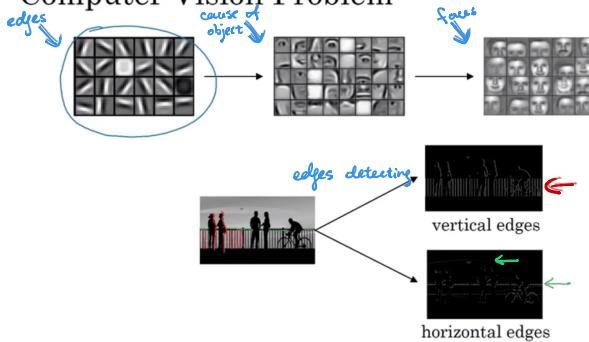


Convolutional Operations

The convolution operation is one of the fundamental building blocks of a convolutional neural network. To see how convolutional operation works, let's use an example of edge detection.

We have talked about how early layers of a NN might edges, then coarse of objects later, and finally complete objects.

Computer Vision Problem



Andrew Ng

Let's see how we can detect vertical edges in an image by using a simple example of a 6 by 6 grayscale image.

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	3
2	4	5	2	3	9

↑ image

(6 by 6 by 1 because it's greyscale)

$$\begin{array}{c} * \\ \uparrow \\ \text{convolution} \end{array} \quad \begin{array}{c} \text{filter/kernel} \\ (3 \times 3) \end{array} = \begin{array}{c} \begin{pmatrix} 3 & 0 & 1 \\ 1 & 5 & 8 \\ 2 & 7 & 2 \end{pmatrix} * \begin{pmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{pmatrix} \\ \begin{pmatrix} -5 & -4 & 0 & 8 \\ \vdots & \vdots & \vdots & \vdots \\ -5 & -4 & 0 & 8 \end{pmatrix} \end{array}$$

↑ edge detector (4 by 4)

(do not mistake it for multiplication, or element-wise multiplication)

python: conv-forward, tensorflow: tf.nn.conv2d, keras: Conv2D)

$$\begin{pmatrix} 3 & 0 & 1 \\ 1 & 5 & 8 \\ 2 & 7 & 2 \end{pmatrix} * \begin{pmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{pmatrix} = \begin{pmatrix} 3 \times 1 + 0 \times 0 + 1 \times -1 \\ + 1 \times 1 + 5 \times 0 + 8 \times -1 \\ + 2 \times 1 + 7 \times 0 + 2 \times -1 \end{pmatrix} = -5$$

So why does the matrix on the right turn out to be vertical edge detector? Let's use an even simpler example:

note:
the dimensions here seem a little bit wrong
that the detected edge seems thicker. That's
only because the image on the left is very small.
For large image (e.g. 1000 by 1000), this does
a really great job.

There are more edge detectors than the above one, including positive and negative edges (i.e. difference between light to dark versus dark to light edge transitions), as well as algorithm-learned edge detector instead of the hand-coded one like we did.

- Positive and negative edge:

→ same as above

↑ represent light to dark edge

→ reversed

↑ represent dark to light edge

Andrew Ng

- Horizontal edge:

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

** flip it 90°*

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

* note: applicable to other filters as well, i.e. flip the vertical filter 90 degrees to get the horizontal filter.

- Other hand-picked filters

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

more weight
to the central
pixel

3	0	-3
10	0	-10
3	0	-3

- Algorithm - learned filter

treat these 9 numbers as parameters, and let the algorithm learn them by back-propagation:

W ₁	W ₂	W ₃
W ₄	W ₅	W ₆
W ₇	W ₈	W ₉

→ other than detecting vertical/horizontal edges,
it's also able to edges at 45°, 75° etc.

But before we go about using back propagation to learn these 9 numbers, let's first talk about some other details that are important in building CNNs, including padding and different strides for convolutions.

Padding

In order to build deep CNN, one modification to the basic convolutional operation is padding.

