

BST 261: Data Science II

Lecture 7

Heather Mattie

Department of Biostatistics
Harvard T.H. Chan School of Public Health
Harvard University

April 9, 2018

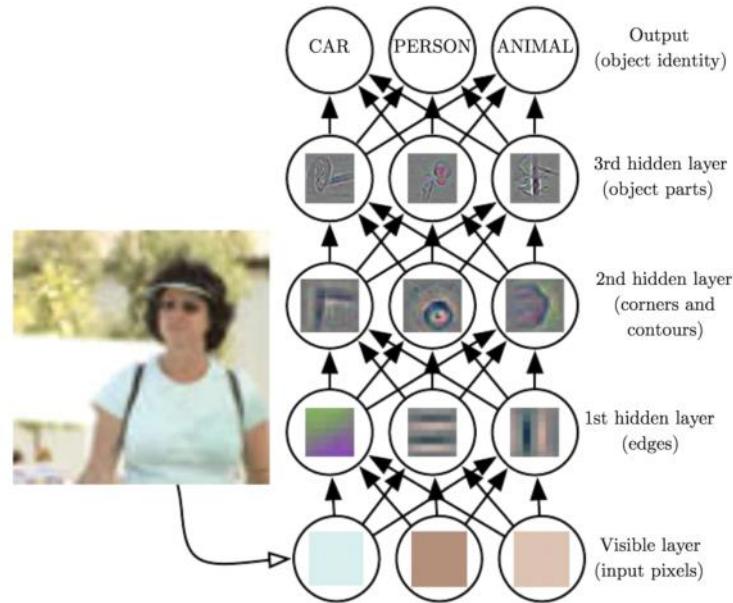
Convolutional Neural Networks (CNNs)

CNNs

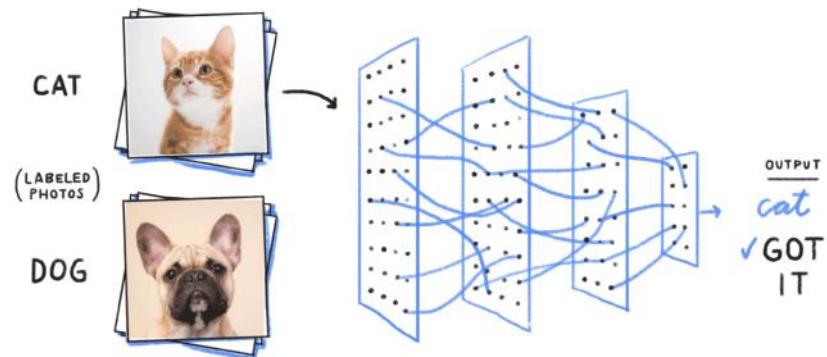
CNNs, or convnets for short, are at the heart of deep learning, emerging in recent years as the most prominent strain of neural networks in research. They have revolutionized computer vision, achieving state-of-the-art results in many fundamental tasks, as well as making strong progress in natural language processing, computer audition, reinforcement learning, and many other areas. CNNs have been widely deployed by tech companies for many of the new services and features we see today. They have numerous and diverse applications, including:

- Detecting and labeling objects, locations, and people in images
- Converting speech into text and synthesizing audio of natural sounds
- Describing images and videos with natural language
- Tracking roads and navigating around obstacles in autonomous vehicles
- Analyzing videogame screens to guide autonomous agents playing them
- “Hallucinating” images, sounds, and text with generative models

