

CS512 - COMPUTER VISION
ASSIGNMENT - 4

SAI MANOHAR YEMURI
A20514848

① Convolution Layers

a) 4x4 RGB Image

Channel R

1	1	1	1
1	1	1	1
1	1	1	1
1	1	1	1

*

1	1	1
1	1	1
1	1	1

3x3 filter

=>

9	9
9	9

ans

(1)

$$10 \times 10 = 100$$

1 = 4 examples

$$\text{filter} = (1+1+1+1) * 3 \\ \text{Conn} = 9$$

Similarly, we need to move the filter by 1 position as the stride is 1.
Since, all the values are same we get the same output i.e "9".

Channel G

2	2	2	2
2	2	2	2
2	2	2	2
2	2	2	2

*

1	1	1
1	1	1
1	1	1

=>

18	18
18	18

for ans,
ans

$$= (2*1 + 2*1 + 2*1) * 3 \\ = 18$$

Similarly, we need to apply the filter for the rest of the image pixels with the stride 1.

channel B

$$\begin{array}{|c|c|c|c|} \hline 1 & 1 & 1 & 1 \\ \hline 2 & 2 & 2 & 2 \\ \hline 3 & 3 & 3 & 3 \\ \hline 4 & 4 & 4 & 4 \\ \hline \end{array} * \begin{array}{|c|c|c|c|} \hline 1 & 1 & 1 & 1 \\ \hline 1 & 1 & 1 & 1 \\ \hline 1 & 1 & 1 & 1 \\ \hline \end{array}$$

first Element = $\begin{pmatrix} 1x1 + 1x1 + 1x1 \\ 2x1 + 2x1 + 2x1 \\ 3x1 + 3x1 + 3x1 \end{pmatrix} = 3+6+9 = 18$

Second Element = 18

Third Element = $\begin{array}{|c|c|c|} \hline 2 & 3 & 1 \\ \hline 3 & 3 & 3 \\ \hline 4 & 4 & 4 \\ \hline \end{array} * \begin{array}{|c|c|c|c|} \hline 1 & 1 & 1 & 1 \\ \hline 1 & 1 & 1 & 1 \\ \hline 1 & 1 & 1 & 1 \\ \hline \end{array}$

$$= \begin{pmatrix} 2x1 + 2x1 + 2x1 \\ + 3x1 + 3x1 + 3x1 \\ + 4x1 + 4x1 + 4x1 \end{pmatrix} = 27$$

Fourth Element = 27

So, the matrix will be

$$\begin{array}{|c|c|} \hline 18 & 18 \\ \hline 27 & 27 \\ \hline \end{array}$$

So, we have

$$R = \begin{array}{|c|c|} \hline 9 & 9 \\ \hline 9 & 9 \\ \hline \end{array}, \quad G = \begin{array}{|c|c|} \hline 18 & 18 \\ \hline 18 & 18 \\ \hline \end{array}, \quad B = \begin{array}{|c|c|} \hline 18 & 18 \\ \hline 27 & 27 \\ \hline \end{array}$$

⇒ Element-wise addition of all 3-channels

$$= \begin{bmatrix} 9+18+18 & 9+18+18 \\ 9+18+27 & 9+18+27 \end{bmatrix} = \begin{bmatrix} 45 & 45 \\ 54 & 54 \end{bmatrix}$$

⇒ On. Averaging each Element with $\frac{1}{9}$, The New Elements will be

$$= \begin{bmatrix} 5 & 5 \\ 6 & 6 \end{bmatrix}$$

b)

$$R = \begin{array}{|c|c|c|c|c|c|} \hline 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 1 & 1 & 1 & 1 & 0 \\ \hline 0 & 2 & 2 & 1 & 1 & 0 \\ \hline 0 & 1 & 1 & 1 & 1 & 0 \\ \hline 0 & 1 & 1 & 1 & 1 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 & 0 \\ \hline \end{array} * \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$

