

Course 4, week 1: Foundations of Convolutional Neural Networks

Computer vision

- rapid advances in computer vision are enabling brand new applications that weren't possible a few years ago.
 - ↳ includes image classification, object detection, neural style transfer (transforming a content image in a certain artistic style).
- one of the challenges of computer vision problems is that the inputs can get really big.
 - ~ color channels
 - ↳ $64 \times 64 \times 3 = 12288$ ~ input feature size for example
 - ↳ larger images: $1000 \times 1000 \times 3 = 3\text{ million}$ input features
 - ↳ if we utilize 1000 hidden units in our first hidden layer, our weight matrix $w^{(1)} = (1000, 3\text{million})$ in a standard fully connected network
 - ↳ the matrix would have 3 billion parameters.
 - ↳ This is computationally expensive and it's unlikely that we'll find enough data to train our network.
- To scale up, we must use the convolution operation.

Edge Detection Example

6x6

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	8
2	4	5	2	3	9

Convolution operator

*

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

=

-5	-4	0	8
-10	-2	2	3
0	-2	-4	-7
-3	-2	-3	-16

3x3 : filter / Kernel

4x4

i. take the element-wise product of the matrix and the filter overlaid

on top of it: $3 \times 1 + 1 \times 1 + 2 \times 1 + 0 \times 0 + 5 \times 0 + 7 \times 0 + 1 \times -1 + 8 \times -1 + 2 \times -1 = -5$

ii. shift the overlaid square one step to the right and repeat

the process: $0 \times 1 + 1 \times 0 + 2 \times -1 + 5 \times 1 + 8 \times 0 + 9 \times -1 + 7 \times 1 + 2 \times 0 + 5 \times -1 = -4$

... and you continue as such ...

... after shifting to the right, we begin shifting down

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

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

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

... and we continue until we have covered all the elements of the matrix ...

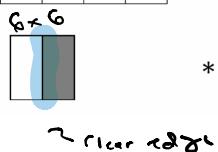
* By convolving in this way, we obtain a vertical edge detection

* Why does this work as a vertical edge detector?

10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0

$$\begin{matrix} * & \begin{matrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{matrix} \\ & 3 \times 3 \end{matrix}$$

$$\begin{matrix} = & \begin{matrix} 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \end{matrix} \\ & 4 \times 4 \end{matrix}$$



$$\begin{matrix} * & \begin{matrix} \text{filter} \end{matrix} \end{matrix}$$



* the white region in the middle corresponds to the detected vertical edge.
Andrew Ng

- the line in the resulting image is thick because we are working with a 6x6 image of low resolution. If our image were 1000x1000 instead, our filter would do a more granular job in detecting the vertical edges of the image.

more edge detection

- positive edge: light to dark vs. negative edge: dark to light

$$\begin{array}{c}
 \begin{array}{|c|c|c|c|c|c|} \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 \end{array} * \begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array} = \begin{array}{|c|c|c|c|c|c|} \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 \end{array} \begin{array}{|c|c|} \hline
 \text{light to dark} \\ \hline
 \text{matrix resultant} \\ \hline
 \end{array} \\
 \begin{array}{c} \blacksquare \\ \square \end{array} \\
 \begin{array}{|c|c|c|c|c|c|} \hline
 0 & 0 & 0 & 10 & 10 & 10 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 \\ \hline
 \end{array} * \begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array} = \begin{array}{|c|c|c|c|c|c|} \hline
 0 & -30 & -30 & 0 \\ \hline
 0 & -30 & -30 & 0 \\ \hline
 0 & -30 & -30 & 0 \\ \hline
 0 & -30 & -30 & 0 \\ \hline
 \end{array} \begin{array}{|c|c|} \hline
 \text{dark to light} \\ \hline
 \text{values.} \\ \hline
 \end{array} \\
 \begin{array}{c} \blacksquare \\ \blacksquare \end{array}
 \end{array}$$

- horizontal edge detector:

$$\begin{array}{|c|c|c|} \hline
 1 & 1 & 1 \\ \hline
 0 & 0 & 0 \\ \hline
 -1 & -1 & -1 \\ \hline
 \end{array}$$

