb{R}$$

$$\mathbf{z} = g(x)$$

Select all that apply

$$\frac{\partial y}{\partial x} = \dots$$

A. $\frac{\partial y}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial x}$

B. $\left(\frac{\partial y}{\partial \mathbf{z}}\right)^T \frac{\partial \mathbf{z}}{\partial x}$

C. $\frac{\partial y}{\partial \mathbf{z}} \left(\frac{\partial \mathbf{z}}{\partial x}\right)^T$

D. $\left(\frac{\partial y}{\partial \mathbf{z}}\right)^T \left(\frac{\partial \mathbf{z}}{\partial x}\right)^T$

E. $\left(\frac{\partial y}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial x}\right)^T$

F. None of the above

Piazza Poll 2

$$y = f(\mathbf{z}) \quad y \in \mathbb{R}, \mathbf{z} \in \mathbb{R}^N, \mathbf{x} \in \mathbb{R}^M$$

$$\mathbf{z} = g(\mathbf{x})$$

Select all that apply

$$\frac{\partial y}{\partial \mathbf{x}} = \dots$$

A. $\frac{\partial y}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{x}}$

B. $\left(\frac{\partial y}{\partial \mathbf{z}}\right)^T \frac{\partial \mathbf{z}}{\partial \mathbf{x}}$

C. $\frac{\partial y}{\partial \mathbf{z}} \left(\frac{\partial \mathbf{z}}{\partial \mathbf{x}}\right)^T$

D. $\left(\frac{\partial y}{\partial \mathbf{z}}\right)^T \left(\frac{\partial \mathbf{z}}{\partial \mathbf{x}}\right)^T$

E. $\left(\frac{\partial y}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{x}}\right)^T$

F. None of the above

Piazza Poll 2

$$y = f(\mathbf{z})$$

$$\mathbf{z} = g(\mathbf{x})$$

Select all that apply

$$\frac{\partial y}{\partial \mathbf{x}} = \dots$$

A. $\frac{\partial y}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{x}}$

B. $\left(\frac{\partial y}{\partial \mathbf{z}}\right)^T \frac{\partial \mathbf{z}}{\partial \mathbf{x}}$

C. $\frac{\partial y}{\partial \mathbf{z}} \left(\frac{\partial \mathbf{z}}{\partial \mathbf{x}}\right)^T$

D. $\left(\frac{\partial y}{\partial \mathbf{z}}\right)^T \left(\frac{\partial \mathbf{z}}{\partial \mathbf{x}}\right)^T$

E. $\left(\frac{\partial y}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{x}}\right)^T$

F. None of the above

Network Optimization

$$J(\mathbf{w}) = z_4$$

$$z_4 = f_4(w_D, w_E, z_2, z_3)$$

$$z_3 = f_3(w_C, z_1)$$

$$z_2 = f_2(w_B, z_1)$$

$$z_1 = f_1(w_A, x)$$

Need multivariate chain rule!

Network Optimization

$$J(\mathbf{w}) = z_4$$

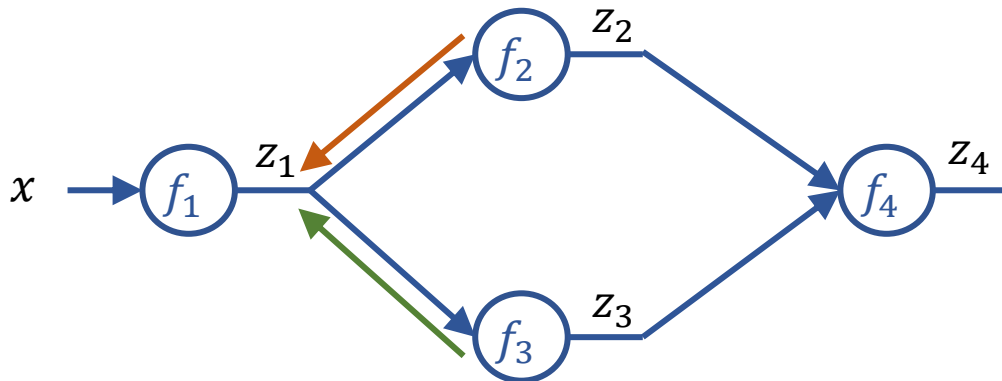
$$z_4 = f_4(w_D, w_E, z_2, z_3)$$

$$z_3 = f_3(w_C, z_1)$$

$$z_2 = f_2(w_B, z_1)$$

$$z_1 = f_1(w_A, x)$$

Need multivariate chain rule!



$$\frac{\partial J}{\partial w_E} = \frac{\partial J}{\partial z_4} \frac{\partial z_4}{\partial w_E}$$

$$\frac{\partial J}{\partial w_D} = \frac{\partial J}{\partial z_4} \frac{\partial z_4}{\partial w_D}$$

$$\frac{\partial J}{\partial z_3} = \frac{\partial J}{\partial z_4} \frac{\partial z_4}{\partial z_3}$$

$$\frac{\partial J}{\partial z_2} = \frac{\partial J}{\partial z_4} \frac{\partial z_4}{\partial z_2}$$

$$\frac{\partial J}{\partial w_C} = \frac{\partial J}{\partial z_3} \frac{\partial z_3}{\partial w_C}$$

$$\frac{\partial J}{\partial w_B} = \frac{\partial J}{\partial z_2} \frac{\partial z_2}{\partial w_B}$$

$$\frac{\partial J}{\partial z_1} = \frac{\partial J}{\partial z_2} \frac{\partial z_2}{\partial z_1} + \frac{\partial J}{\partial z_3} \frac{\partial z_3}{\partial z_1}$$

$$\frac{\partial J}{\partial w_A} = \frac{\partial J}{\partial z_1} \frac{\partial z_1}{\partial w_A}$$

Backpropagation (updated)

Compute derivatives per layer, utilizing previous derivatives

Objective: $J(\mathbf{w})$

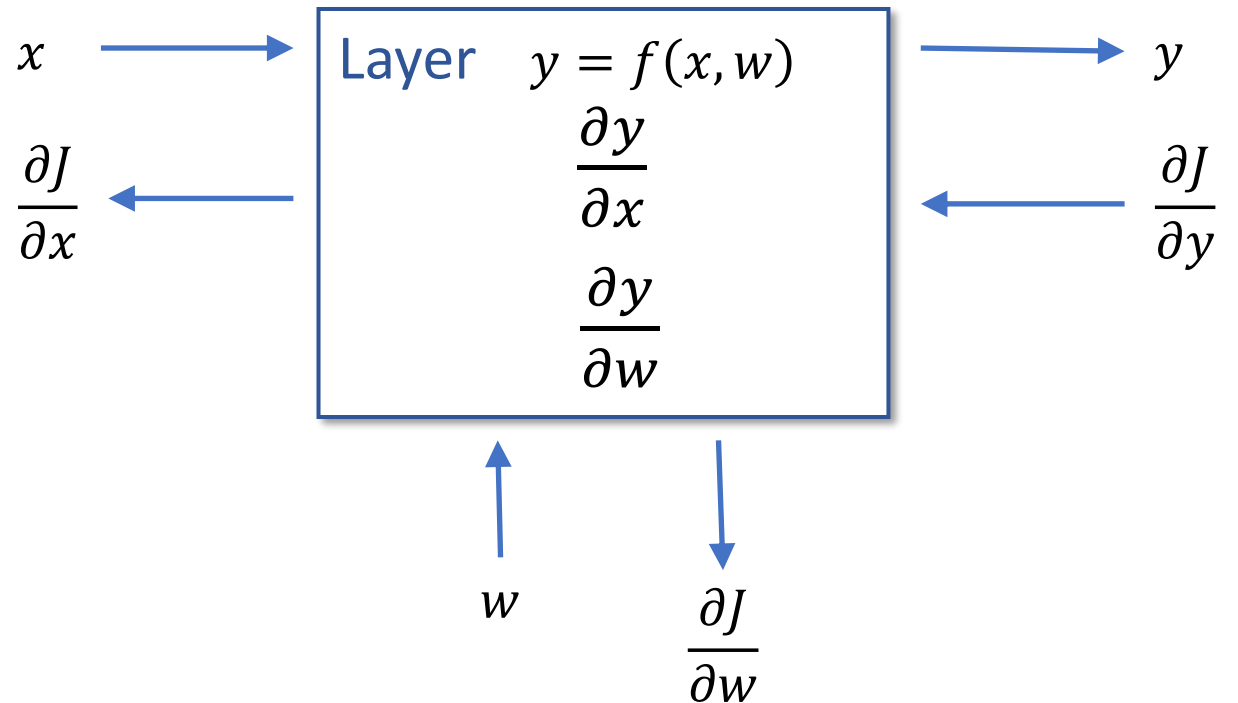
Arbitrary layer: $y = f(x, w)$

Init:

- $\frac{\partial J}{\partial x} = 0$
- $\frac{\partial J}{\partial w} = 0$

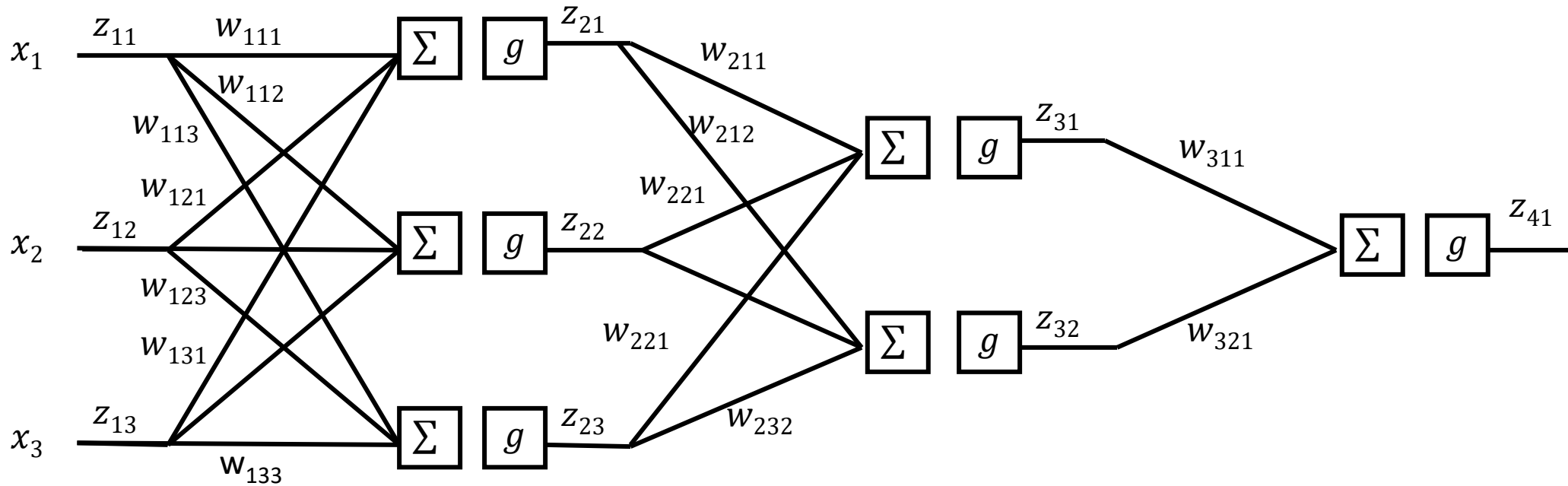
Compute:

- $\frac{\partial J}{\partial x} \text{ + } = \frac{\partial J}{\partial y} \frac{\partial y}{\partial x}$
- $\frac{\partial J}{\partial w} \text{ + } = \frac{\partial J}{\partial y} \frac{\partial y}{\partial w}$



Neural Network Implementation

Which pieces to we treat as functions?



Neural Networks Properties

Practical considerations

- Large number of neurons
 - Danger for overfitting
- Modelling assumptions vs data assumptions trade-off
- Gradient descent can easily get stuck local optima

What if there are no non-linear activations?

- A deep neural network with only linear layers can be reduced to an exactly equivalent single linear layer

Universal Approximation Theorem:

- A two-layer neural network with a sufficient number of neurons can approximate any continuous function to any desired accuracy.