∴ final Convolution matrix with zero padding

$$\begin{array}{|c|c|c|c|} \hline 4 & 6 & 6 & 4 \\ \hline 6 & 9 & 9 & 6 \\ \hline 6 & 9 & 9 & 6 \\ \hline 4 & 6 & 6 & 4 \\ \hline \end{array}$$

channel G

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

$$* \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array} \Rightarrow$$

final Convolution matrix
with zero Padding

$$\begin{array}{|c|c|c|c|} \hline 8 & 12 & 12 & 8 \\ \hline 12 & 18 & 18 & 12 \\ \hline 12 & 18 & 18 & 12 \\ \hline 8 & 12 & 12 & 8 \\ \hline \end{array}$$

channel B

final convolution matrix
with zero padding

0	0	0	0	0	0
0	1	1	1	1	0
0	2	2	2	2	0
0	3	3	3	3	0
0	4	4	4	4	0
0	0	0	0	0	0

$$* \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array} \Rightarrow \begin{array}{|c|c|c|c|} \hline 6 & 9 & 9 & 6 \\ \hline 12 & 18 & 18 & 12 \\ \hline 18 & 27 & 27 & 18 \\ \hline 14 & 21 & 21 & 14 \\ \hline \end{array}$$

Element wise addition of all channels

$$= \begin{array}{|c|c|c|c|} \hline 4 & 6 & 6 & 4 \\ \hline 6 & 9 & 9 & 6 \\ \hline 6 & 9 & 9 & 6 \\ \hline 4 & 6 & 6 & 4 \\ \hline \end{array} + \begin{array}{|c|c|c|c|} \hline 8 & 12 & 12 & 8 \\ \hline 12 & 18 & 18 & 12 \\ \hline 12 & 18 & 18 & 12 \\ \hline 8 & 12 & 12 & 8 \\ \hline \end{array} + \begin{array}{|c|c|c|c|} \hline 6 & 9 & 9 & 6 \\ \hline 12 & 18 & 18 & 12 \\ \hline 18 & 27 & 27 & 18 \\ \hline 14 & 21 & 21 & 14 \\ \hline \end{array}$$

$$= \begin{array}{|c|c|c|c|} \hline 18 & 27 & 27 & 18 \\ \hline 30 & 45 & 45 & 30 \\ \hline 30 & 54 & 54 & 30 \\ \hline 26 & 39 & 39 & 26 \\ \hline \end{array} \Rightarrow \text{Resulting matrix after the convolution across } 3 \text{ channels and element wise addition.}$$

In case of Averaging the result, the resulting matrix would be

$$\Rightarrow \frac{1}{9} \begin{array}{|c|c|c|c|} \hline 18 & 27 & 27 & 18 \\ \hline 30 & 45 & 45 & 30 \\ \hline 30 & 54 & 54 & 30 \\ \hline 26 & 39 & 39 & 26 \\ \hline \end{array}$$

stepwise averaging of layers

$$= \begin{array}{|c|c|c|c|} \hline 2 & 3 & 3 & 2 \\ \hline 3.33 & 5 & 5 & 3.33 \\ \hline 3.33 & 6 & 6 & 3.33 \\ \hline 2.88 & 4.3 & 4.3 & 2.08 \\ \hline \end{array}$$

c) Dilated Kernel (with rate = 2)

$$\begin{array}{|c|c|c|c|} \hline & 81 & 81 & 81 \\ \hline 81 & f_1 & f_2 & f_3 \\ \hline f_1 & 81 & 81 & 81 \\ \hline f_2 & 81 & 81 & 81 \\ \hline f_3 & 81 & 81 & 81 \\ \hline \end{array} \quad * \quad \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$

$$K = \begin{bmatrix} 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 \end{bmatrix}$$

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

elements 0 for variables with kernel

R channel :- (Initial zero Padding)

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

$$\begin{array}{|c|c|c|c|c|} \hline 8 & 81 & 81 & 81 & 8 \\ \hline 81 & 81 & 81 & 81 & 81 \\ \hline 81 & 81 & 81 & 81 & 81 \\ \hline 8 & 81 & 81 & 81 & 8 \\ \hline \end{array}$$