- Sobel filter (vertical edge detector):

$$\begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 2 & 0 & -2 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array}$$

↳ this filter has a little more weight on the central pixel

making the processing a bit more robust

In the last training era, we've learned that hard coding nine numbers might not be as useful as just learning them. We can treat these nine numbers as parameters and use backpropagation such that we can convolve it with our image and obtain a good edge detector. (Can better learn the statistics of our dataset compared to something hard coded.)

w_1	w_2	w_3
w_4	w_5	w_6
w_7	w_8	w_9

Instead of learning some horizontal and vertical edges, we might be able to learn edges that are 45° , 70° or the orientation that's best suited for our problem.

padding (modification to the convolution operator)

- When we applied the 3×3 filter to our 6×6 image (through convolution) our output image shrunk to $\frac{n-f+1}{4} \times \frac{n-f+1}{4}$.

The downside to this technique is that we can only convolve our image a few times before our image starts getting very small (reducing down to 1x1) (we don't want our image to shrink every time we do an edge detection operation.)

The second downside is that the edge pixels (in the image) are only used for one direction, while the inner pixels overlap more (so they're used for more directions in the output)

0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0

↳ we are throwing away a lot of information near the edges of the picture.

- to solve both these problems, we can pad the image before applying the convolutional operation. (imagine adding a border of one pixel to the image)

↳ if the amount of padding at the border = p , then the output after the convolution operation will be

$$[n + 2p - f + 1 \times n + 2p - f + 1]$$

- how much should we pad our images? (valid and same convolutions)

↳ valid: no padding. $n \times n * f \times f \rightarrow n - f + 1 \times n - f + 1$

↳ same: pad so that the output size is the same as the input size.

$$\text{↳ } n + 2p - f + 1 \times n + 2p - f + 1$$

$$\text{↳ for } n + 2p - f + 1 = n \rightarrow p = \frac{f-1}{2}$$

↳ if which is the filter rows or columns \Rightarrow

↳ usually kept as an odd # in computer

vision to find a whole number p

(symmetric). When we have an odd 3x3

or 5x5, there's an absolute central position

to the filter (central pixel)

Strided Convolutions

- Strided convolutions happen when we jump to the right or up by more than one step \Rightarrow we convolve the

filter with the image. stride is denoted by s.

- When adding stride, the output is denoted as,

$$\lfloor \frac{n + 2p - f}{s} + 1 \times \frac{n + 2p - f}{s} + 1 \rfloor$$

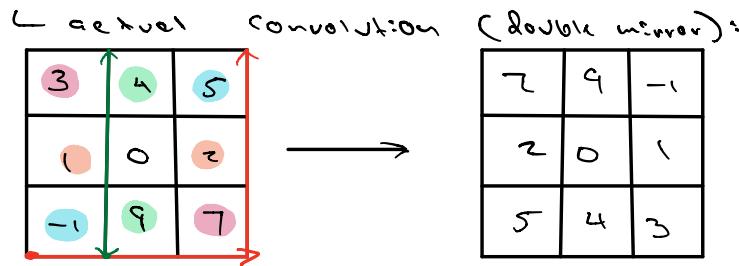
Where, n is the image's pixels, p is the # of padding pixels at an edge, f is the filter's pixels and s is the size of the stride.

What if the fraction is not an integer?

We floor the output such that if our input image (image+padding) has regions

where the filter lies outside, these specific computations aren't performed.

- Convolution vs Cross Correlation:



In actual convolution for signal processing, we double mirror the filter matrix before overlaying it onto the image matrix.