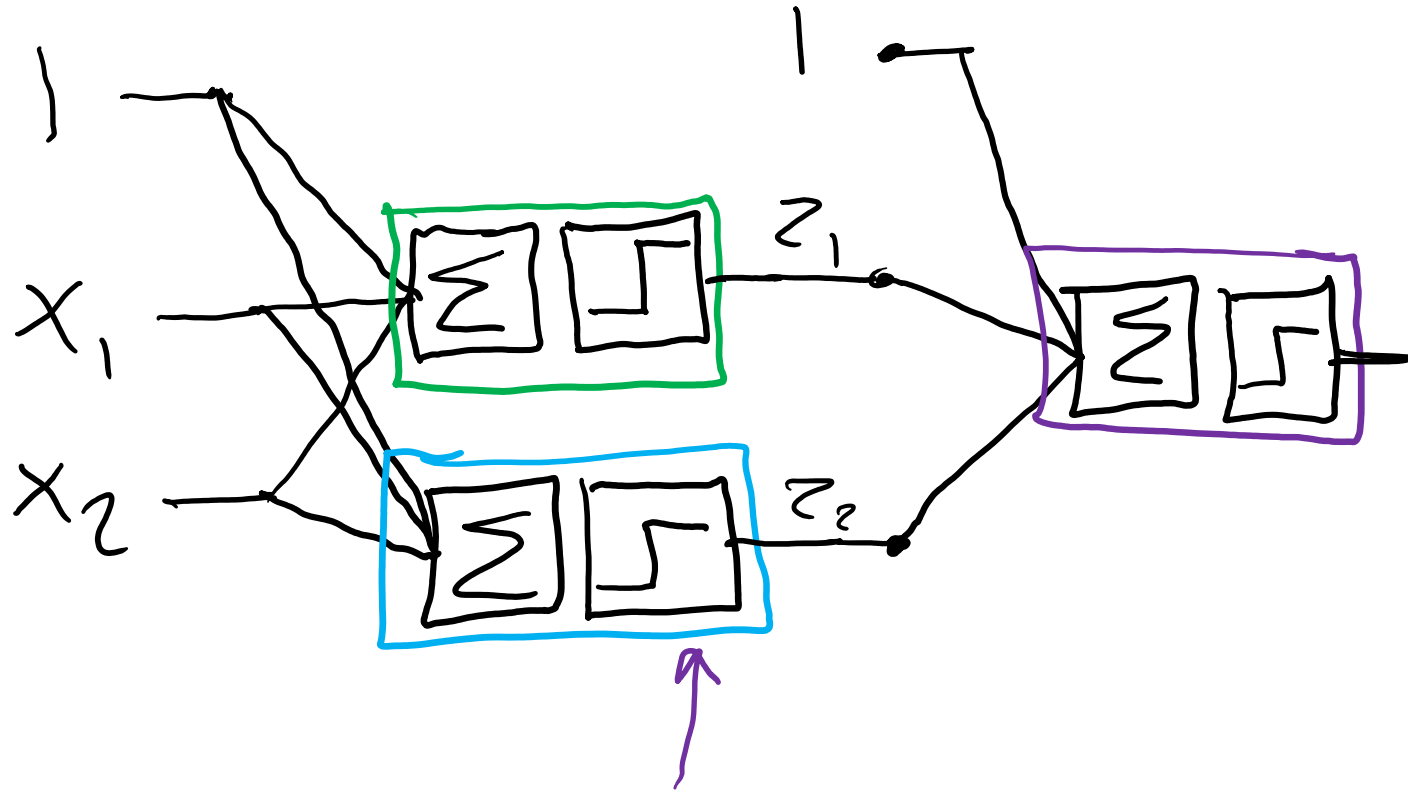
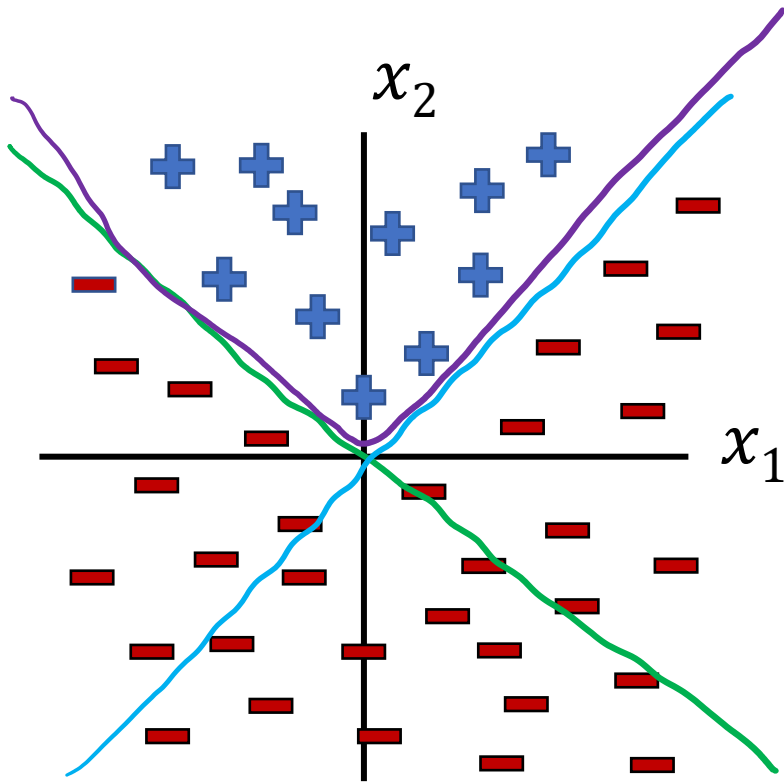
Classification Design Challenge

How could you configure three specific perceptrons to classify this data?

$$h_A(\mathbf{x}) = \text{sign}(\mathbf{w}_A^T \mathbf{x} + b_A)$$

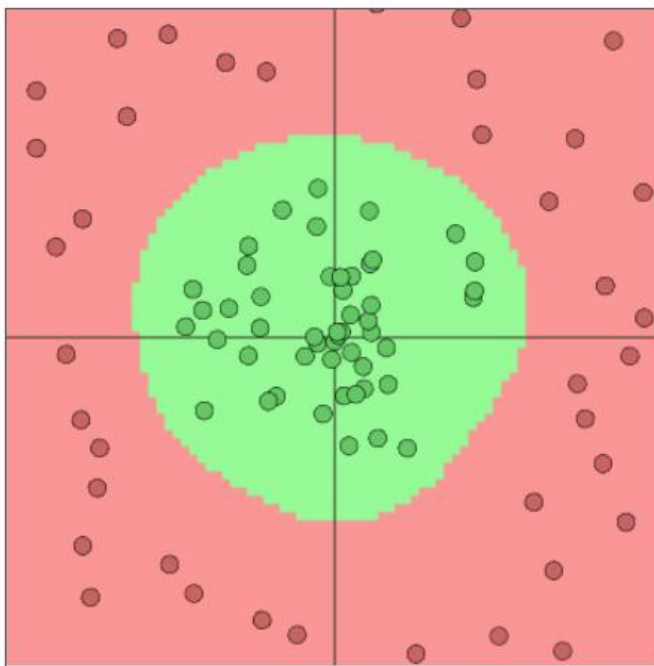
$$h_B(\mathbf{x}) = \text{sign}(\mathbf{w}_B^T \mathbf{x} + b_B)$$

$$h_C(\mathbf{z}) = \text{sign}(\mathbf{w}_C^T \mathbf{z} + b_C)$$



Network to Approximate Binary Classification

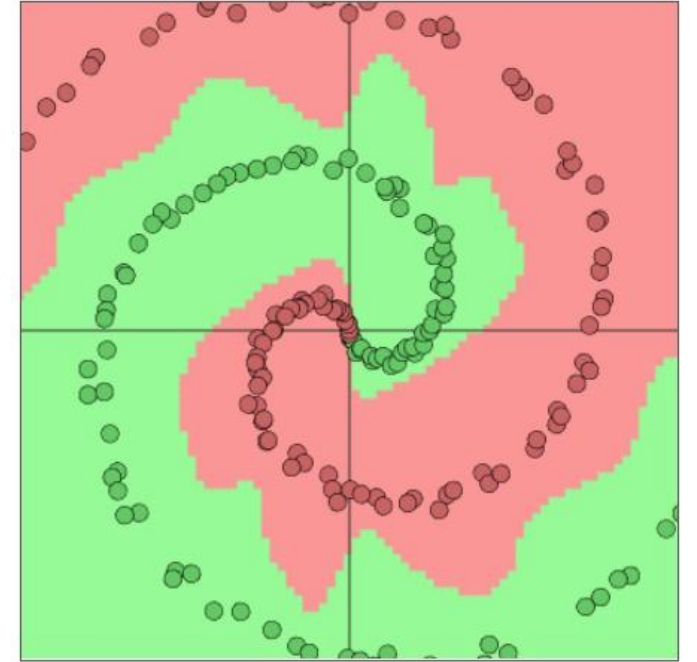
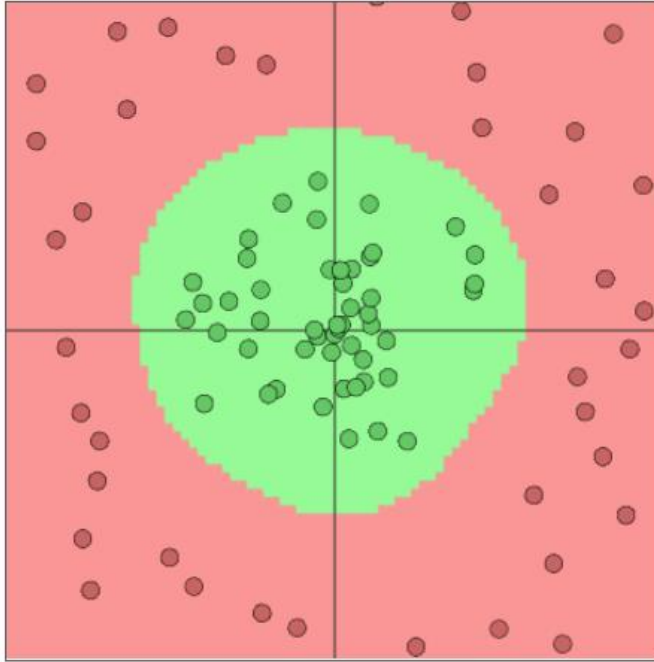
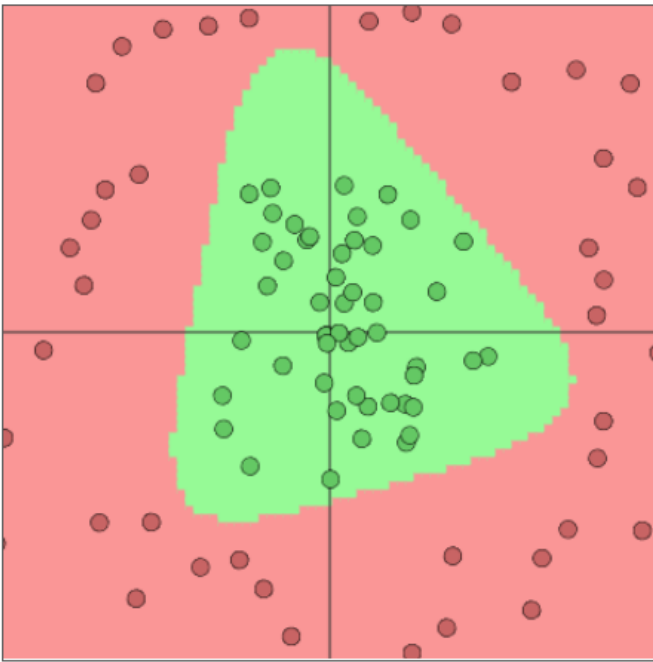
Approximate arbitrary decision boundary



<https://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>

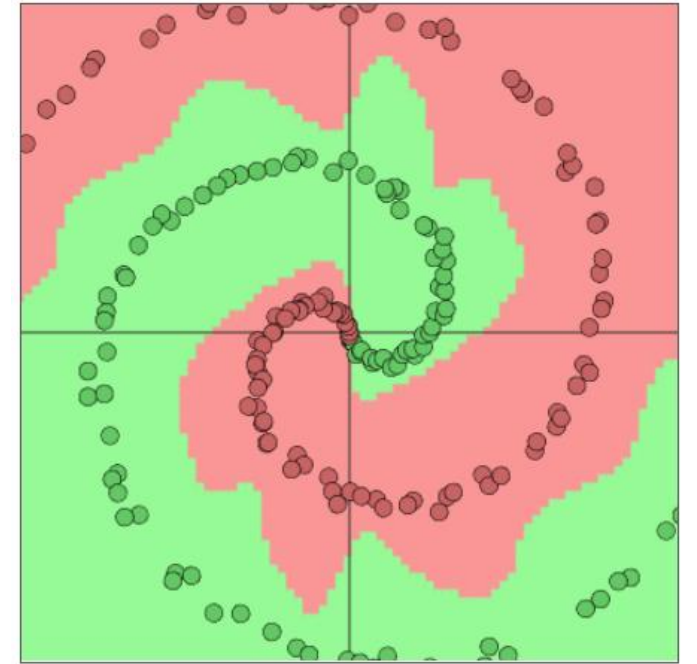
Network to Approximate Binary Classification