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

5×5

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

$$\begin{array}{|c|c|c|c|c|} \hline 81 & f_1 & f_2 & f_3 & 81 \\ \hline f_1 & 81 & f_2 & f_3 & 81 \\ \hline f_2 & f_3 & 81 & 81 & 0 \\ \hline f_3 & 81 & 0 & 81 & 0 \\ \hline 81 & 0 & 0 & 81 & 0 \\ \hline \end{array}$$

1st element = ~~0+0+0+0+0~~ $0+0+2+0+2$ (padding is zero)

$$= 4$$

Similarly we need to compute for the rest of the matrix
Resulting matrix would be

$$\begin{array}{|c|c|} \hline 4 & 4 \\ \hline 4 & 4 \\ \hline \end{array}$$

G channel :-

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 2 & 2 & 2 & 0 \\ 0 & 2 & 2 & 2 & 2 & 0 \\ 0 & 2 & 2 & 2 & 2 & 0 \\ 0 & 2 & 2 & 2 & 2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} * \begin{bmatrix} 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 \end{bmatrix}$$

1st Element = $4 + 4 = 8$

Similarly we need to compute for the rest of the matrix
and the feature Map will be =

$$\begin{bmatrix} 8 & 8 \\ 8 & 8 \end{bmatrix}$$

B channel :-

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 1 & 0 \\ 0 & 2 & 2 & 2 & 2 & 0 \\ 0 & 3 & 3 & 3 & 3 & 0 \\ 0 & 4 & 4 & 4 & 4 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} * \begin{bmatrix} 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 \end{bmatrix}$$

Feature Map =

$$\begin{bmatrix} 12 & 12 \\ 8 & 8 \end{bmatrix}$$

Element wise addition of all channels

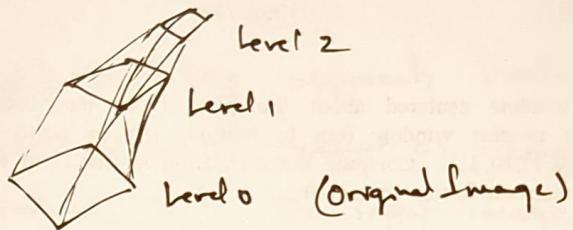
$$\begin{bmatrix} 4 & 4 \\ 4 & 4 \end{bmatrix} + \begin{bmatrix} 8 & 8 \\ 8 & 8 \end{bmatrix} + \begin{bmatrix} 12 & 12 \\ 8 & 8 \end{bmatrix} = \begin{bmatrix} 24 & 24 \\ 20 & 20 \end{bmatrix}$$

d) Template matching :-

Template matching approach involves defining a small image or kernel, which is known as template. The template is then convolved with the input image by sliding it across the image, one pixel at a time and computing the similarity between the template ~~and~~ and covered window on the image using the dot product.

It is implemented through 2D convolution. In convolution, the value of the off pixel is computed by multiplying elements of two matrices & summing the results.

e) MultiScale Analysis :-



To achieve the multiscale analysis with a fixed window size, a pyramid can be constructed by down-sampling the original image. By using the fixed size window at different scale, we can analyze the image at different scales and detect features more accurately.

f) In order to compensate the spatial resolution decrease, we can increase the depth or number of channels in a neural network to extract more complex features and compensate for the loss of information.

g) Input tensor = $(128 \times 128 \times 32) \Rightarrow (h \times w \times \text{channels})$

Convolution filter size = ~~(3×3)~~ $\Rightarrow (h \times w)$

Output Width =
$$\frac{\text{Input Tensor width} - \text{filterWidth} + 2(\text{Padding})}{\text{Stride}} + 1$$

$$= \frac{128 - 3 + 2(0)}{1} + 1$$

$$= 125 + 1$$

$$= 126$$

Similarly, Output height is also 126. Since there are 16 convolution filters, the channels in the off tensor would be 16.

