

# *CS-577 Deep Learning. Assignment 2*

*Guillermo Lopez-Areal*

*Illinois Institute of Technology*



ILLINOIS INSTITUTE  
OF TECHNOLOGY

1. Let  $I$  be a  $4 \times 4$  RGB image where the R channel is all 1-s and G channel is all 2-s. The B channel has a value of 1 in its first row, a value of 2 in its second row, a value of 3 in its third row, and a value of 4 in its 4th row. Compute the convolution of this image with a  $3 \times 3$  filter having all ones without zero padding.

$I = 4 \times 4$  RGB  $3 \times 3$  filter

$R \rightarrow 1$

$G \rightarrow 2$

$B \rightarrow \begin{cases} 1 \\ 2 \\ 3 \\ 4 \end{cases}$

121	121	121	121
122	122	122	122
123	123	123	123
124	124	124	124

A	B
C	D

Since  $A = B$  &  $C = D$ , we'll compute only A and C. They will be the same because the convolutions for both, take the same exact values.

$$\begin{aligned} A: & 121 + 121 + 121 \\ & 122 + 122 + 122 \\ & 123 + 123 + 123 \end{aligned}$$

$$\{9, 18, 18\} = 45$$

$$\begin{bmatrix} 45 & 45 \\ 54 & 54 \end{bmatrix}$$

$$\begin{aligned} C: & 122 + 122 + 122 \\ & 123 + 123 + 123 \\ & 124 + 124 + 124 \end{aligned}$$

$$\{9, 18, 27\} = 54$$

2. Repeat the previous question with zero padding.

zero padding

0	0	0	0	0	0
0					0
0		M			0
0					0
0					0
0	0	0	0	0	0

A	B	C	D
E	F	G	H
I	J	K	L
M	N	O	P

$$A: (121) \cdot 2 + (122) \cdot 2 = 18$$

$$B: (121) \cdot 3 + (122) \cdot 3 = 27$$

$$E: 2(121 + 122 + 123) = 30$$

$$I: 2(122 + 123 + 124) = 36$$

$$M: 2(123 + 124) = 26$$

$$F: 3(121 + 122 + 123) = 45$$

$$J: 3(122 + 123 + 124) = 54$$

$$\{A = D\}, \{B = C\}, \{E = H\}, \{I = L\}, \{M = P\}$$

$$\{F = G\}, \{J = K\}$$

$\Rightarrow$	18	27	27	18
	30	45	45	30
	36	54	54	36
	26	39	39	26

3. Repeat the previous question when using dilated (atrous) convolution with a dilation rate of 2.

Dilated convolution w dilation Rate 2

It is pretty much the same in terms of the matrix, so we'll compute the answer without tracing the matrices

$$\Delta: 2(121) + 2(123) = 20 \quad \left( \begin{array}{cc} 21 & 24 \\ 20 & 20 \end{array} \right)$$

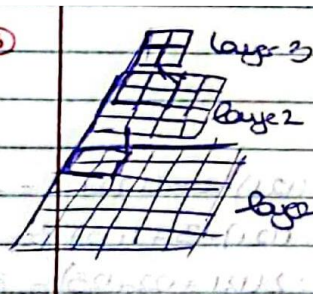
$$E: (122) \cdot 2 + 2(124) = 24$$

4. Explain the template matching interpretation of convolution.

Template matching of convolution

Convolutional filters could be equivalent to an artificial neuron activation in every region scanned by the filter. When there exists a resemblance in the image and the filter, there's expected a high response.

5. Explain how multiple scale analysis can be achieved with a fixed window size (using a pyramid).



Spatial dimensions are sampled and we get a shape as follows; a pyramid shape where each layer corresponds to a network with different resolutions.

6. Explain how to compensate for spatial resolution decrease using depth (number of channels) and the purpose for doing so.

As spatial dimensions decrease, depth increase to compensate for reduced coefficients. This is done to help, to facilitate the learning of abstract and complex structures and make them smaller to fit them in the network.