Approximate arbitrary decision boundary



Network to Approximate Binary Classification

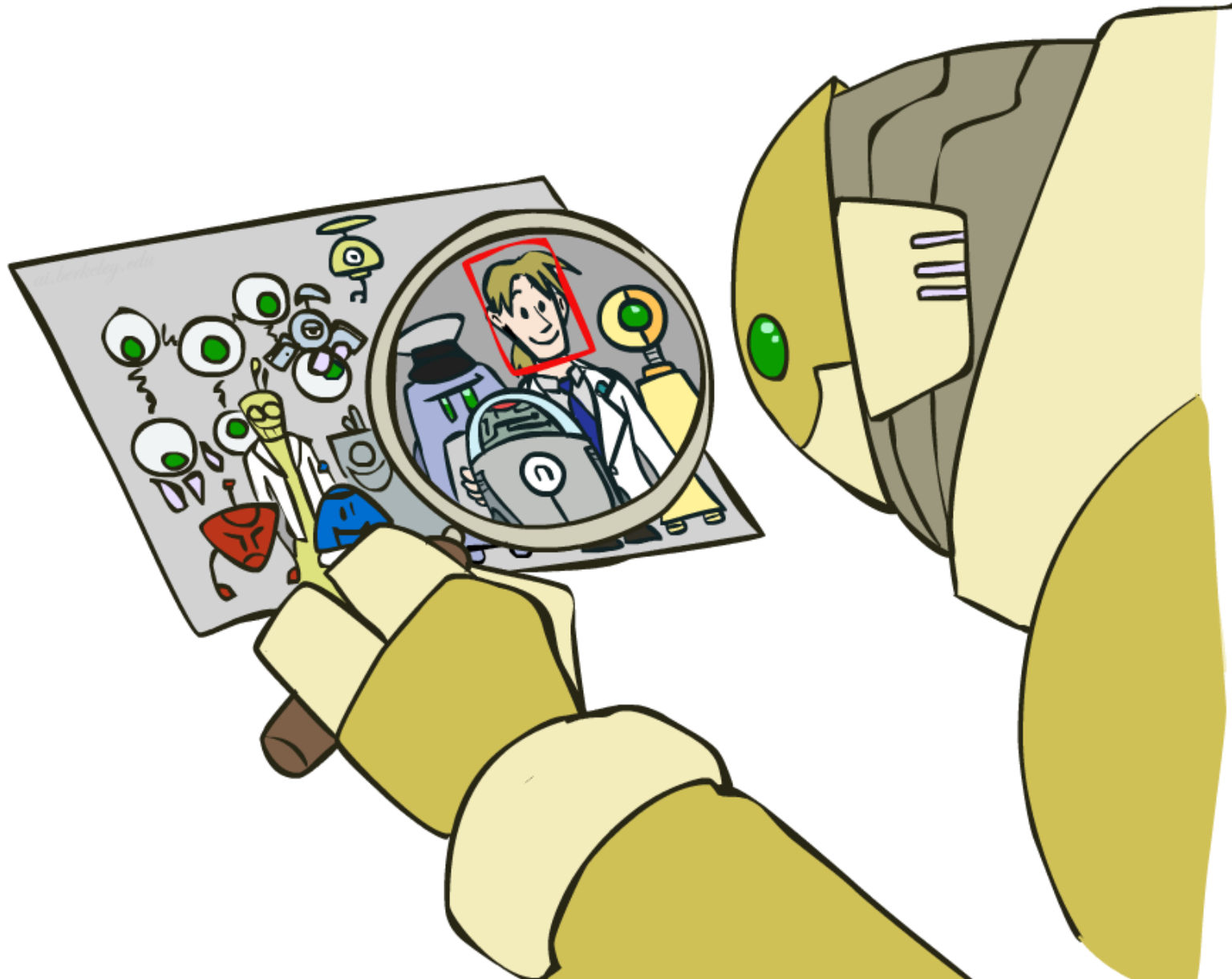
Approximate arbitrary decision boundary



<https://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>

Convolutional Neural Nets

Computer Vision: How far along are we?

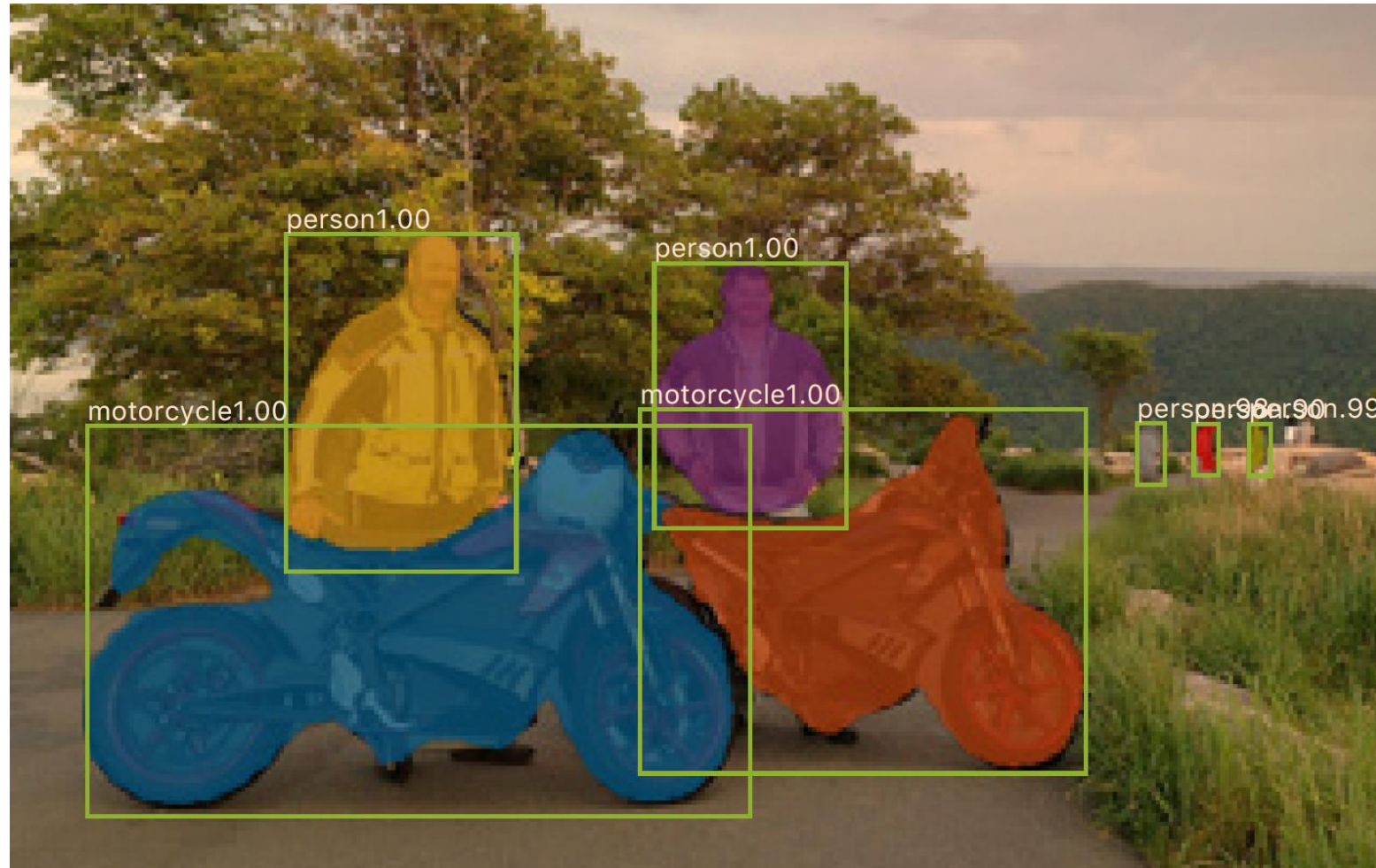


Computer Vision: How far along are we?



Terminator 2, 1991 <https://www.youtube.com/watch?v=9MeaaCwBW28>

Computer Vision: How far along are we?



0.2 seconds
per image

Mask R-CNN

He, Kaiming, et al. "Mask R-CNN." *Computer Vision (ICCV), 2017 IEEE International Conference on*. IEEE, 2017.

Computer Vision: How far along are we?



“My CPU is a neural net processor, a learning computer”

Terminator 2, 1991

Computer Vision: Autonomous Driving



Tesla, Inc: <https://vimeo.com/192179726>

Computer Vision: Domain Transfer

CycleGAN



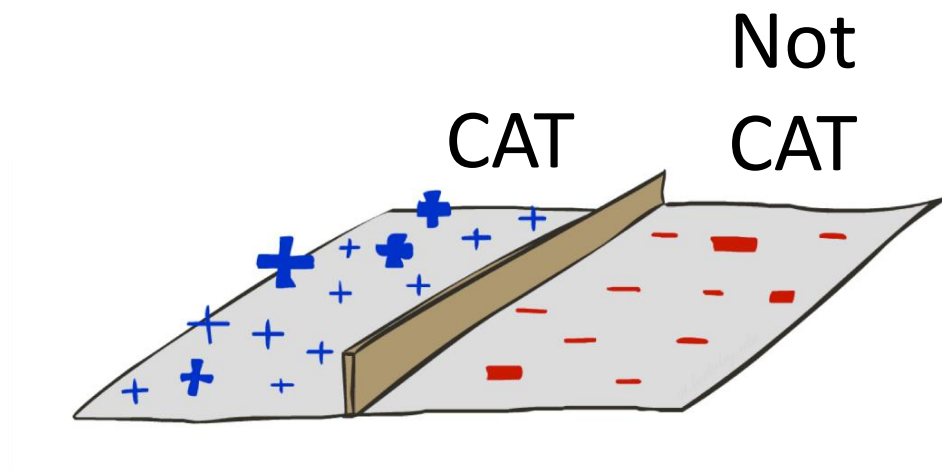
Jun-Yan Zhu*, Taesung Park*, Phillip Isola, and Alexei A. Efros. "Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks", ICCV 2017.

Outline

1. Measuring the current state of computer vision
2. Why convolutional neural networks
 - Old school computer vision
 - Image features and classification
3. Convolution “nuts and bolts”

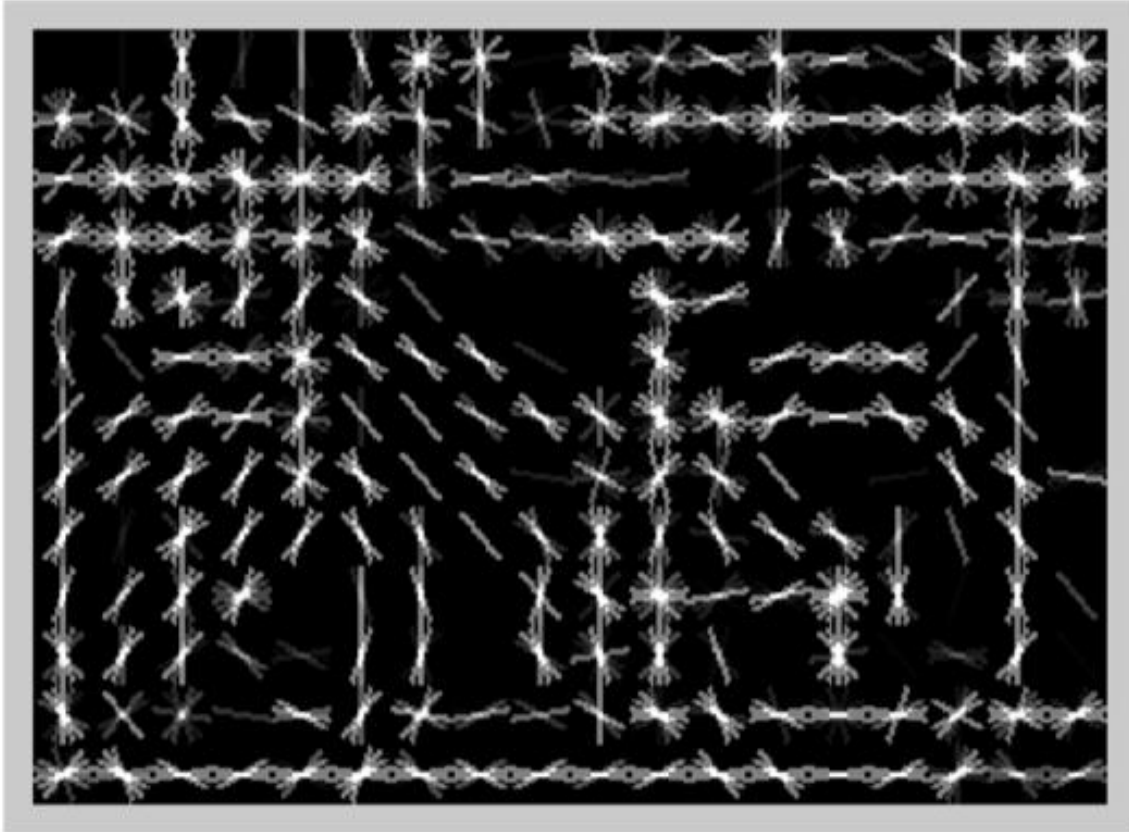
Image Classification

What's the problem with just directly classifying raw pixels in high dimensional space?