Output shape = $(126 \times 126 \times 16)$

h) $= \frac{128 - 3 + 2(0)}{2} + 1$

$= \frac{125}{2} + 1$

$= 62.5 + 1$

$= 63.5$

≈ 63

o/p tensor shape = $(63 \times 63 \times 16)$

- i) The 1×1 Convolution layer can reduce the no. of channels in a feature map by using a 1×1 convolution with fewer filters.

ex:-

Input shape - $64 \times 64 \times 64$

filter - $1 \times 1 \times 64$, filters = 3

output shape - $64 \times 64 \times 3$

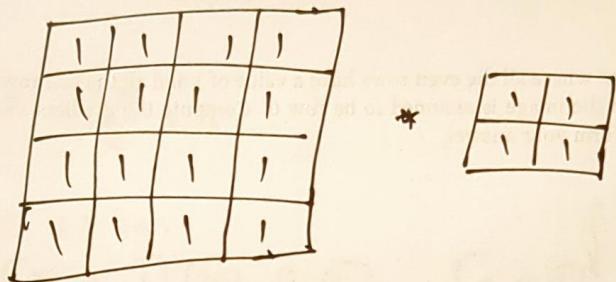
- j) A convolution layers generally used with ConvNet is a class of neural network that specializes in processing data that has a grid like topology such as an image

Convolution layers is the core building block of CNN. It performs the dot product b/w kernel and receptive field of image.

Difference b/w Early & deep layers:-

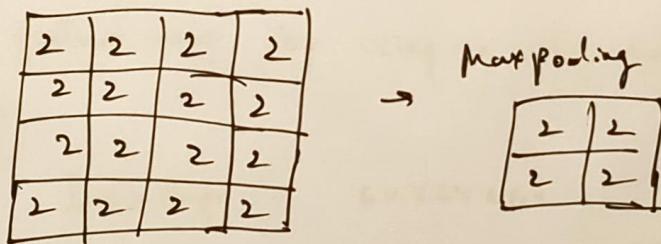
- * Early Convolution layers can learn features like line, edge etc.
- * Deep Convolution layers will learn more complex features such as face detection, object detection etc.

K) R channel :-

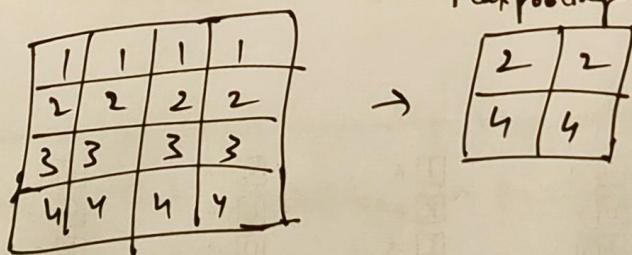


$$\text{Max pooling with stride 2} = \begin{array}{|c|c|} \hline 1 & 1 \\ \hline 1 & 1 \\ \hline \end{array}$$

G channel :-



B channel :-



l) Pooling! - down sampling spatial dimensions without changing the depth in the image.

Purpose of Pooling :-

1. It is used for Multi-scale Analysis
2. Makes Convolutions Computationally less expensive
3. helps in reducing the no. of network coefficients

m) Data Augmentation :- It is a technique which is used for increasing the size of a dataset by applying various transformations to the original data.

Purpose of data augmentation :-

- * It helps the Machine learning model to improve the performance and generalization ability.
- * It helps by adding more diversity to the training data which helps in reducing the overfitting.
- * It is useful when training dataset is small.
- * Helps to address the class imbalance problem in the training data.

d. CNN

a) Transfer learning :- It is a technique in Machine learning where a pretrained model is used instead of training the model from scratch.

The purpose of transfer learning is to improve the performance and reduce the training time of a new model by using the knowledge and features learned from a pretrained model on some related task.

Transfer learning is useful when

- * You have limited data available for training the model.
- * The new task is related to that the pretrained model was trained on. In such case, the pretrained models' learned features can be adopted to our task.
- * New task requires high computational resources and we have limited time to train.

- b) freezing the coefficients of the pretrained network refers to setting the parameters of a layer as non-trainable to prevent them from being updated during the training process.