Basic Convolutional Operation

The original image has shrunk from 6x6 to 4x4 and after a few convolutions it might shrink to 1 by 1.
 Pixels on the corners are used much less in the output, i.e. throwing away a lot of information on the edge of the image.
 the size of the output is determined by: $n-f+1 \times n-f+1$ (n : original size, f : filter size)

padding with 1 broader ($p=1$)

$\frac{6 \times 6}{n+2p-f+1 = 6}$

How much to pad → 2 common choices:

- Valid padding : no padding

$$n \times n \text{ image} * f \times f \text{ filter} = n-f+1 \times n-f+1 \text{ output}$$

- Same padding : pad so that output size is the same as the input size

$$n+p \times n+p \text{ image} * f \times f \text{ filter} = n \times n \text{ image}$$

$$\text{where } p: n+p-f+1 = n \Rightarrow p = \frac{f-1}{2}$$

when f is odd, we can make sure that output size is preserved.
 By convention, f is odd, because if f is even we would have some asymmetric padding to have same padding; second, odd filter has a central position

• Strided Convolutions

Strided convolution is another piece of basic building block of convolutions as used in CNNs.

In earlier convolutions, we take the boxes in the matrix sequentially, but in strided convolutions, we step the boxes by 2 steps:

$\begin{matrix} 2 & 3 & 7 & 4 & 6 & 2 & 9 \\ 6 & 6 & 9 & 8 & 7 & 4 & 3 \\ 3 & 8 & 3 & 8 & 9 & 7 \end{matrix} \quad * \quad \begin{matrix} 3 & 4 & 4 \\ 1 & 0 & 2 \\ -1 & 0 & 3 \end{matrix} = \begin{matrix} 21 & 9 & 8 \\ 160 & 83 \\ 69 & 91 & 127 \\ 44 & 72 & 74 \end{matrix}$

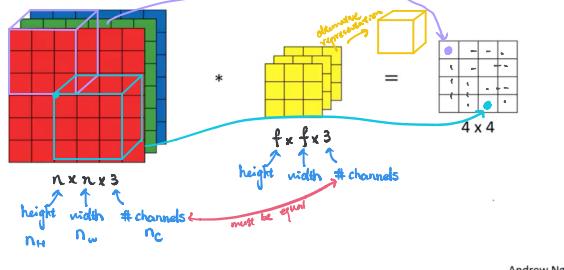
$n \times n$ image with p padding 3×3 $f \times f$ filter with S strides 3×3
determined by: $\left\lfloor \frac{n+2p-f+1}{s} + 1 \right\rfloor = \text{floor} \left(\frac{n+2p-f}{s} - 1 \right) + 1$
so $\frac{7+0-3}{2} + 1 = 3$

ensures that if the box steps outside of the matrix, it will not do that computation, i.e. the box is fully contained within the image plus padding.

note: so far we call this operation as convolution operations, but really it should be cross-correlation. In convention we call it convolution operation.

• Convolutions over volume

So far we've only talk about greyscale image. Let's see how convolution works on RGB images.



→ convolutions operations of each corresponding channel (i.e. R and 1st channel in filter, G & 2nd channel of filter, B & 3rd channel of filter) and then add them up together

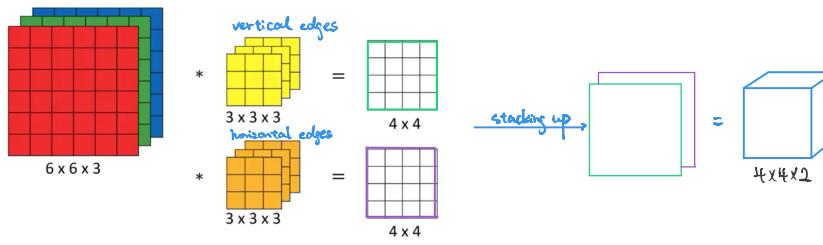
→ detect edges in red colour:

detect edges in any colour:

Andrew Ng

By stacking up output matrices, we can detect more features at the same time. For example, if we want to detect vertical edges and horizontal edges at the same time:

Multiple filters



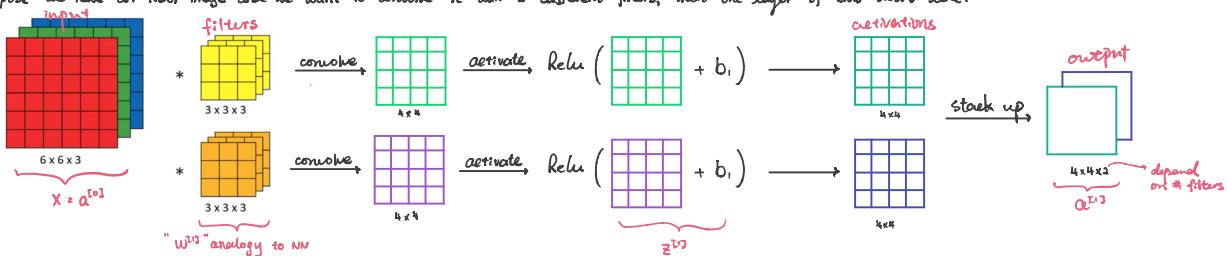
Andrew Ng

→ this example has neither padding nor stride. If its done with padding or stride, output size is determined by: $\left\lfloor \frac{n+2p-f+1}{s} + 1 \right\rfloor$

→ # channels of output is determined by the # filters used, so it could be up to hundreds of channels depending on how many features we want to detect at that layer

• CNN structure

Suppose we have an RGB image and we want to convolve it with 2 different filters, then one layer of CNN looks like:



So in short, $X = \begin{matrix} 6 \times 6 \times 3 \\ \xrightarrow{\text{3 layers}} \\ 4 \times 4 \times 3 \end{matrix} = a^3$

Andrew Ng

parameters = $(\text{width} \times \text{height} \times \# \text{channels} + 1) \times \# \text{filters} \xrightarrow{\text{+ bias}}$ makes CNN less prone to overfitting as compared to the images, # parameters is very small and fixed

Notation: if layer l is a convolutional layer,

$f^{[l]}$ = filter size

$p^{[l]}$ = # padding

$s^{[l]}$ = # stride

$n_h^{[l]}$ = # filters

$n_h^{[l]} / n_w^{[l]}$ = height / width

$$\text{for } n_h^{[l]} \text{ or } n_w^{[l]}, n_h^{[l]} = \frac{n_h^{[l-1]} + 2p^{[l-1]} - f^{[l-1]}}{s^{[l-1]}} + 1$$

SIZES:

$$\text{input: } n_h^{[L-1]} \times n_w^{[L-1]} \times n_c^{[L-1]}$$

$$\text{output: } n_h^{[L]} \times n_w^{[L]} \times n_c^{[L]}$$

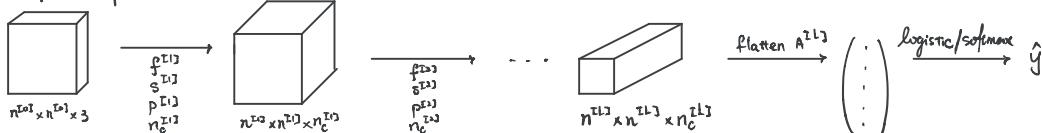
each filter is $f^{[L-1]} \times f^{[L-1]} \times n_c^{[L-1]}$ → to match input

activation matrix $A^{[L]} : m \times n_h^{[L]} \times n_w^{[L]} \times n_c^{[L]}$ stack of output

weights: $f^{[L-1]} \times f^{[L-1]} \times n_c^{[L-1]} \times n_c^{[L]}$ stack of filters

