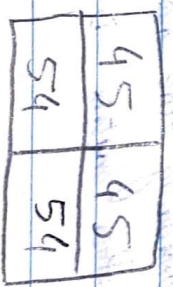
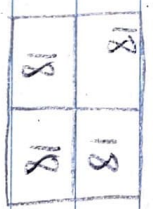
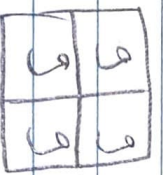
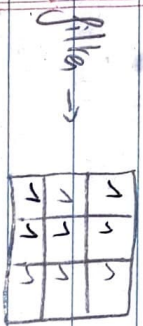
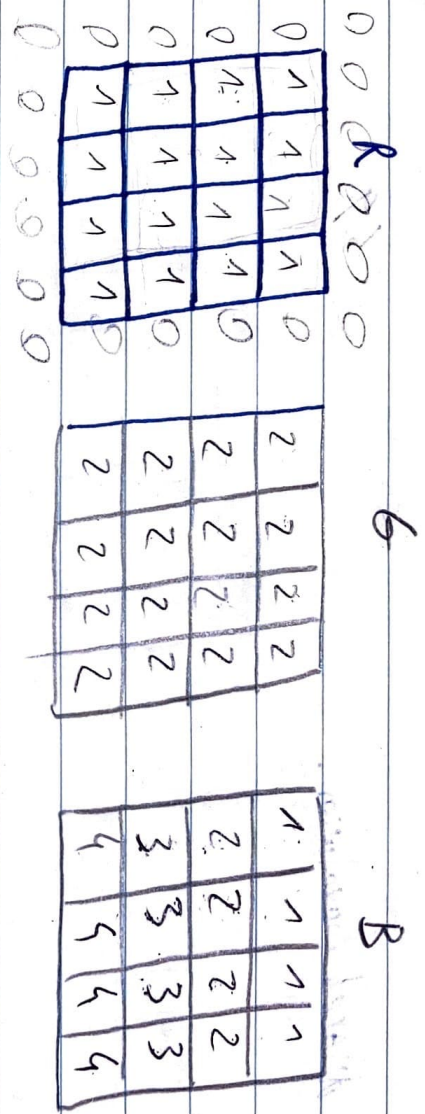


GABRID  
BARANS

Deep Learning HW4 CS577

Q.1:



convolution of the image  
without zero padding.

2) with zero padding:



convolution of the image  
with zero padding.



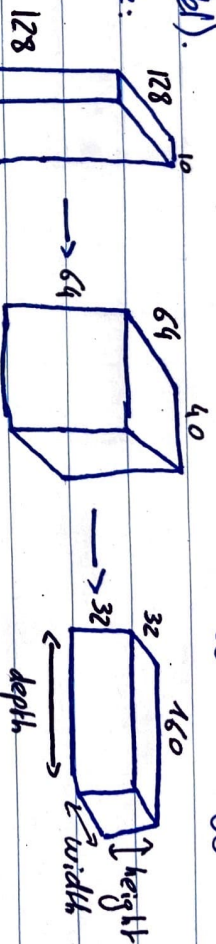


Q6:

We use multiple convolutional layers and make the image's resolution smaller and smaller to be able to detect bigger and bigger objects.

(Pyramid model).

Take Example:



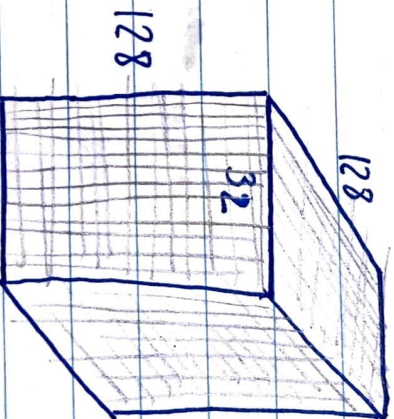
As we go deeper and deeper into the network, the resolution, height and width decreases. At the same time that we reduce the spatial resolution, we increase the depth.

When we reduce the spatial resolution, we lose information.