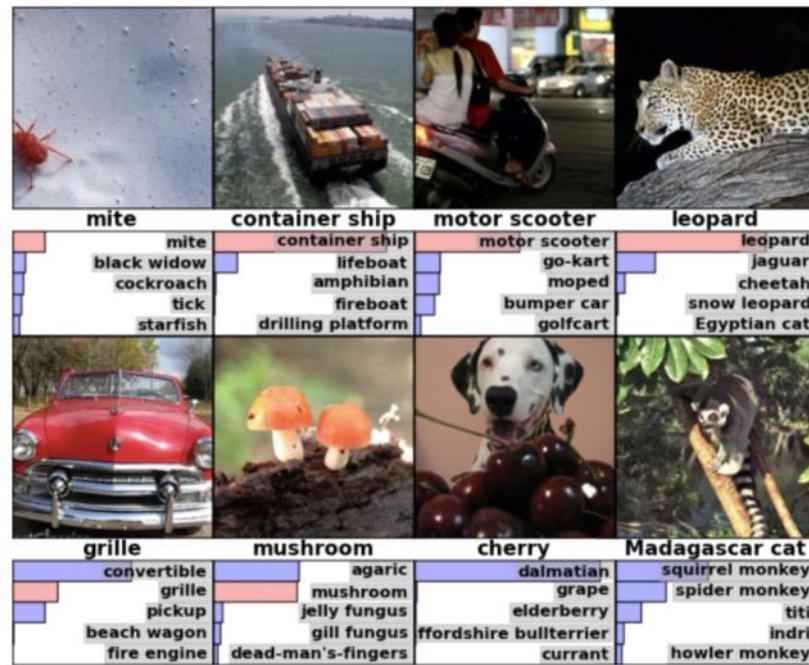
Computer Vision



Computer Vision

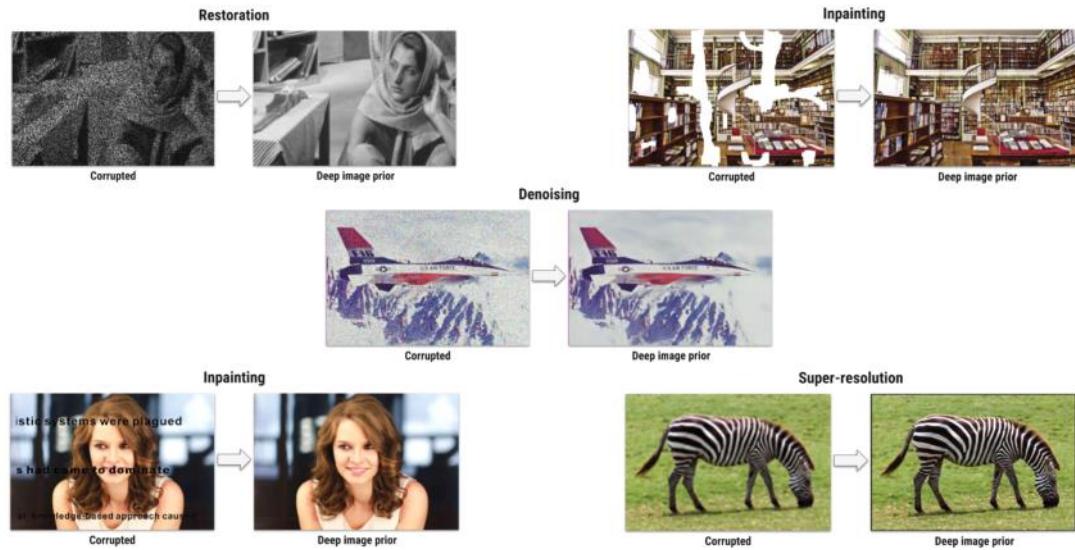


Computer Vision



Top: 4 correctly classified examples. Bottom: 4 incorrectly classified examples. Each example has an image, followed by its label, followed by the top 5 guesses with probabilities. From Krizhevsky *et al.* (2012).

Restoration, Inpainting, Denoising and Super-resolution



Autonomous Cars

- YouTube video: <https://www.youtube.com/watch?v=bD0nn0-4Nq8>

MNIST Example

```
1 import keras
2 from keras import layers
3 from keras import models
4 from keras.datasets import mnist
5 from keras.utils import to_categorical
6
7 (train_images, train_labels), (test_images, test_labels) = mnist.load_data()
8
9 train_images = train_images.reshape((60000, 28, 28, 1))
10 train_images = train_images.astype('float32') / 255
11
12 test_images = test_images.reshape((10000, 28, 28, 1))
13 test_images = test_images.astype('float32') / 255
14
15 train_labels = to_categorical(train_labels)
16 test_labels = to_categorical(test_labels)
```

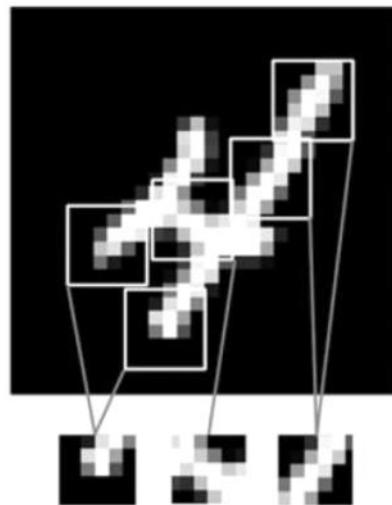
MNIST Example

new

```
1 model = models.Sequential()
2 model.add(layers.Conv2D(32, (3, 3), activation='relu', input_shape=(28, 28, 1)))
3 model.add(layers.MaxPooling2D((2, 2)))
4 model.add(layers.Conv2D(64, (3, 3), activation='relu'))
5 model.add(layers.MaxPooling2D((2, 2)))
6 model.add(layers.Conv2D(64, (3, 3), activation='relu'))
7
8 model.add(layers.Flatten())
9 model.add(layers.Dense(64, activation='relu'))
10 model.add(layers.Dense(10, activation='softmax'))
11
12 model.compile(optimizer='rmsprop',
13                 loss='categorical_crossentropy',
14                 metrics=['accuracy'])
15 model.fit(train_images, train_labels, epochs=5, batch_size=64)
```

CNNs vs Deep Feedforward Networks

- The fundamental difference between the feedforward networks we've seen and CNNs, is that dense layers learn **global** patterns while convolution layers learn **local** patterns



The Convolution Operation

- The primary purpose of the convolution operation is to extract features from the input
- To do this, the input is split into several different areas, and the convolution operation applied to each area
- A summary value is then calculated and kept as part of the output
- Here, I is the input image (sometimes called **feature map**), K is a two-dimensional $f \times g$ **kernel**, or more commonly, **filter**, and is a weighting function
- The output resulting from this operation is called the **response map**

$$S(i, j) = (I * K)(i, j) = \sum_f \sum_g I(i - f, j - g)K(f, g) \quad (1)$$

- Note: some neural networks implement a related function called the **cross-correlation**, but call it convolution. They are very similar and in some instances equivalent