bias: $n_c^{[L]}$ or $1 \times 1 \times n_c^{[L]}$ stack of bias

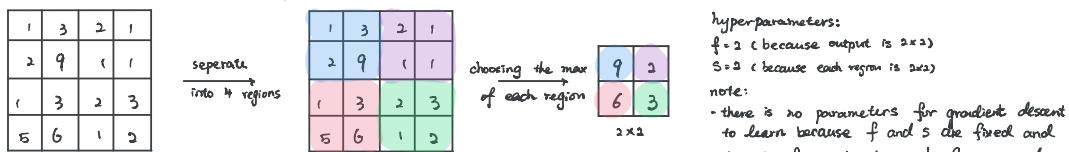
An simple example of CNN:



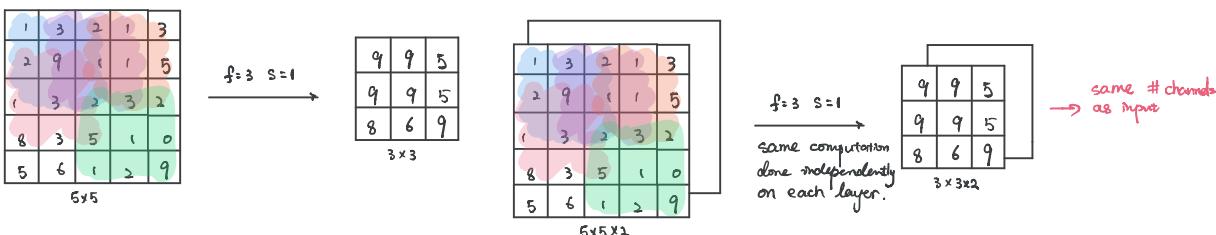
So far we have only talked about one type of layer in a CNN, which is convolution. There are other types of layer as well: pooling and fully-connected.

- pooling: can reduce size of the representation, to speed up computation as well as make some of the features that detects a bit more robust.

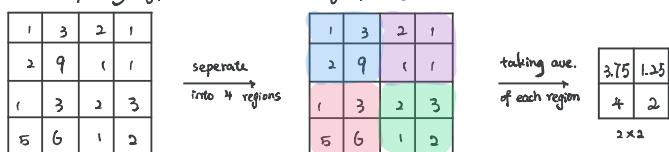
suppose we have a 4x4 input, max pooling is as follow:



another example,

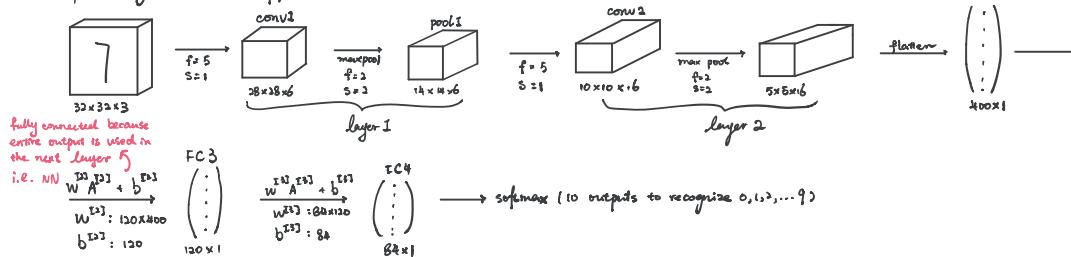


another pooling type known as average pooling is used less often.



- fully-connected

inspired by LeNet-5, suppose we have a CNN:



Why is CNN so useful can be mainly explained by 2 reasons:

- parameter sharing: a feature detector (e.g. vertical edge detector $\begin{pmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{pmatrix}$) that's useful in one part of image is probably useful in other part of the image \Rightarrow a detector be used multiple times in one layer \Rightarrow reduce # parameters used
- sparsity of connections: in each layer, each output value depends only on a small number of input (outputs don't affect each other)