7. Given a  $128 \times 128 \times 32$  tensor and 16 convolution filters of size  $3 \times 3 \times 32$ , what will be the size of the resulting tensor when convolving without zero padding.

$$\text{Size: } w - k + 2p = 1 \Rightarrow 126 \times 126 \times 16$$

8. Repeat the previous question when using a stride of 2.

$$63 \times 63 \times 16$$



9. Explain how the number of channels can be reduced using a  $1 \times 1$  convolution.

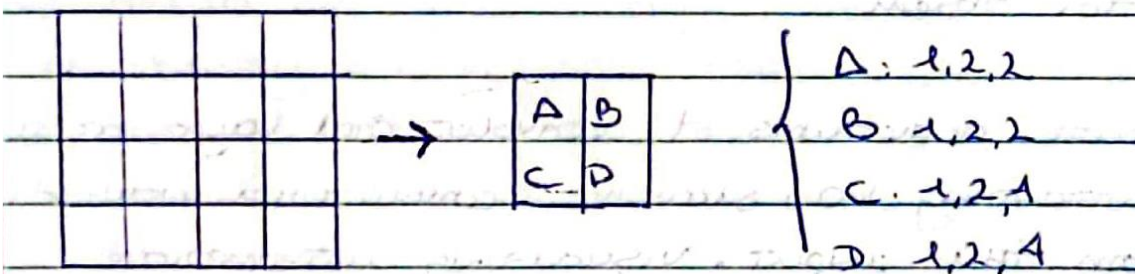
convolution with  $N$   $1 \times 1$  filter results  $N$  number of channels output. with a fewer  $n$  of filters it can also be reduced.

10. Explain the interpretation of convolution layers and the difference between early and deeper convolution layers.

used to extract image patterns, vectorize image windows and filter extending full length of image

11. Let  $I$  be an image as in question 1. Write the result obtained using max pooling with a  $2 \times 2$  filter with a stride of 2.

Pooling with  $2 \times 2$  filter with stride 2



12. Explain the purpose of pooling.

Pooling is made for down sampling dimensions (depth unchanged).

13. Explain the purpose of data augmentation and when it is most useful.

Data augmentation: augment data for better generalization using modifications of available data. It is commonly used when there are small datasets.

14. Explain the purpose of transfer learning and when it is most useful.

Transfer learning is the use of pre-trained model's weights. It is commonly used to get a better performance when training new models, especially in convolutional neural networks for image recognition.

15. Explain the need for freezing the coefficients of the pre-trained network.

Freezing is done to not push the gradients through the pre-trained base when starting to train the new classifier on the top of the network.

16. Explain how the coefficients of a pre-trained network can be fine-tuned.

Fine-tuning is done after training the classifier, by unfreezing some layers on top of the conv-base and retraining to let the model to fit the data.

17. Explain the purpose of inception blocks.

Inception blocks allow to use multiple filter sizes in a single image block in convolutional layers.

18. Explain the advantage of residual blocks.

Residual blocks help to avoid gradient vanishing in deep convolutional neural networks.

Zero weights in the block produce identity instead of destroying the signal.

Residual blocks create an identity mapping to activations earlier in the network, in order to thwart the performance degradation problem.

19. Explain how intermediate activations of convolution layers can be visualized given an input. What is the purpose for doing so?

Intermediate activations of convolutional layers are useful for understanding how successive convolutional network layers transform their input. Visualizing intermediate activations consists of displaying the feature maps that are outputs by various convolution layers in a network. This gives a view into how an input is decomposed into the different filters learned by the network.

20. Explain how the filter weight of the trained convolution layers can be visualized (using gradient ascent to find the input with maximal response). What is the purpose for doing so?

The way in which the filter weights are determined to achieve the most efficient outcome, is by performing gradient descent, so as to find the input features from which the outcome of the filters is maximized. The purpose of doing so, is so we can determine which filter to use, and how such filter looks over each epoch.



21. Explain how the heatmap of class activation can be visualized for a specific image and class. explain how pooled gradients can be used to weight channels in this visualization. Explain the purpose of this visualization.

After sending a certain image to our neural network, ~~and comparing it~~ ~~to~~ we compute gradient descent ~~over~~ in each channel. And what it's interesting, is that the neural network decides which CNN channel is more crucial for predicting the results. And this ~~loss~~ is done by comparing the gradients in each channel and comparing them to others, so as the one that is more significant, will receive a higher weight. Similarly those who are less significant will receive a lower weight. And finally, we perform global average pooling, so as we add up on ~~combine~~ every possible channel and, we have finally an output. This is done, so as we can visualize which ~~are~~ of the image helped the most. The neural network to determine to which class a certain image belongs.