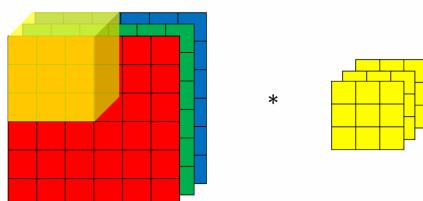
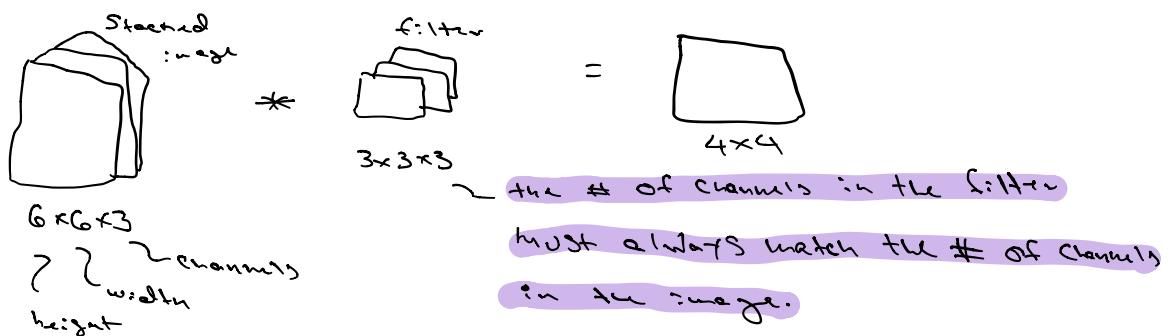
This enables the conservation of associativity such that, $(A * B) * C = A * (B * C)$

This property is not important in deep learning, therefore, we don't do this operation and just directly overlay

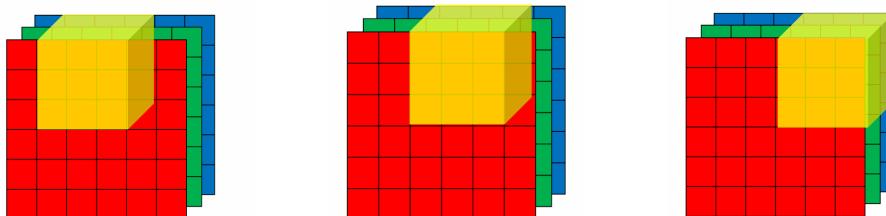
the filter on the image (aving computational cost). Hence, Cross-correlation is just called convolution in deep learning applications as a convention.

Convolution over Volumes (3D)

- detecting convolution on RGB images ($6 \times 6 \times 3$)



- we overlay the $3 \times 3 \times 3$ filter on the volume (similar to the 2D grayscale case) and shift it to the right and down sequentially.



- the $3 \times 3 \times 3$ filter has a total of 27 numbers as we overlay them on top of the volume, we do an element wise multiplication and then sum. we obtain a resultant

4×4 matrix (similar to the 2D case).

- what's so special about analyzing volumes?

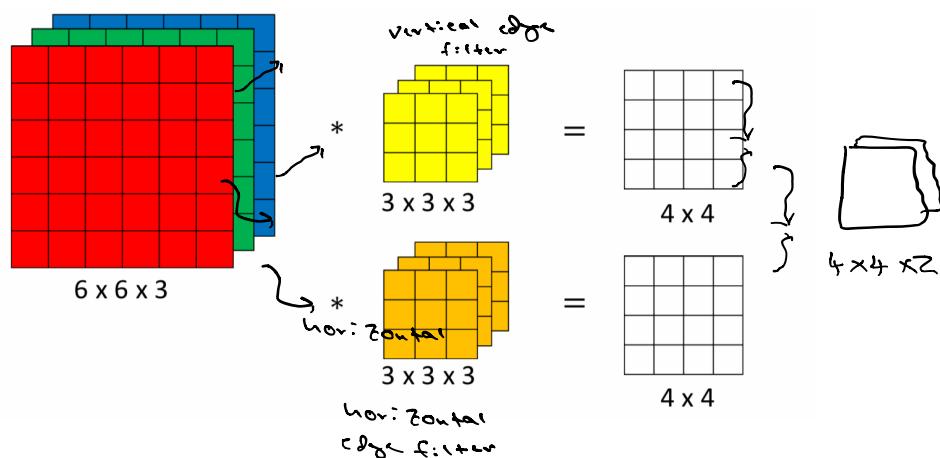
↳ we can choose to find edges in particular channel, for example to detect vertical edges in the red channel ...

$$\begin{array}{c} R \\ \boxed{\begin{matrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{matrix}} \\ G \\ \boxed{\begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix}} \\ B \\ \boxed{\begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix}} \end{array}$$

We stitch this together to form the $3 \times 3 \times 3$ filter.