The Convolution Operation

$$S(i,j) = (I * K)(i,j) = \sum_f \sum_g I(i-f, j-g) K(f,g)$$

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

input image (I): 5×5
central pixel position: $(0,0)$

1,1	1,2	1,3
2,1	2,2	2,3
3,1	3,2	3,3

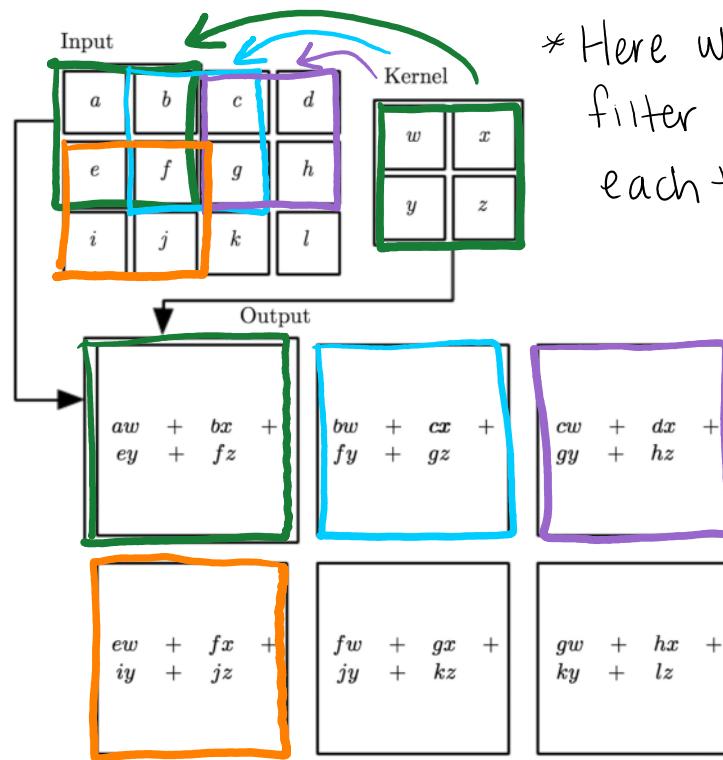
filter (K): 3×3

$S(1,1) : i=1, j=1$

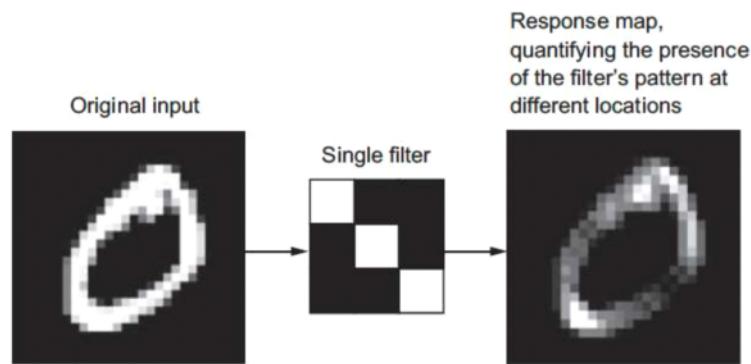
Output (S) : 3×3

$$\begin{aligned} S(1,1) &= \sum_{f=1}^3 \sum_{g=1}^3 I(i-f, j-g) K(f,g) = I(1-1, 1-1)K(1,1) + I(1-1, 1-2)K(1,2) + I(1-1, 1-3)K(1,3) \\ &\quad + I(1-2, 1-1)K(2,1) + I(1-2, 1-2)K(2,2) + I(1-2, 1-3)K(2,3) \\ &\quad + I(1-3, 1-1)K(3,1) + I(1-3, 1-2)K(3,2) + I(1-3, 1-3)K(3,3) \\ &= \underline{I(0,0)K(1,1)} + \underline{I(0,-1)K(1,2)} + \underline{I(0,-2)K(1,3)} \\ &\quad + \underline{I(-1,0)K(2,1)} + \underline{I(-1,-1)K(2,2)} + \underline{I(-1,-2)K(2,3)} \\ &\quad + \underline{I(-2,0)K(3,1)} + \underline{I(-2,-1)K(3,2)} + \underline{I(-2,-2)K(3,3)} \end{aligned}$$

The Convolution Operation



The Convolution Operation



Filters

- The most common filter sizes are 3×3 and 5×5
- Sometimes 7×7 filters are also used
- A 1×1 filter is a special case that we will see later in the course
- Odd dimension filters are preferred for computer vision
 - Natural padding (we'll see that in the following slides)
 - Ensures a 'central' position or pixel of the filter

Motivation

Convolution leverages 3 important ideas that can help improve a machine learning system:

- **Sparse interactions** / sparse connectivity / sparse weights
 - Filters are smaller than the input images and thus fewer parameters (weights) need to be stored
 - This reduces computational expense and improves statistical efficiency
- **Parameter sharing**
 - The same parameter is used for more than one function in the model
 - In a CNN, each element of the filter is applied to every position of the input
- **Equivariant representations** → during convolution stage
 - When the filter is shifted to a new patch of the image, the output from convolution will change accordingly.
 - Parameter sharing causes equivariance to translation: if the input changes, the output changes in the same way
 - A function $f(x)$ is invariant to a function g if $f(g(x)) = g(f(x))$
 - When processing time-series data, this means that convolution produces a sort of timeline that shows when different features appear in the input
 - Note that convolution is not equivariant to some other transformations, such as changes in the scale or rotation of an image
- Translation invariant: After learning a certain pattern from one part of an image, it can recognize it anywhere → during pooling stage
- Convolution also provides a way to work with inputs of different size

Motivation

- Example of both sparse interactions and parameter sharing:

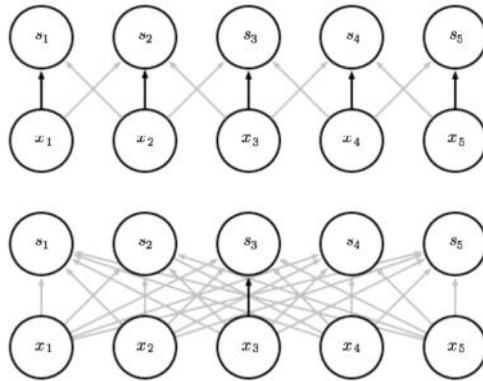
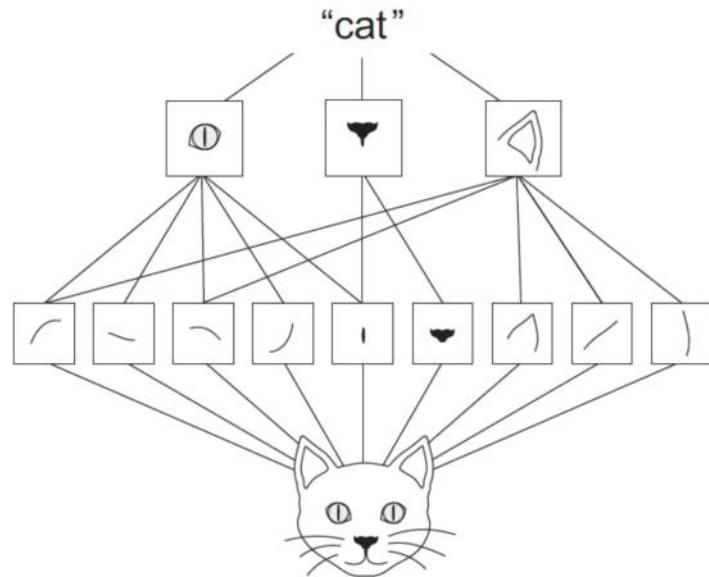


Figure 9.5: Parameter sharing. Black arrows indicate the connections that use a particular parameter in two different models. (Top)The black arrows indicate uses of the central element of a 3-element kernel in a convolutional model. Because of parameter sharing, this single parameter is used at all input locations. (Bottom)The single black arrow indicates the use of the central element of the weight matrix in a fully connected model. This model has no parameter sharing, so the parameter is used only once.

Motivation

- CNNs learn **spacial hierarchies** of patterns: one convolutional layer will learn small patterns and the next larger patterns made of the features of the layer before, and so on

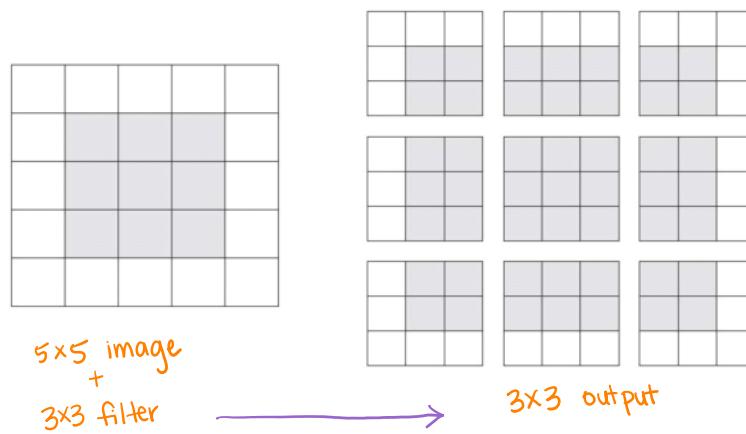


Padding

- Applying a filter to an input image shrinks it - the output dimensions are smaller than the input dimensions
- Additionally, when we perform the convolution operation on an input, the pixels on the edges aren't used as much as the pixels in the middle
- To make the amount of information used more equal across pixels, you can use **padding**
- Padding consists of adding an appropriate number of rows and columns to each side of the image (think a border of pixels)
- This enables an output with the same dimensions as the input

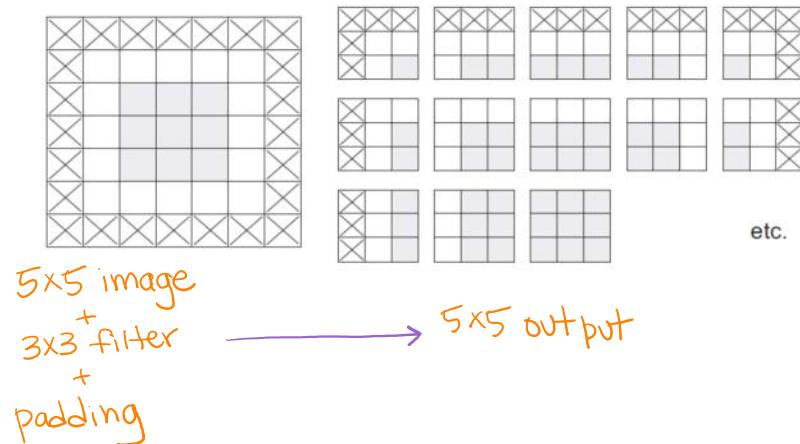
Valid Padding

- In Keras, the default is no padding, or '**valid**' padding
- This means the output will not be the same dimension as the input, and will instead depend on the dimension of the input and the size of the filter
- Below is a 5×5 image
- If we apply a 3×3 filter, the output will also be 3×3



Same Padding

- Below is the same 5×5 image, but with an added border
- If we use a 3×3 filter, we need to add $p = 1$ padding - here, p is the number of rows to add to the border of an image
- The output will then be 5×5
- Because the input and output have the same dimensions, this is called '**same**' padding