CAT

Image Classification



[Dalal and Triggs, 2005]

HoG Filter

HoG: Histogram of oriented gradients

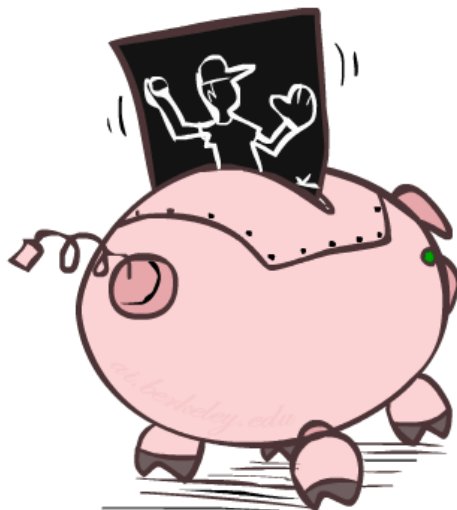
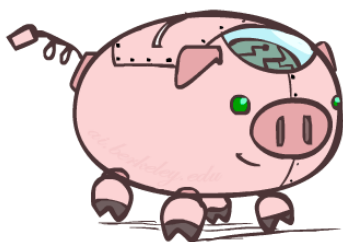
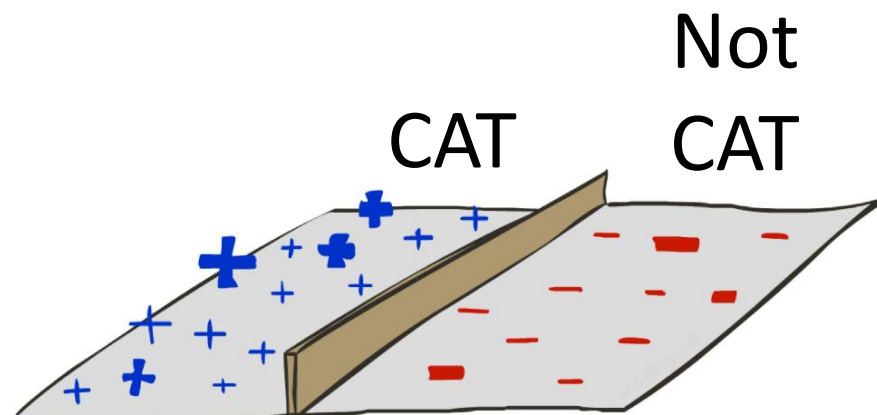


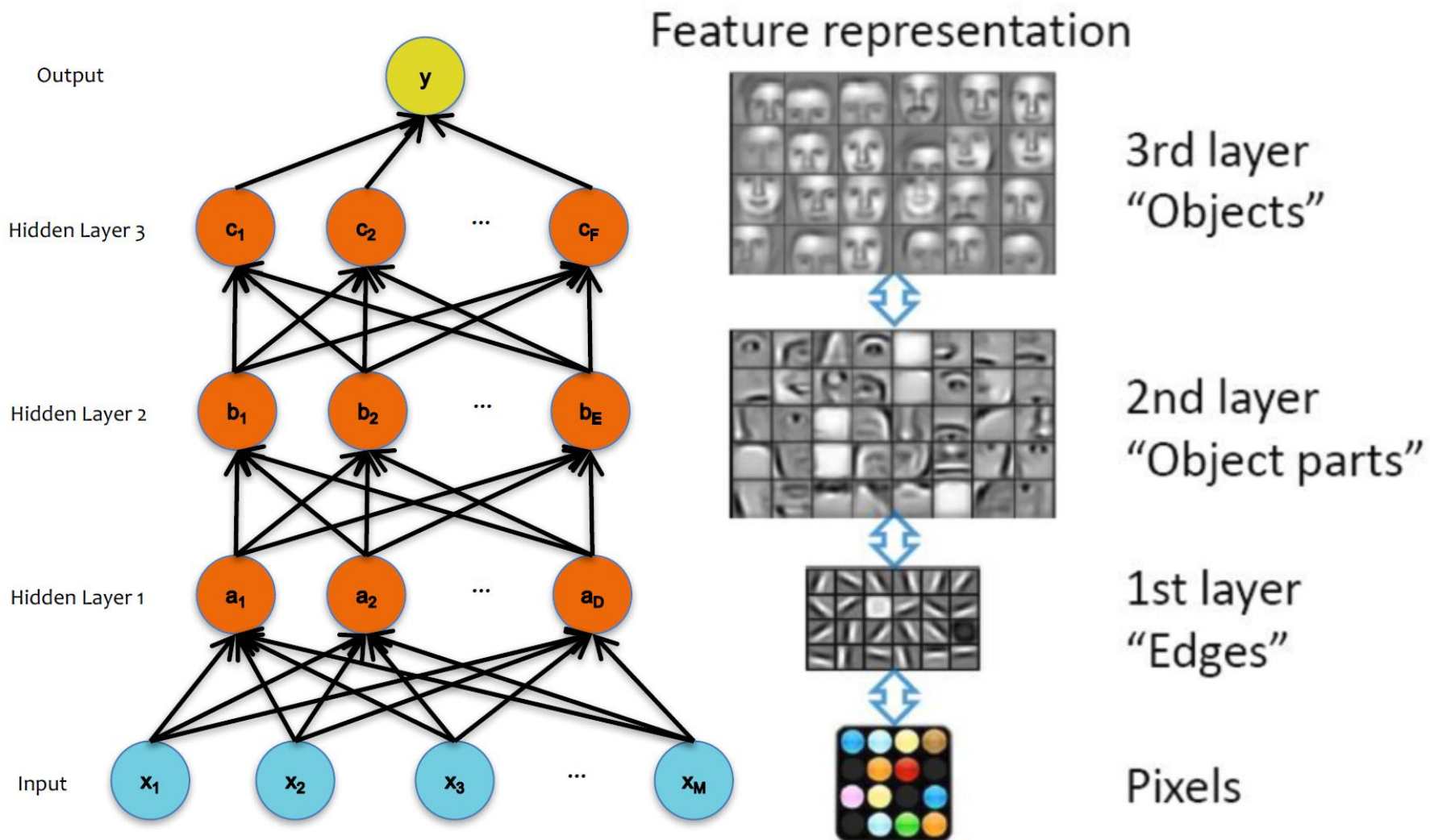
Image Classification

HOG features passed to a linear classifier (logistic regression / SVM)



CAT

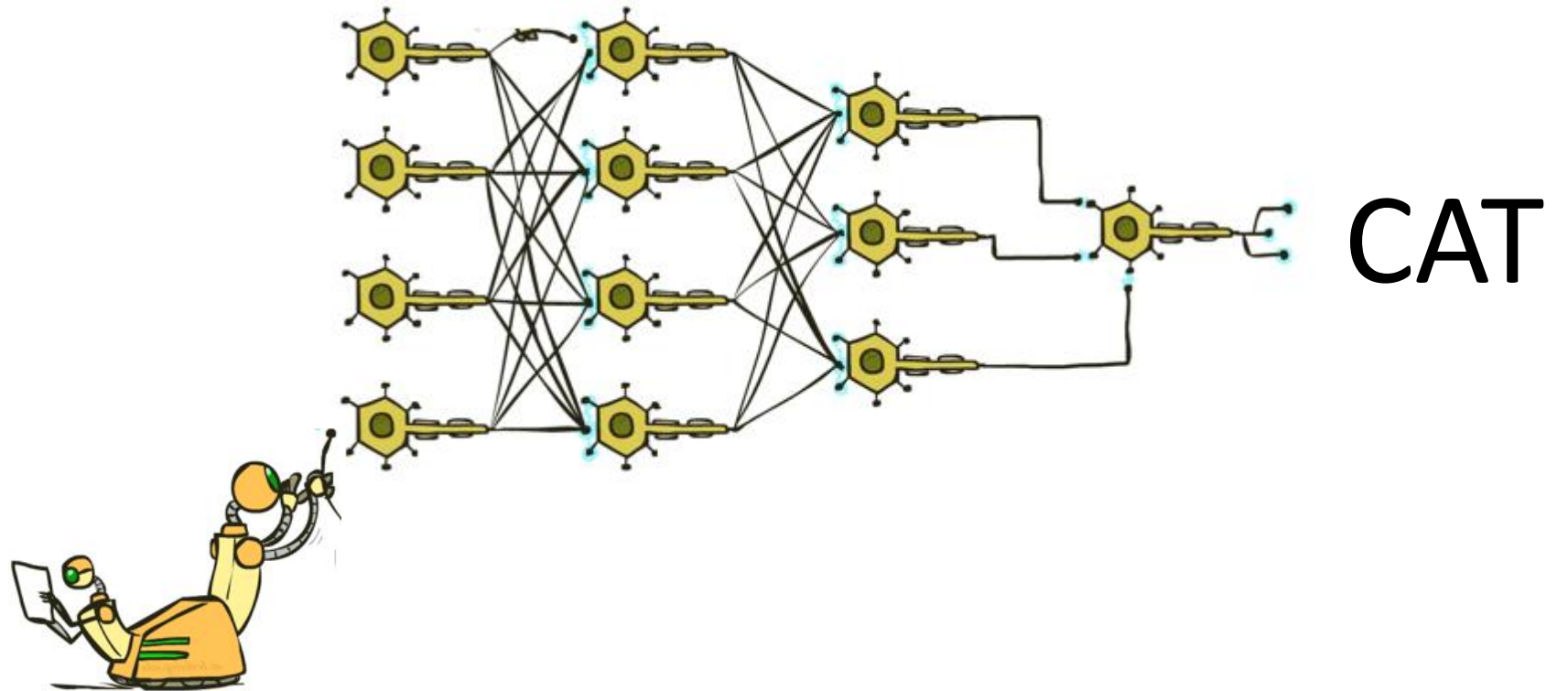
Classification: Learning Features



Example from Honglak Lee (NIPS 2010)

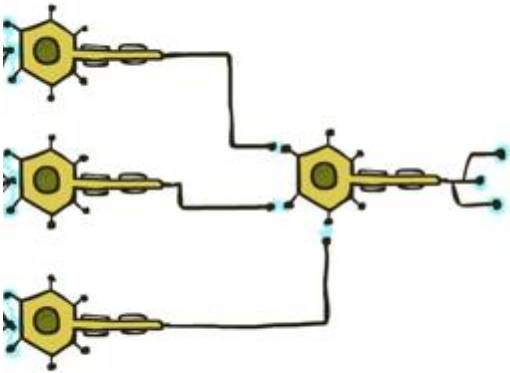
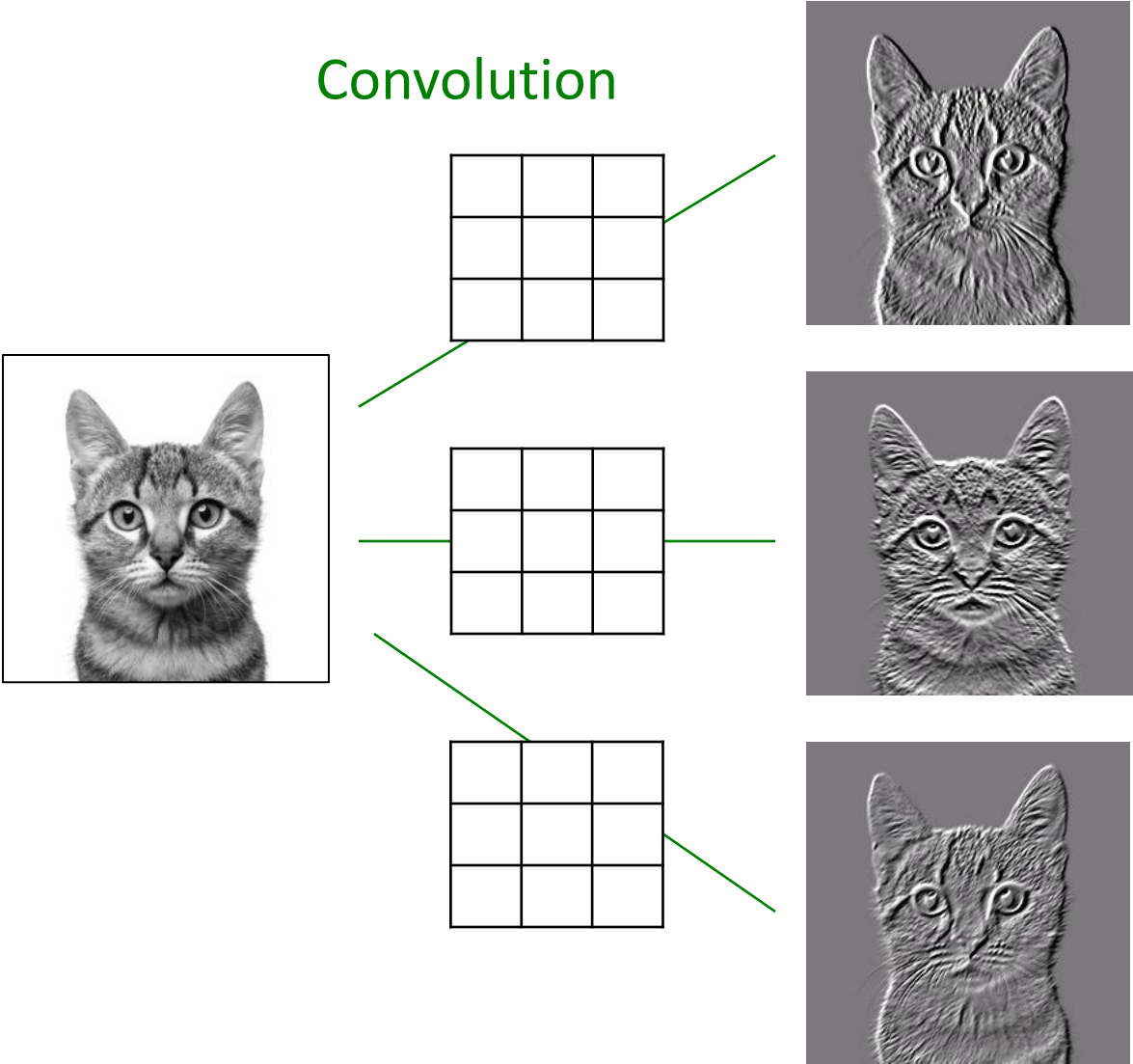
Classification: Deep Learning

Fully connected neural network?