↳ alternatively, we can set all the filters to the same value to find vertical edges in any channel.

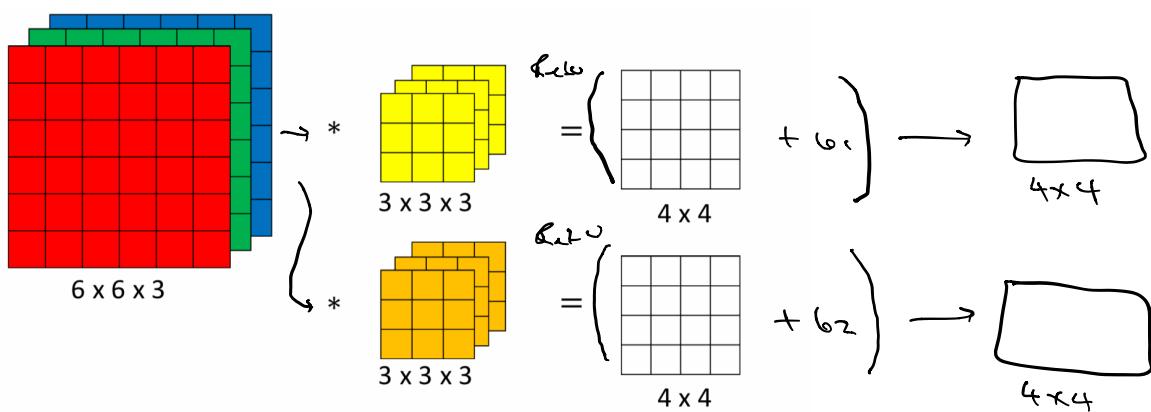
- how do we combine multiple filters to detect both horizontal and vertical (or other types) of edges at the same time?



↳ the # of filters in the output is equal to the # of features being extracted from the volume (in this case

we applied two filters and summed them together, hence, our depth = 2).

One layer implementation of a convolutional neural network



- to make the output of convolutional network layer, we add a bias (through python broadcasting) and apply a non-linear activation function for each of the outputs.

After this step we stack both outputs, to make a $4 \times 4 \times 2$ volume.

The filters acts similar to the weight matrix in a normal neural network scenario, $\tilde{z}^{[l]} = W^{[l]} \alpha^{[l]} + b^{[l]}$

↳ we transform the $6 \times 6 \times 3$ into a $4 \times 4 \times 2$ after passing one layer of the convolutional neural network.

↳ because the # of parameters is only dependent on the size of the filter, the # of filters and the bias term, the # of parameters that must be learned are fixed for an image of any input size.

↳ so filters that are $3 \times 3 \times 3$ (+1 from the bias)

have a total of ~80 parameters that must be learned regardless of the size of the input image (resolution).

This feature of the convnet makes it less prone to overfitting (the small parameter size)

- Summary of notation (convolutional neural network layer)

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

- $p^{[l]}$ = padding (valid vs same convolution)

- $S^{[l]}$ = Stride

Input = $h \times w \times n_c^{[l-1]} \sim \# \text{ of channel}$
 Activations from previous layer

Output: $h' \times w' \times n_c^{[l]} \quad \begin{matrix} \text{if } l=1, \text{ it's the size of} \\ \text{the input image} \end{matrix}$

The output width is given by,

$$n' = \frac{n^{[l-1]} + 2p^{[l]} - f^{[l]}}{S^{[l]}} + 1 \quad (\text{and we floor this #})$$

* This formula can be used for either the height or the width.

* $n_c^{[l]}$ in the output is only dependent on the # of filters that we process in the layer of the neural network.