General Padding Formula

- There is a general formula that can help you decide how much padding you need or want
- Let the input be $n \times m$ and the filter $f \times f$
- Without padding, the output would be:

$$(n - f + 1) \times (m - f + 1) \quad (2)$$

- For an output with the same dimension as the input, need p such that:

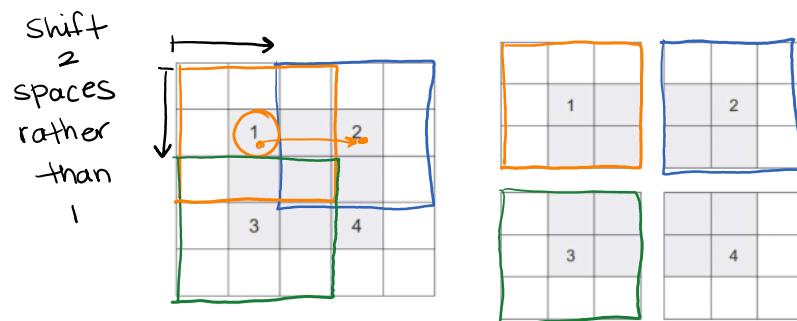
$$n \times m = (n + 2p - f + 1) \times (m + 2p - f + 1) \quad (3)$$

- Which boils down to $p = \frac{f-1}{2}$

Convolution Strides

- Another factor that can influence output size is convolution **strides**
- So far we have assumed that we slide our filter over a single space to extract a new patch - but what if we wanted to slide over 2 spaces, 3 spaces, or more?
- Convolutional strides are convolutions with a stride greater than 1
- If your stride is set equal to 2, you will downsample the width and height of the input image by a factor of 2
- If s is the size of the stride, the output will have dimensions

$$\left\lfloor \frac{n + 2p - f}{s} + 1 \right\rfloor \times \left\lfloor \frac{m + 2p - f}{s} + 1 \right\rfloor \quad (4)$$



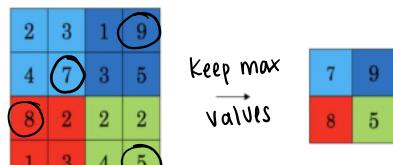
- Strided convolutions are rarely used in practice - **max-pooling** is more common

Pooling

- A typical layer of a CNN consists of 3 stages:
 - Convolution stage - linear transformation of the input
 - Detector stage - Nonlinear activation (ex: `relu` activation function is applied)
 - Pooling stage - further modification (downsizing) of the output
- A pooling function replaces the output at a certain location with a **summary statistic** of the nearby outputs
- Pooling greatly reduces the computational expense of the network by decreasing the number of parameters to be learned
- There are different types of pooling
 - Max pooling
 - Outputs the maximum value from a patch for each channel
 - Similar to convolution, but instead of transforming patches via a learned linear transformation, they're transformed via a hardcoded `max` tensor operation
 - Very common
 - Usually done with 2×2 windows
 - Average pooling
 - Weighted average pooling
 - Based on the distance from the central pixel
- Max pooling tends to work better than other pooling functions and convolutional strides

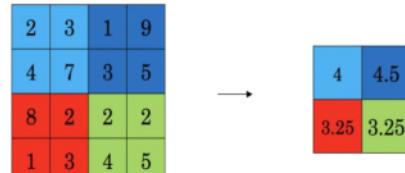
Pooling Functions

Max Pool



Max-Pool with a 2 by 2 filter and stride 2.

Average Pool



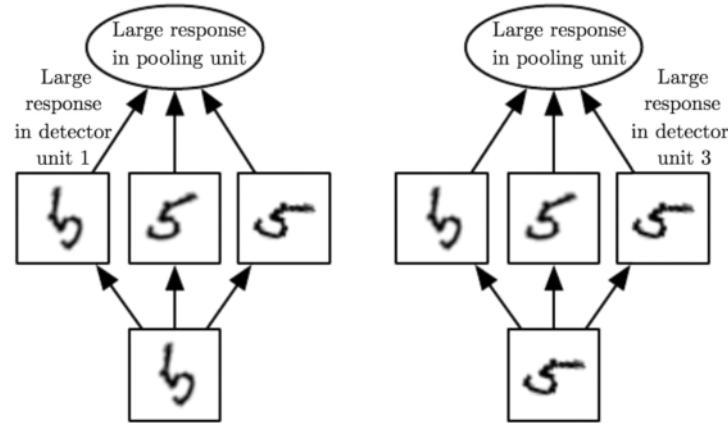
Invariant to translation:

2	3	1	9
4	7	3	5
8	2	2	2
1	3	4	5

$\max(\square) = 7$] the input has been translated, but the
 $\max(\square) = 7$ output hasn't changed.
 → This detects if a feature is present,
 not exactly where the feature is.