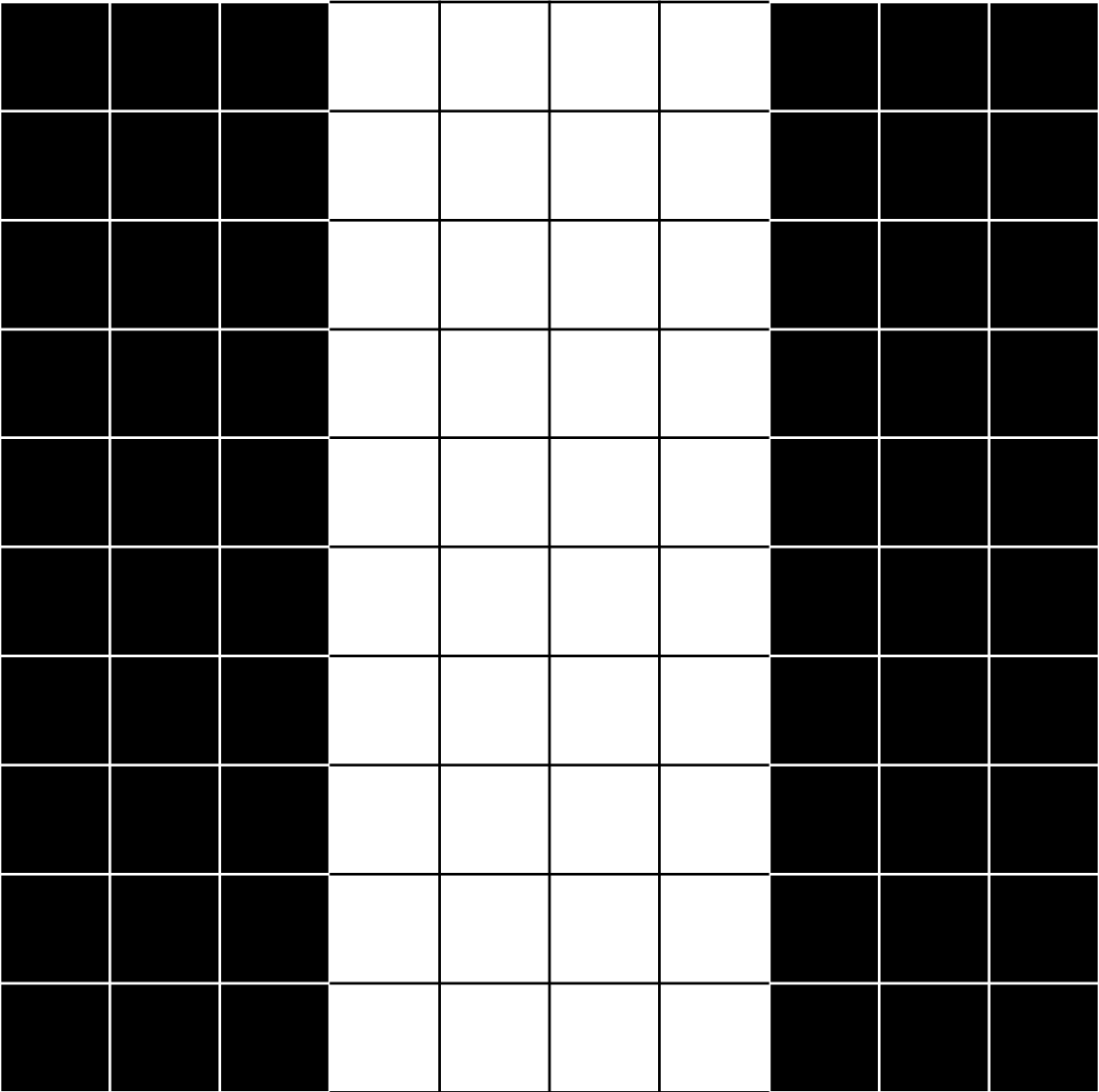
Convolutional Neural Networks

Convolution



CAT

Convolution



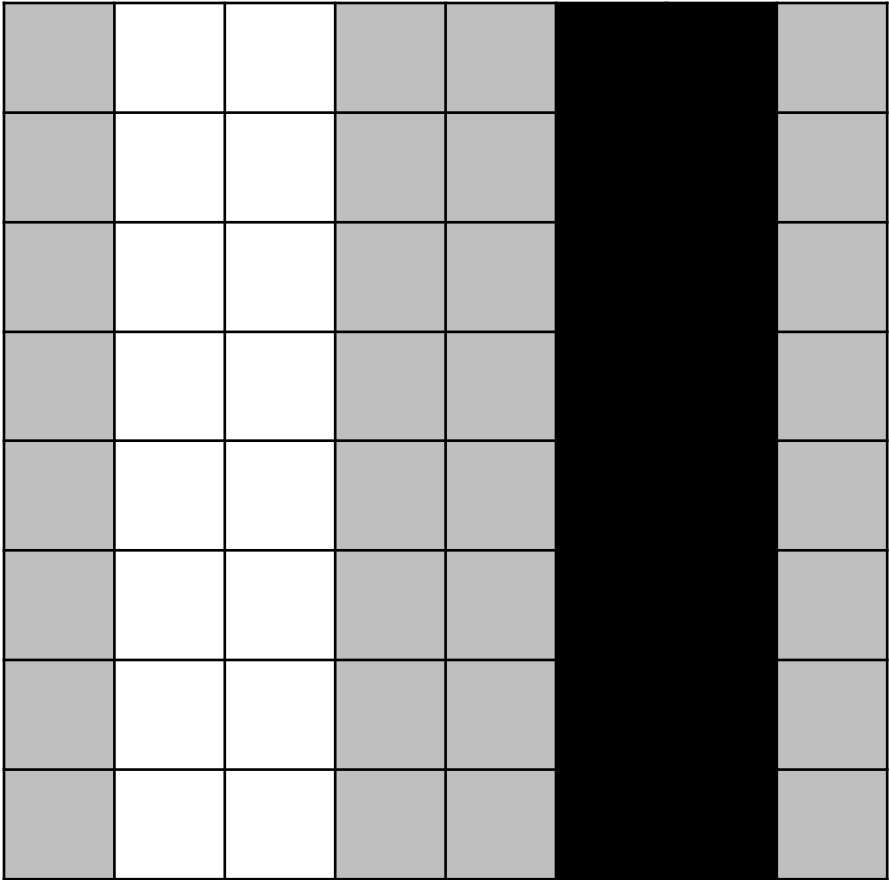
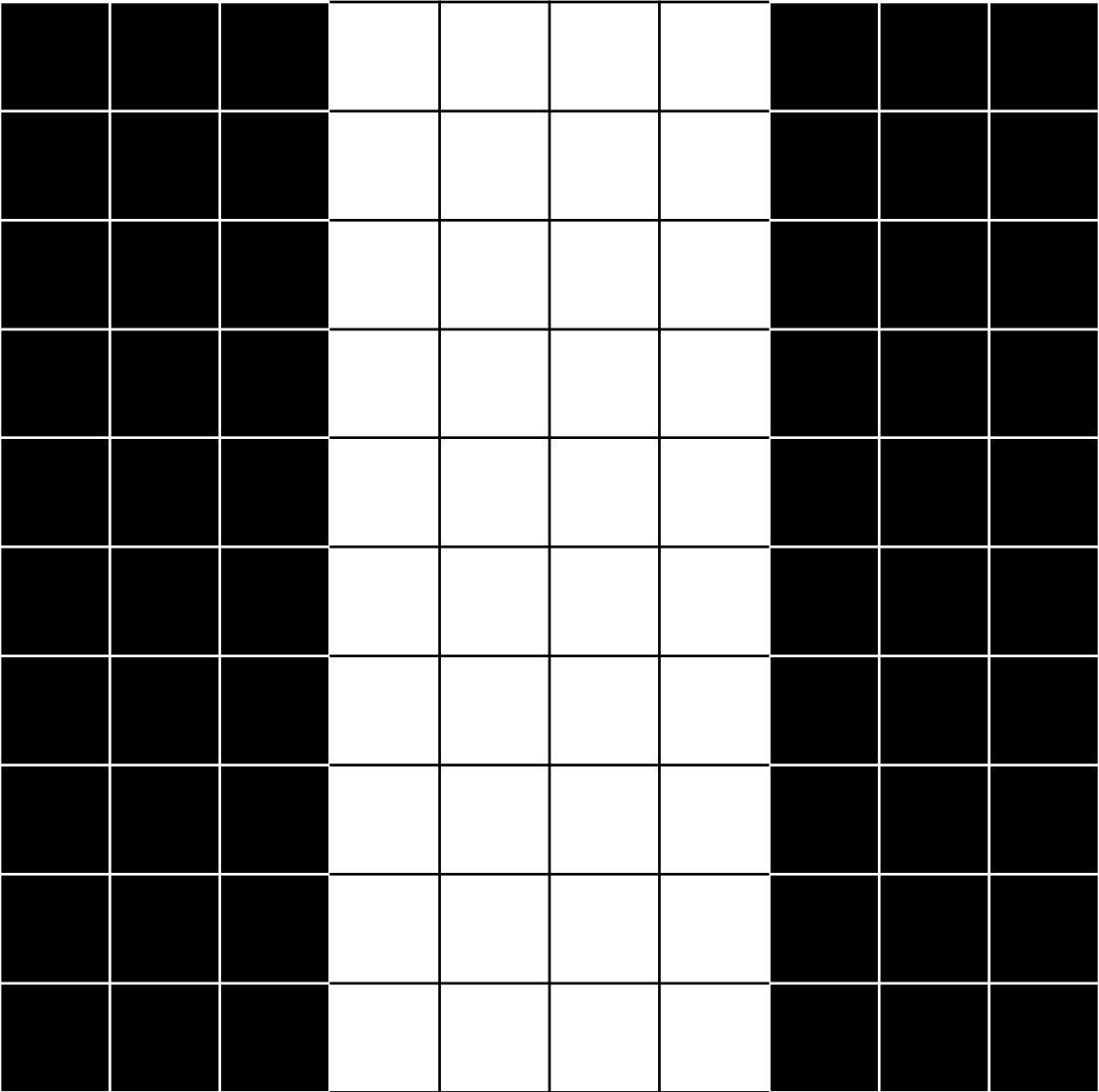
-1	0	1
-1	0	1
-1	0	1

Convolution

0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0

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

Convolution



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

Convolution

Signal processing definition

$$z[i, j] = \sum_{u=-\infty}^{\infty} \sum_{v=-\infty}^{\infty} x[i-u, j-v] \cdot w[u, v]$$

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

Relaxed definition

- Drop infinity; don't flip kernel

$$z[i, j] = \sum_{u=0}^{K-1} \sum_{v=0}^{K-1} x[i+u, j+v] \cdot w[u, v]$$

A 6x6 grid of squares. The top-left 3x3 subgrid is outlined with a thick blue border. The remaining cells in the grid are outlined with thin black borders.