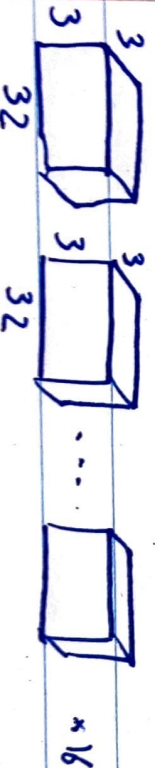
To compensate the low resolution spatial dimension, we get higher depth.

As spatial dimension decrease, depth increase to compensate for reduction coefficient (Keep the same number of coefficient).

2) Reason  $128 \times 128 \times 32$  and 16 convolution layers  $3 \times 3 \times 32$

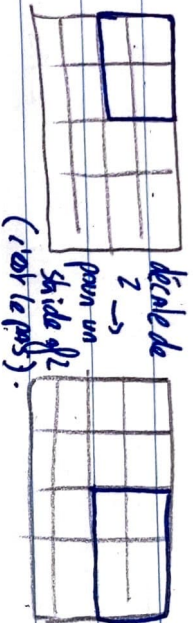


Size of the Resulting tensor without zero padding:  $128 \times 128 \times 16$





8) with a stride of two, without zero padding, the size of the resulting tensor is  $4 \times 3 \times 3 \times 16$ .



9)

A  $1 \times 1$  convolution simply maps an input pixel with all it's channels to an output pixel not looking at anything around itself. It is often used to reduce the number of depth channels, since it is often very slow to multiply volumes with extremely large depth.

Input 256 depth  $\rightarrow$   $1 \times 1$  convolution: 64 depth  $\rightarrow$   $1 \times 1$  convolution: 256 depth.

10) The early layers (in the network) extract simple features, and as we progress, we extract more and more complex features (i.e. layers). A convolution layer network then have early layers and deep layers.

11) max pooling with image from Q1)

$$\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} + \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix} + \begin{bmatrix} 2 & 2 \\ 4 & 4 \end{bmatrix} = \begin{bmatrix} 5 & 5 \\ 7 & 7 \end{bmatrix}$$



12) Pooling is used to reduce the <sup>spatial</sup> size of the image. Pooling reduces the height and width but not the depth of the image.

13) Data augmentation is used to augment data for better generalization using the image data generation. It is useful when we work with images to reduce the overfitting.

14) Transfer learning: we use a pre-trained convnet trained on a large set (1 million images for example). Used for classification but also for object segmentation and other tasks.

Transfer learning is a way to speed up deep learning training. It helps solve complex problems with pre-existing knowledge.

15) The task of fine-tuning a network is to tweak the parameters of an already existing trained network so that it adapts to the new task at hand. The initial layer learns very general features and as we go higher up the network, the layers tend to learn ~~more~~ patterns more specific to the task it is being trained on. Thus, for fine-tuning, we want to keep the initial layers intact and freeze them, ~~and retain the layers to be changed~~.

16) After training the fully connected layer, we unfreeze some top layers in "conv-base" and retain to allow the model to fit better the data.

Steps:

- all custom network on top of trained layers.
- jointly train the custom network and unfreeze layers.
- Freeze trained layers
- train custom network
- unfreeze top layers in the base network



17)

An inception block aims to approximate an optimal local sparse structure in a CNN. It allows for us to use multiple type of filter size instead of being restricted to a single filter size.

18)

Essentially, residual blocks allow memory or information to flow from the initial to last layers.

19) To visualize intermediate activation function, we simply plot what each filter has extracted given an input.

Visualizing intermediate activation gives a view into how an input is decomposed into the different filters learned by the network.

20) A way of learning about what the convolution network is looking for in the images is to visualize the convolution layer filter.

By displaying the network layer filter, we can learn about the pattern to which each filter will correspond to.

21) To visualize the heatmap, we use Grad-CAM (Gradient class Activation Map). The idea behind it is to find the importance of a certain class in the model, uptake its gradient with respect to the final convolutional layer and then weight it against the output of this layer.

It is useful because it helps understand if the neural network is looking at appropriate parts of the image, or if the neural network is cheating.