Pooling and Invariance

- A pooling unit that pools over multiple features that are learned with separate parameters can learn to be invariant to transformations of the input
- Here, a set of 3 learned filters and a max pooling unit can learn to become invariant to rotation

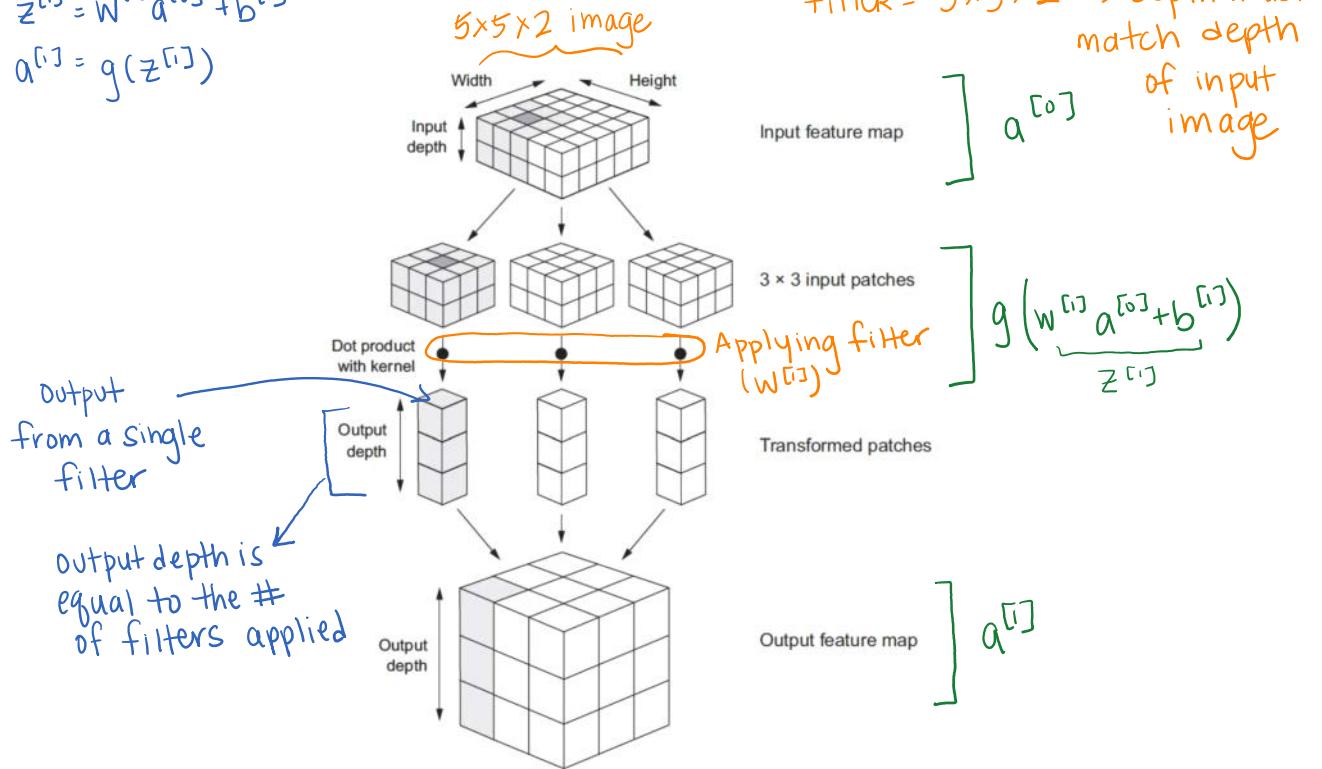


Terminology

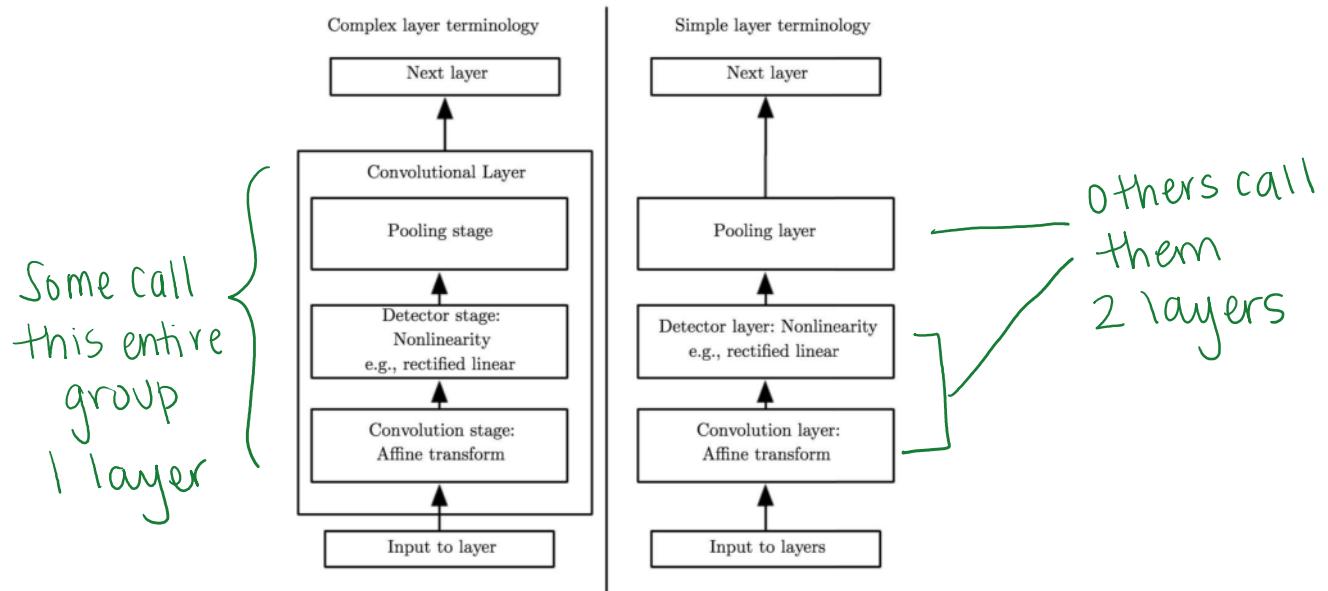
- Convolutions operate over 3D tensors with two spatial axes, **height** and **width**, as well as a **depth** axis (or **channels** axis)
- For a color (RGB) image, the depth is equal to 3
- For black and white (greyscale) images, the depth is equal to 1
- The convolution operation extracts different **patches** from the input image and applies the same transformation to each of them, resulting in a response map that is also a 3D tensor
- The response map has a width, height, and depth, all of which depend on the input image, filter, padding and stride
- Convolutions are defined by 2 key parameters:
 - Size of the filter
 - Depth of the output response map, i.e., how many filters are applied to the input
- Convolution works by sliding the filter over the 3D input image, stopping at every possible location, and extracting the 3D patch at each location
- Each 3D patch is transformed into a 1D vector
- All 1D vectors are then reassembled into a 3D output map