Convolution

Relaxed definition

$$z[i, j] = \sum_{u=0}^{K-1} \sum_{v=0}^{K-1} x[i+u, j+v] \cdot w[u, v]$$

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

```
for i in range(0, im_width - K + 1):
    for j in range(0, im_height - K + 1):
        im_out[i,j] = 0
        for u in range(0, K):
            for v in range(0, K):
                im_out[i,j] += im[i+u, j+v] * kernel[u,v]
```

GPU!!

A 6x6 grid of squares. The top-left 3x3 subgrid is outlined with a thick blue border. The remaining cells in the grid are outlined with thin black borders.

Convolution: Padding

0	0	1	1	1	1	0	0
0	0	1	1	1	1	0	0
0	0	1	1	1	1	0	0
0	0	1	1	1	1	0	0
0	0	1	1	1	1	0	0
0	0	1	1	1	1	0	0
0	0	1	1	1	1	0	0
0	0	1	1	1	1	0	0

0	2	2	0	0	-2	-2	0
0	3	3	0	0	-3	-3	0
0	3	3	0	0	-3	-3	0
0	3	3	0	0	-3	-3	0
0	3	3	0	0	-3	-3	0
0	3	3	0	0	-3	-3	0
0	3	3	0	0	-3	-3	0
0	2	2	0	0	-2	-2	0

Piazza Poll 3 : Which kernel goes with which output image?

Input



K1

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

K2

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

K3

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

Im1



Im2



Im3



Piazza Poll 3: Which kernel goes with which output image?

Input



K1

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

K2

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

K3

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

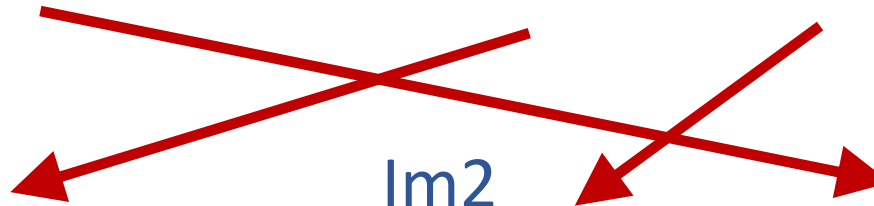
Im1



Im2



Im3



Convolutional Neural Networks

Convolution

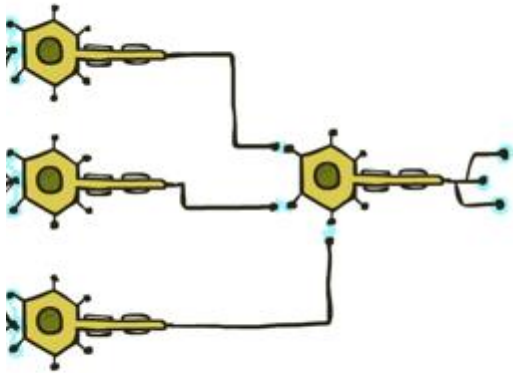
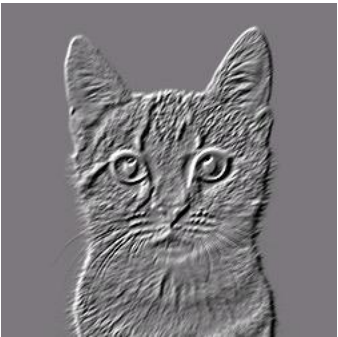
-1	0	1
-2	0	2
-1	0	1



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

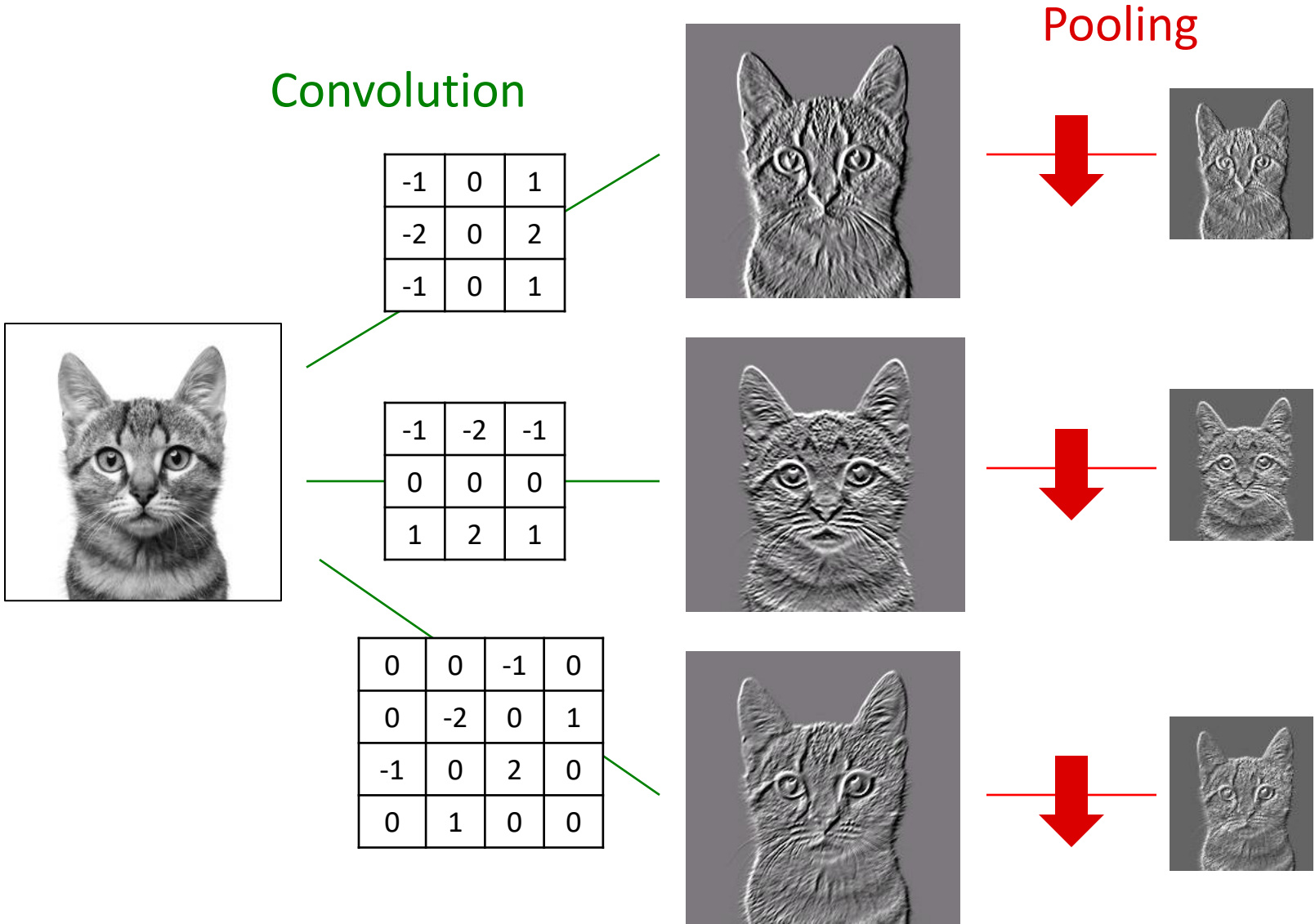


0	0	-1	0
0	-2	0	1
-1	0	2	0
0	1	0	0



CAT

Convolutional Neural Networks



Convolution: Stride=2

0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0
0	0	0	1	1	1	1	0	0	0

.25	.25
.25	.25

Stride: Max Pooling

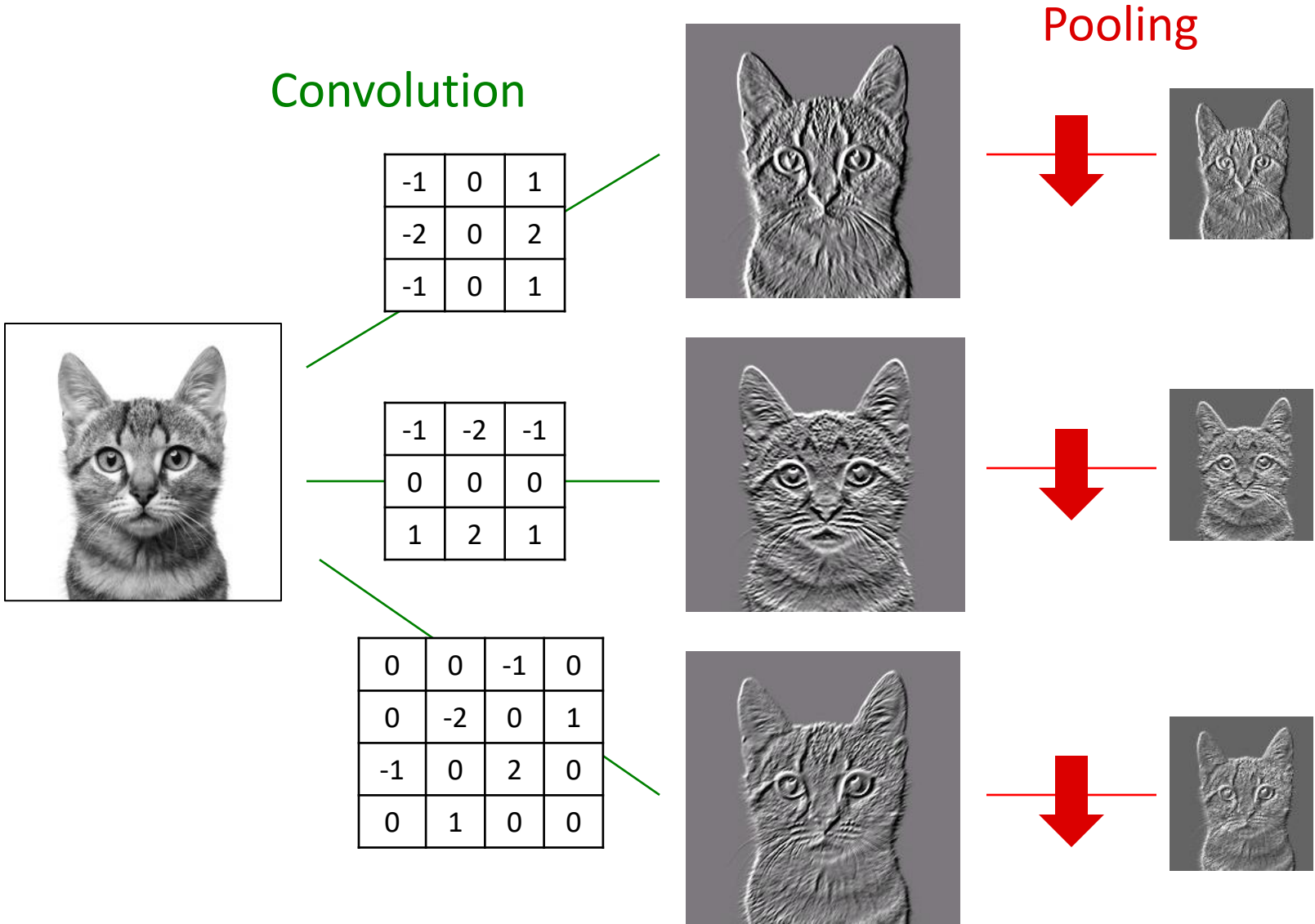
1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pool with 2x2 filters
and stride 2