Each filter is: $f^{[l]} \times f^{[l]} \times n_c^{[l-1]}$

The activation dimensions are, $a^{[l-1]} \rightarrow h^{[l]} \times w^{[l]} \times n_c^{[l]}$

If we batch the examples (for batch gradient descent),

$$A^{[l]} = m \times h^{[l]} \times w^{[l]} \times n_c^{[l]}$$

Where m is the index of the training example

The weights are given as, $f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n^{[l]}$

We take into account both the # of filters being

inputted from the previous layer and the # of filters being processed in the current layer

L 6:2:

↳ we will have one real # for each filter being processed within the layer, hence, the dimensions are $^{(c)} h \times w$

A simple Convolution network Example

After the last step of a convolutional neural network is to flatten the output of the last layer (lets say $7 \times 7 \times 40$) into units (1960 units) and then densely connecting those units to a single final output neuron which uses a logistic regression or softmax to compute the output y .

• a lot of the work in designing convolutional neural networks is selecting the hyperparameters such as the total size, stride, padding in each layer, how many filters to use, and so on.

↳ a suggestion is that as we go deeper in a convnet, we want the height and width to gradually trend down as we go deeper into the network and the # of channels to increase

$$39 \times 39 \times 3 \rightarrow 37 \times 37 \times 10 \rightarrow 17 \times 17 \times 20 \rightarrow 7 \times 7 \times 40$$

• three types of layers in a convolutional neural network

↳ convolution (conv): what we have been using up to this point.

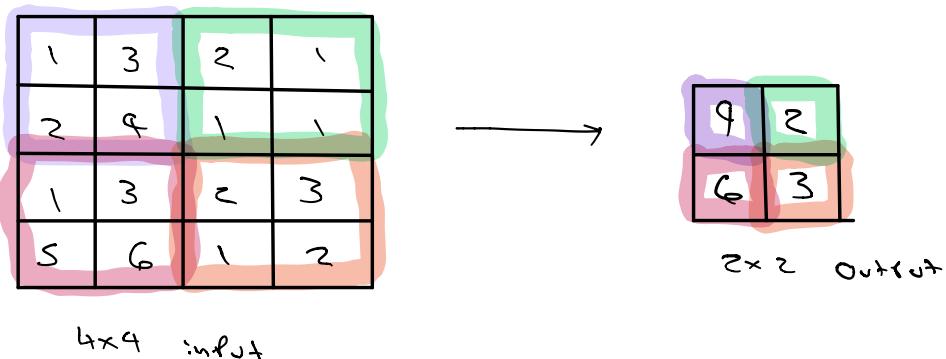
↳ pooling (pool)

↳ fully connected (FC): like the standard neural network layer

* it's possible to design a good neural network with only Conv layers. However, most neural network architectures also have pooling layer and fully connected layers.

Pooling layers

- Pooling layers are used to reduce the size of the representation + to speed up computation as well as make some features detected a bit more robust
- pooling layer: max pooling



Let's take the 4x4 input, break it into different regions and in the output, each of the numbers will be the max from the corresponding shaded region.

Let's take we are applying a filter of $f=2$ with a stride of 2. $(4-2)/2 + 1 = 2$

What's the intuition behind what max pooling is doing?

If we think about the 4x4 as some set of features (activations in some layer of the network) then

a large # might mean that we've detected a particular feature.

An interesting property of max pooling is that it has a set of hyperparameters, but it has no parameters to learn.

Once s and f are set, it's a fixed computation.

If we have a 3D input, the output in max pooling will have the same # of channels/depth as the input. The process of taking the max remains (computed on each channel independently).

- another type of pooling which is used less often is average pooling.

In average pooling we take the average of each region instead of the max.

An exception to when average pooling is used over max pooling is very deep in the network when we want to collapse our representation

$$7 \times 7 \times 1000 \rightarrow 1 \times 1 \times 1000$$

* By using hyperparameters $s=2$ $f=2$, the height and width are roughly shrunk down by a factor of 2. When we do pooling, we almost never use padding.

Why convolution?

- two main advantages of using conv layers instead of only

using the fully connected layers variant.

i. Parameter sharing: a feature detector (such as a vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.

If we define a 3×3 filter to detect edges, we can apply the same 3×3 filter to all the input. This applies to both low level features (edges) and high level features (such as faces).

ii. Sparsity of connections: In each layer, each output value depends on a small # of inputs (the size of the overlaid filter / shaded region of the input where we are performing the convolution).