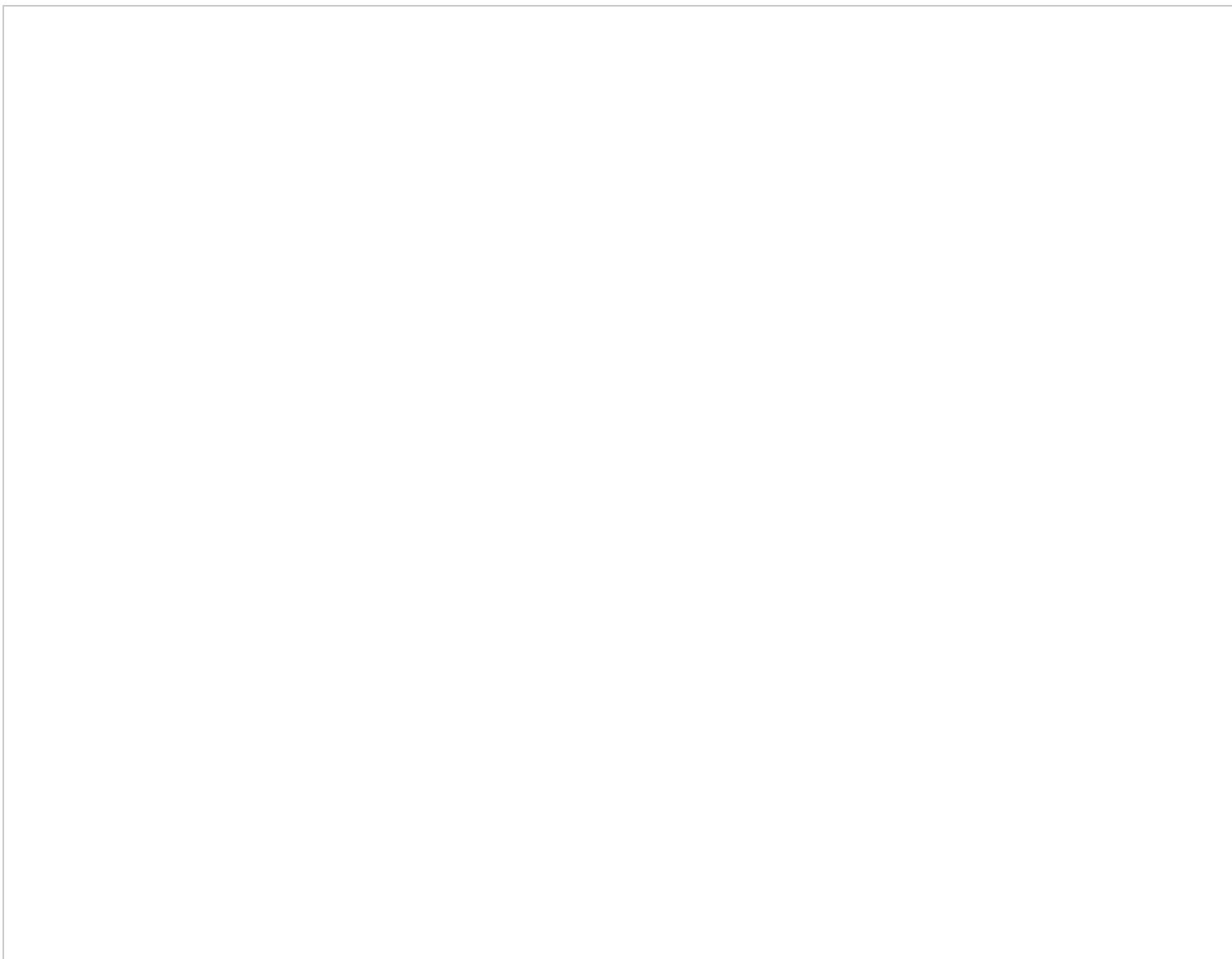
Convolution Schematic

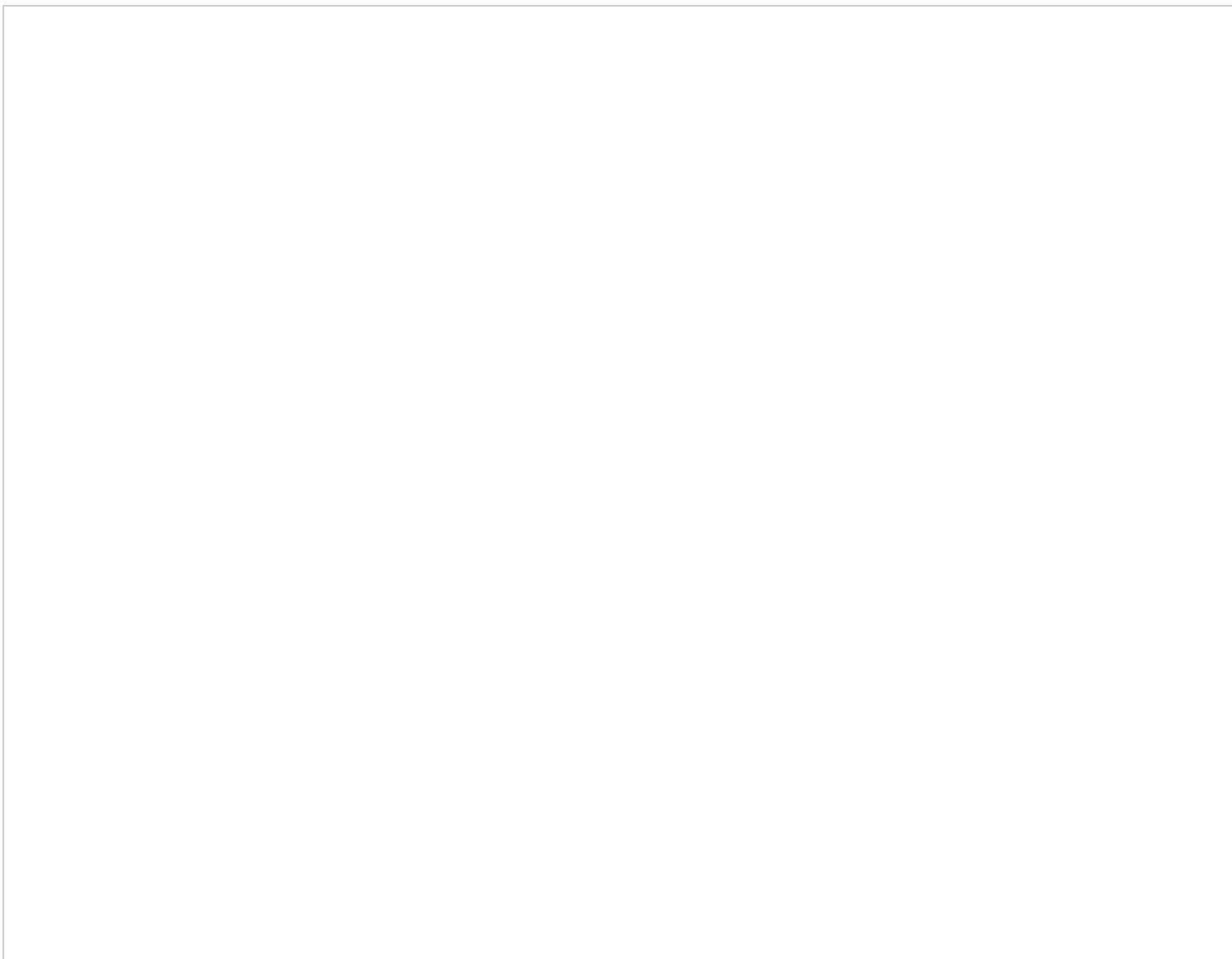
Recall: $z^{[i]} = W^{[i]} a^{[0]} + b^{[i]}$
 $a^{[i]} = g(z^{[i]})$



Layer Terminology







Revisiting the MNIST Example

```
1 # Still want a sequential, stack of layers
2 model = models.Sequential()
3 # Convolutional layer with 32 filters that are 3x3, relu activation function
4 model.add(layers.Conv2D(32, (3, 3), activation='relu', input_shape=(28, 28, 1)))
5 # Pooling layer with 2x2 windows
6 model.add(layers.MaxPooling2D((2, 2)))
7 # Convolutional layer with 64 filters that are 3x3, relu activation function
8 model.add(layers.Conv2D(64, (3, 3), activation='relu'))
9 # Pooling layer with 2x2 windows
10 model.add(layers.MaxPooling2D((2, 2)))
11 # Convolutional layer with 64 filters that are 3x3, relu activation function
12 model.add(layers.Conv2D(64, (3, 3), activation='relu'))
13
14 # Collapse the 3D tensor
15 model.add(layers.Flatten())
16 # Fully connected layer with 64 hidden units, relu activation function
17 model.add(layers.Dense(64, activation='relu'))