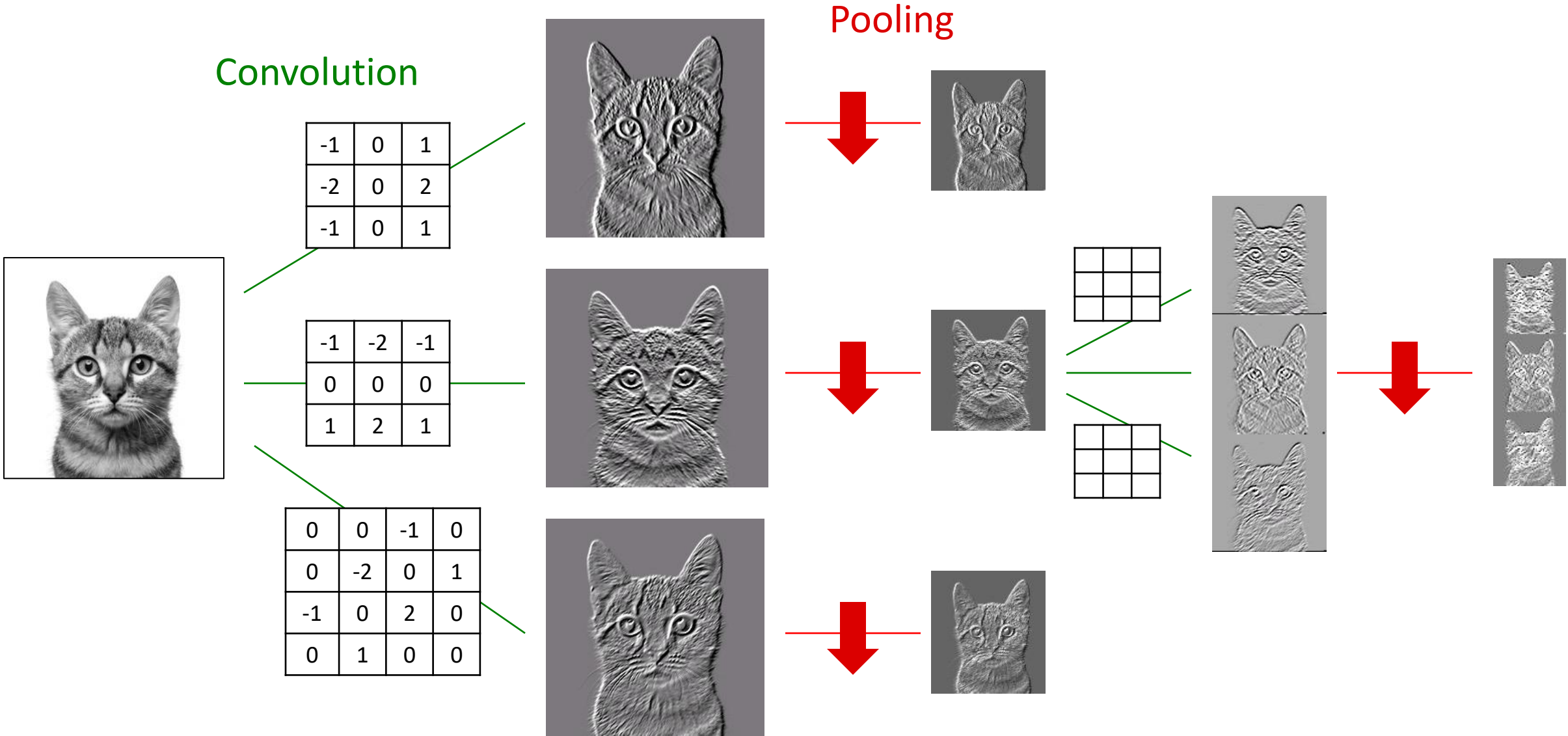


6	8
3	4

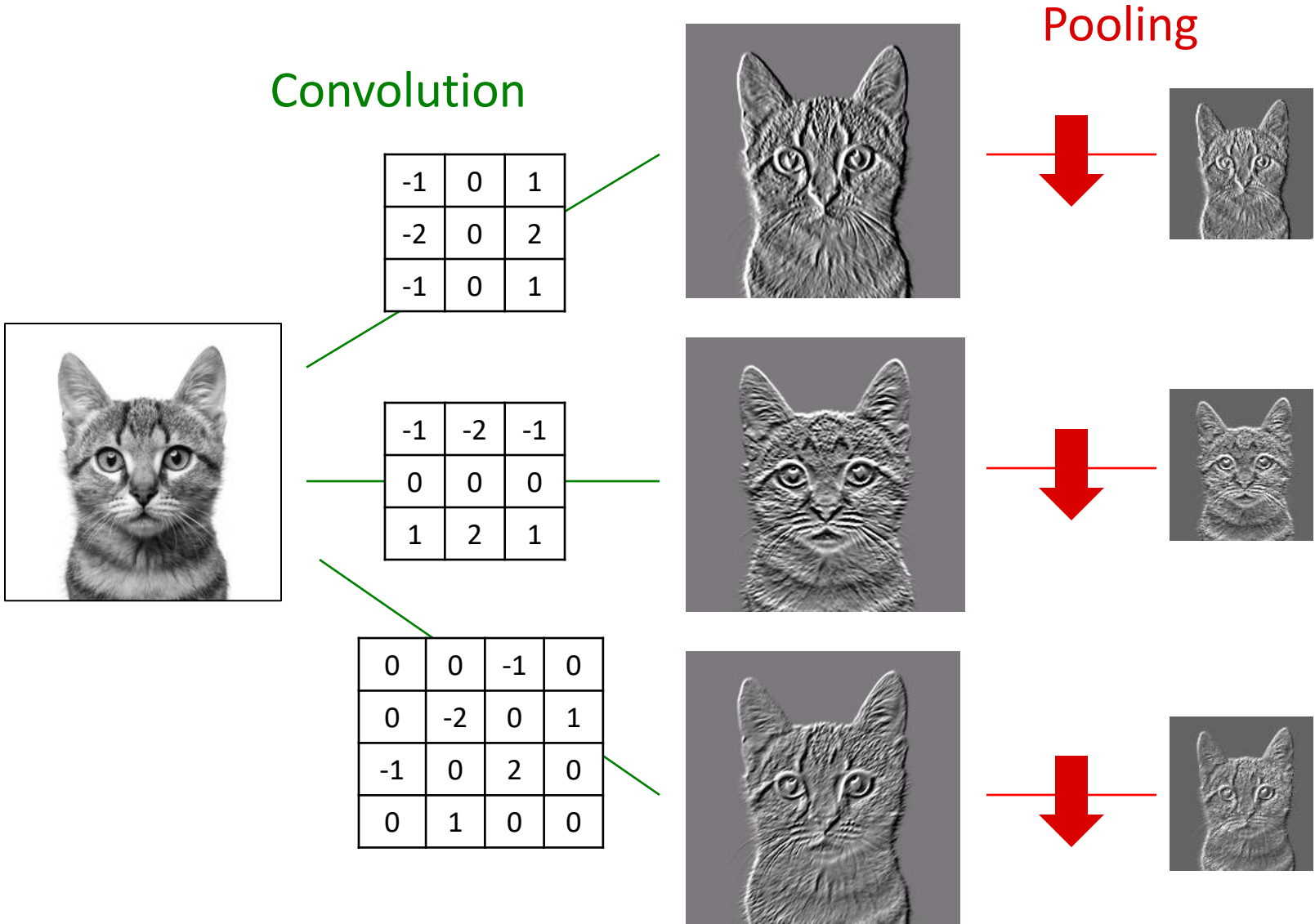
Convolutional Neural Networks



Convolutional Neural Networks



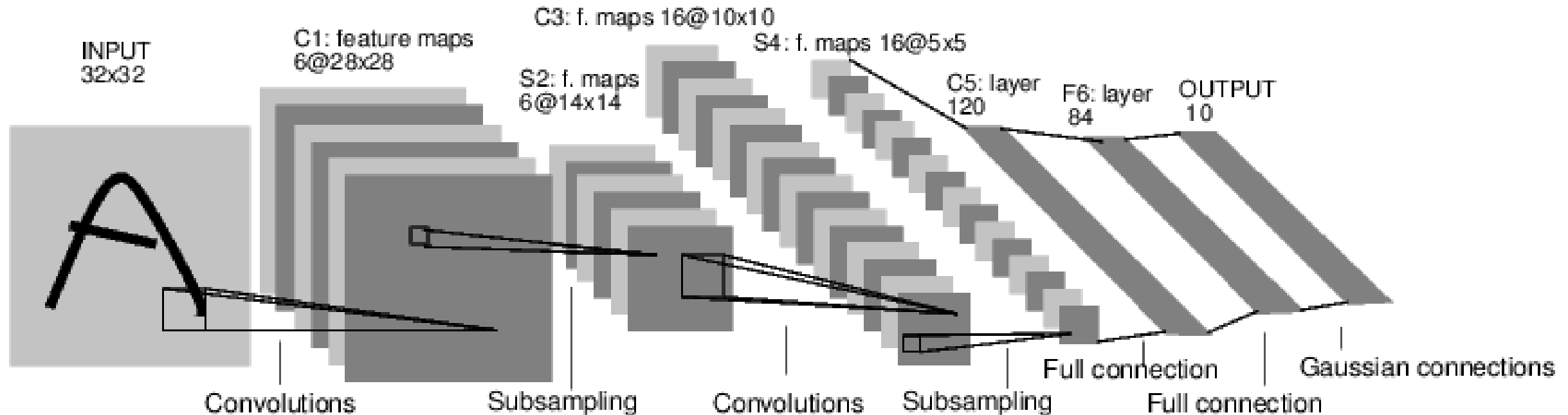
Convolutional Neural Networks



Convolutional Neural Networks

Lenet5 – Lecun, et al, 1998

- Convnets for digit recognition

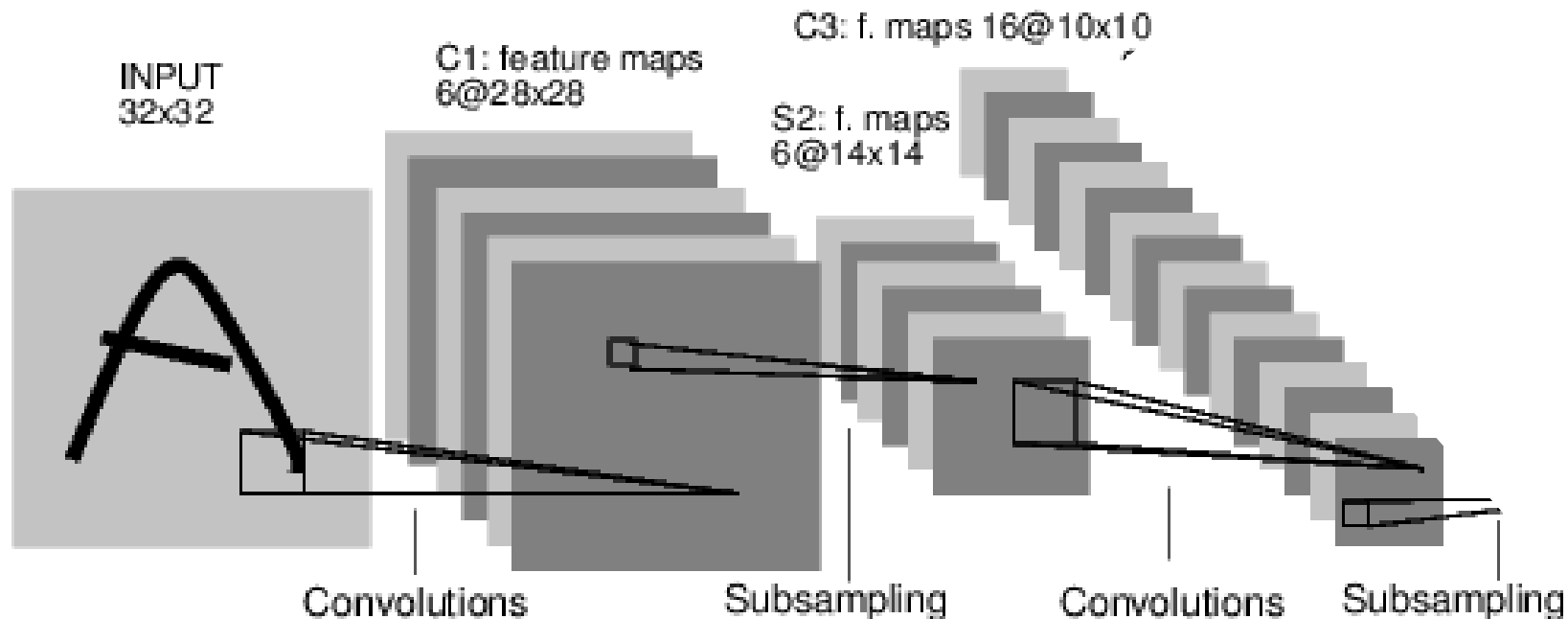


LeCun, Yann, et al. "Gradient-based learning applied to document recognition." Proceedings of the IEEE 86.11 (1998): 2278-2324.

Question:

How big many convolutional weights between S2 and C3?