18 # Softmax output function with 10 classes
19 model.add(layers.Dense(10, activation='softmax'))
20
21 # Same optimizer, loss function and performance measure as before
22 model.compile(optimizer='rmsprop',
23                 loss='categorical_crossentropy',
24                 metrics=['accuracy'])
25 model.fit(train_images, train_labels, epochs=5, batch_size=64)
```

Revisiting the MNIST Example

```
>>> model.summary()
Layer (type)                  Output Shape       Param #
=====
conv2d_1 (Conv2D)             (None, 26, 26, 32)  320
maxpooling2d_1 (MaxPooling2D) (None, 13, 13, 32)  0
conv2d_2 (Conv2D)             (None, 11, 11, 64)  18496
maxpooling2d_2 (MaxPooling2D) (None, 5, 5, 64)   0
conv2d_3 (Conv2D)             (None, 3, 3, 64)   36928
flatten_1 (Flatten)           (None, 576)        0
dense_1 (Dense)               (None, 64)         36928
dense_2 (Dense)               (None, 10)         650
=====
Total params: 93,322
Trainable params: 93,322
Non-trainable params: 0

>>> test_loss, test_acc = model.evaluate(test_images, test_labels)
>>> test_acc
0.9908000000000001
```

Revisiting the MNIST Example