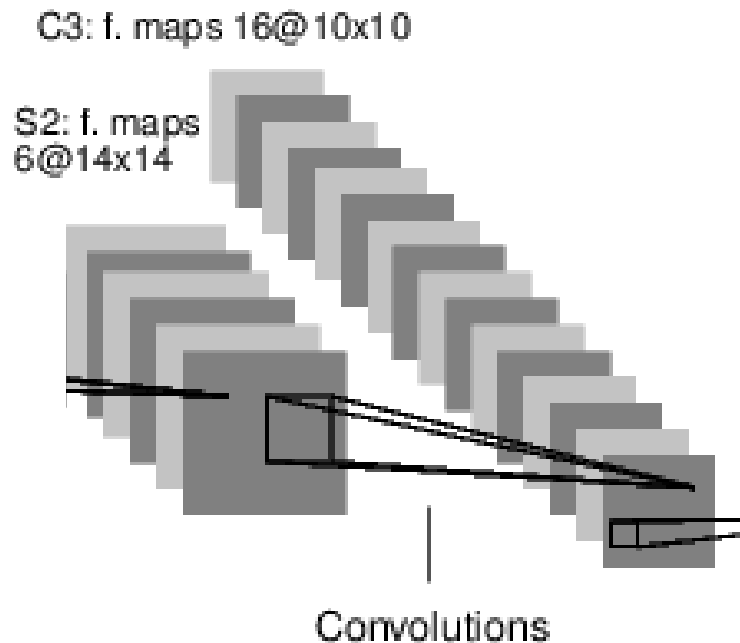
- S2: 6 channels @14x14
- Conv: 5x5, pad=0, stride=1
- C3: 16 channels @ 10x10



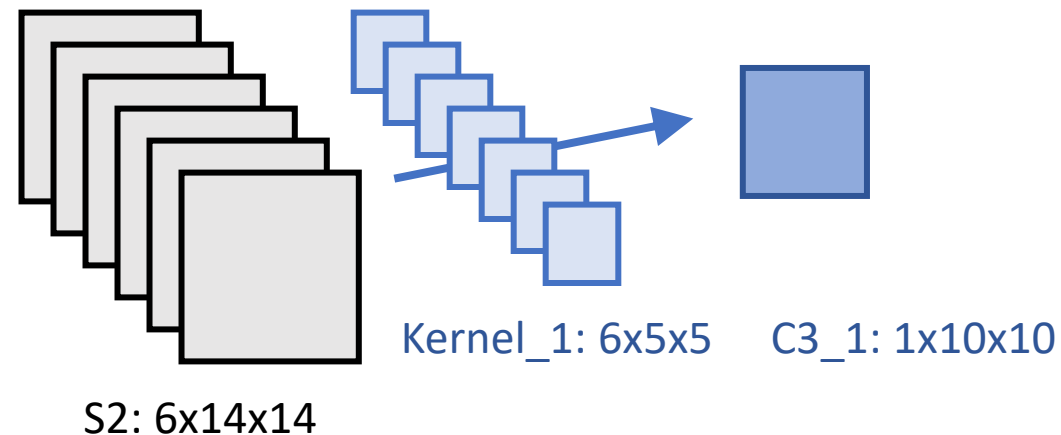
Question:

How big many convolutional weights between S2 and C3?

- S2: 6 channels @14x14
- Conv: 5x5, pad=0, stride=1
- C3: 16 channels @ 10x10



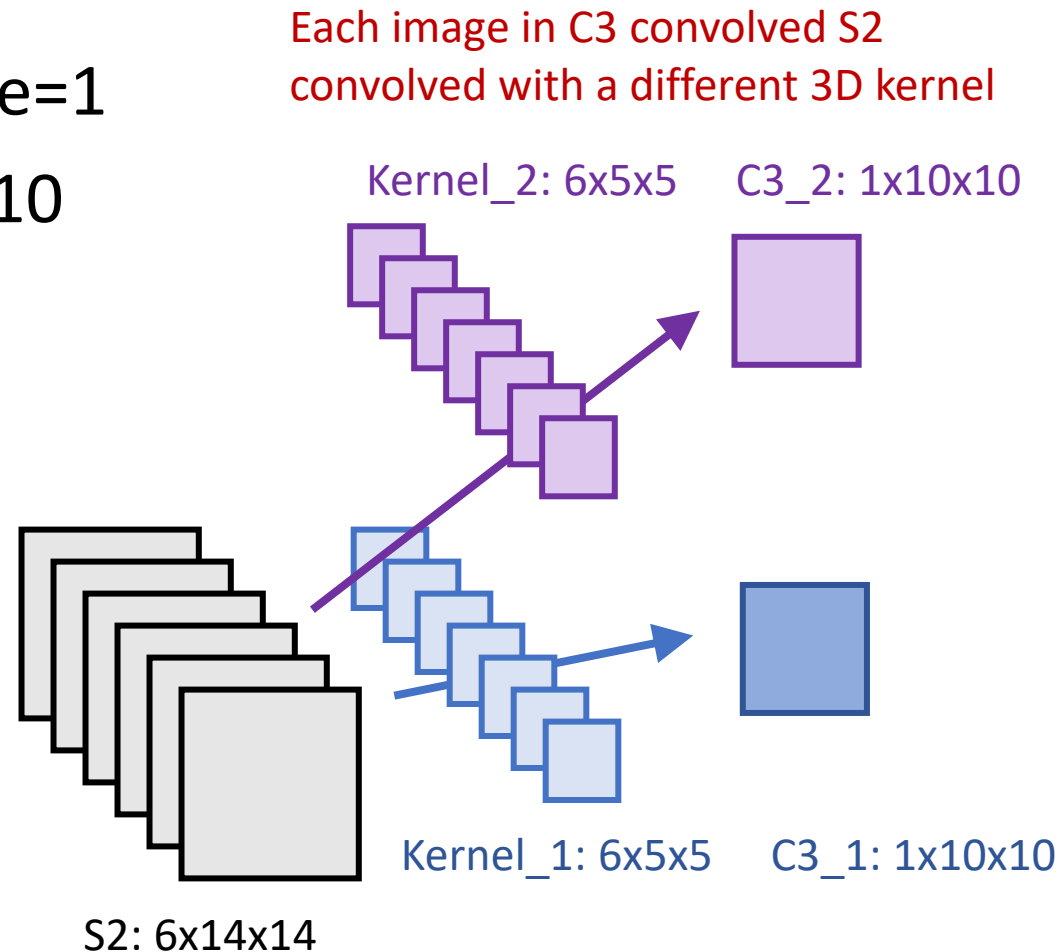
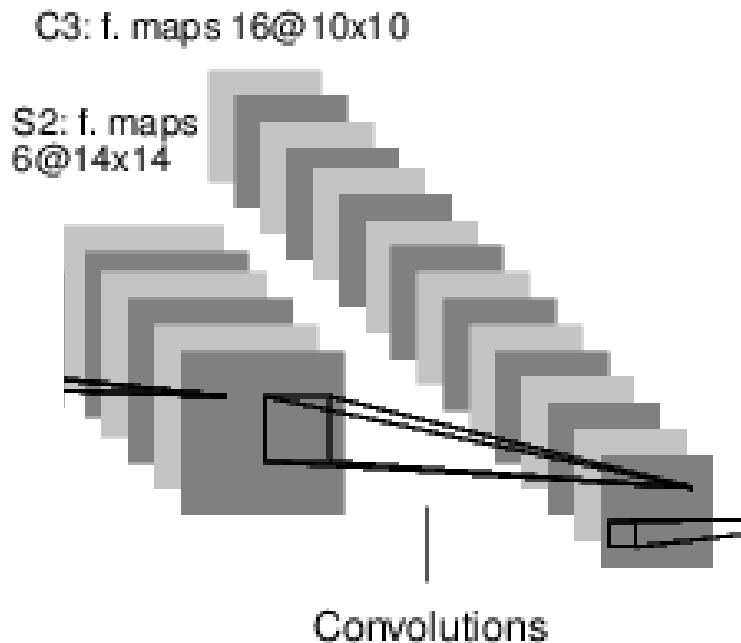
One image in C3 is actually the result of a 3D convolution



Question:

How big many convolutional weights between S2 and C3?

- S2: 6 channels @14x14
- Conv: 5x5, pad=0, stride=1
- C3: 16 channels @ 10x10



Question:

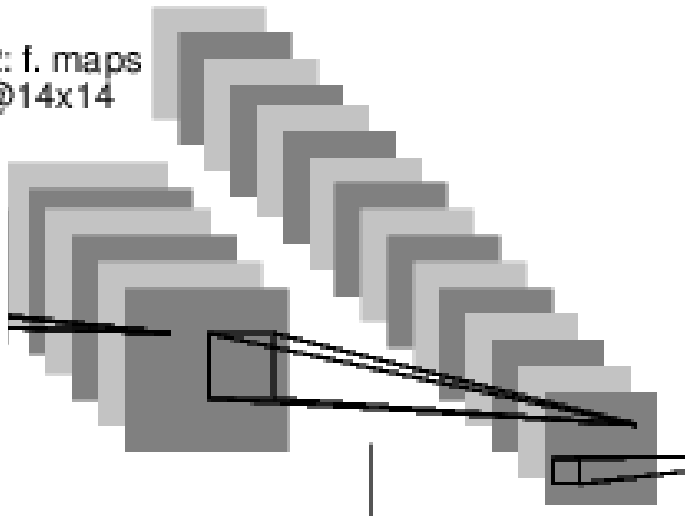
How big many convolutional weights between S2 and C3?

- S2: 6 channels @14x14
- Conv: 5x5, pad=0, stride=1
- C3: 16 channels @ 10x10

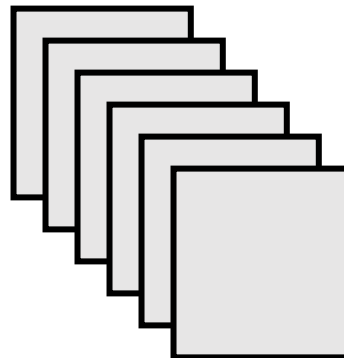
The 16 images in C3 are the result of doing 16 3D convolutions of S2 with 16 different 6x5x5 kernels. Assuming no bias term, this is 16x6x5x5 weights!

C3: f. maps 16@10x10

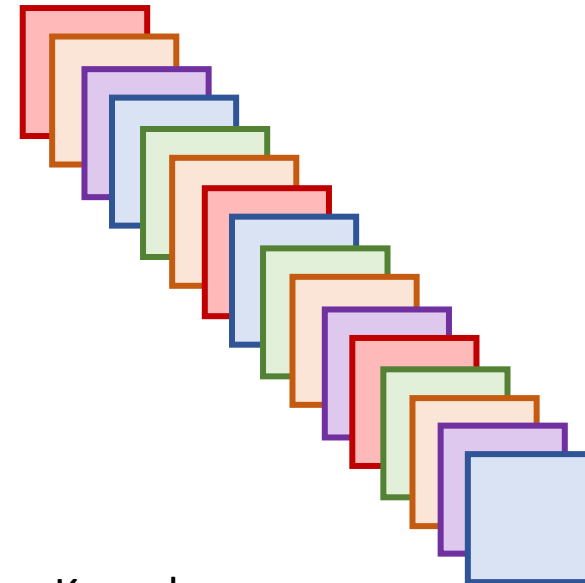
S2: f. maps 6@14x14



Convolutions



S2: 6x14x14



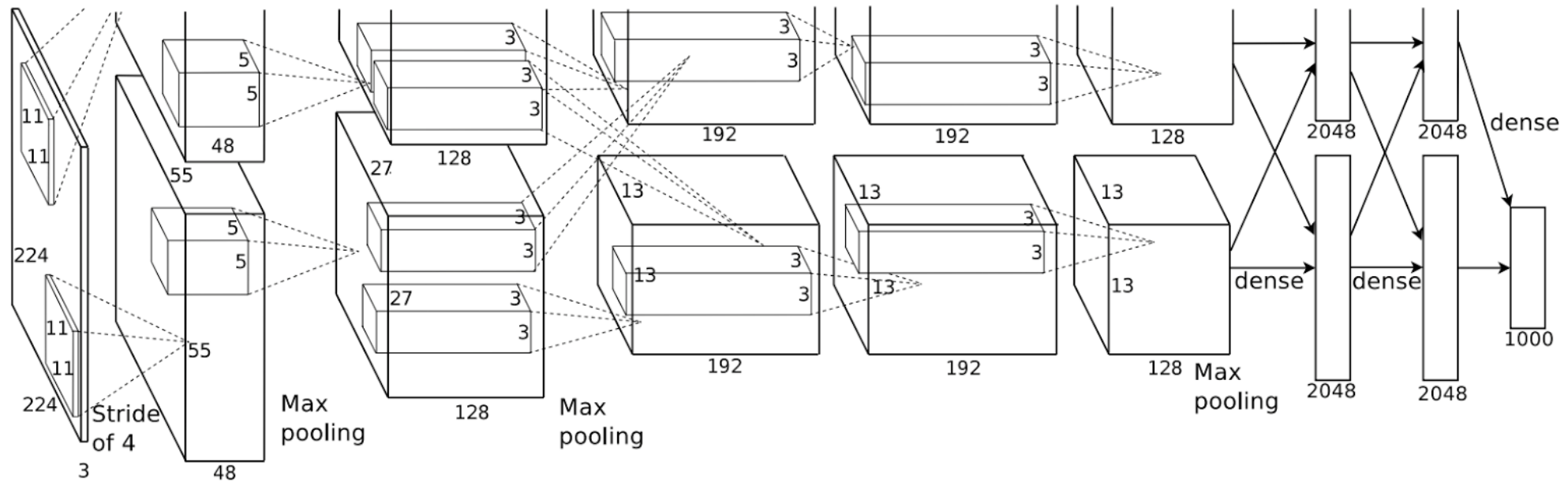
Kernels:
16@6x5x5

C3: 16@10x10

Convolutional Neural Networks

Alexnet – Lecun, et al, 2012

- Convnets for image classification
- More data & more compute power



Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "ImageNet classification with deep convolutional neural networks." NIPS, 2012.

CNNs for Image Recognition