The need for freezing is that pretrained models have learned features specific to the original task they were trained on, and freezing them will prevent the weights being disrupted during training on new task. This allows the new layers added to the network to be trained on the new task without affecting the pretrained layers learned features.

- i) fine tuning a pretrained network involves adjusting its hyperparameters on a new dataset to improve its performance on a specific task. Steps included to fine tune a pretrained network -
- i, Load a pretrained model using tensorflow or pytorch that you want to fine tune.
 - ii, Replace last-layer of the network with a new one that has the appropriate number of outputs for your new task.
 - iii, freeze some layers during the training in order to prevent them from updating. These layers ~~not~~ belong to the pretrained network.
 - iv, after configuring the network, train it on new dataset and Evaluate the performance of the fine tuned network on Validation set.

d) Inception Blocks :- Inception Blocks are a type of Building blocks used in CNN which are used to improve the ability to capture complex patterns in the input data. They are a combination of convolutional layers with different filter sizes and pooling operations to extract features at multiple scales and combine them to form a rich representation of the input.

To address the vanishing gradients problem in Googlenet, the network employed several solutions, including the use of the inception module, auxiliary classifiers placed at intermediate layers, and the ReLU activation function.

The inception module allows the network to learn diverse features at different scales, while auxiliary classifiers force the network to learn more meaningful representations. And ReLU helps prevent output from exploding and vanishing gradient problem.

e) Residual Blocks have skip connections that ~~pass~~ preserve the gradient signal during backpropagation that helps to prevent vanishing gradients problem. This allows the network to learn residual functions and focus on learning the fine grained details of the input. This is because the skip connection ensures that the gradient signal is preserved and not lost as it passes through the network and prevents it from vanishing gradient problems.

f) DenseNet :-

Each layer in the network is connected to the previous layer, which allows to reuse the features and reduce the parameters. DenseNet uses the Bottleneck like Unit to reduce the no. of featuremaps before each dense block which helps to reduce the no. of parameters. Global Average pooling is used to reduce the no. of parameters and improve regularization. It does not use fully connected layer at the end of the net. There is a transition layer in between each block which helps us to reduce the no. of featuremaps and spatial dimensions.

g) Let $R = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}$

 $G = \begin{bmatrix} 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \end{bmatrix}$
 $B = \begin{bmatrix} 3 & 3 & 3 & 3 \\ 3 & 3 & 3 & 3 \\ 3 & 3 & 3 & 3 \\ 3 & 3 & 3 & 3 \end{bmatrix}$

filter 1 = $\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ filter 2 = $\begin{bmatrix} 2 & 2 & 2 \\ 2 & 2 & 2 \\ 2 & 2 & 2 \end{bmatrix}$ filter 3 = $\begin{bmatrix} 3 & 3 & 3 \\ 3 & 3 & 3 \\ 3 & 3 & 3 \end{bmatrix}$

$R * \text{filter 1} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} + \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 9 & 9 \\ 9 & 9 \end{bmatrix}$

$G * \text{filter 2} = \begin{bmatrix} 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \end{bmatrix} * \begin{bmatrix} 2 & 2 & 2 \\ 2 & 2 & 2 \\ 2 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 36 & 36 \\ 36 & 36 \end{bmatrix}$

$B * \text{filter 3} = \begin{bmatrix} 3 & 3 & 3 & 3 \\ 3 & 3 & 3 & 3 \\ 3 & 3 & 3 & 3 \end{bmatrix} * \begin{bmatrix} 3 & 3 & 3 \\ 3 & 3 & 3 \\ 3 & 3 & 3 \end{bmatrix} = \begin{bmatrix} 81 & 81 \\ 81 & 81 \end{bmatrix}$

* Both Depthwise Conv & regular Conv yields same result.

ii) MobileNets :-

- i) parameter reduction :- MobileNets uses fewer parameters than other networks by using depthwise separable convolution and by reducing the no. of filters in each layer.
- ii) MobileNets uses low resolution input which reduces the no. of computations required to process the image.
- iii) Small filter sizes :- MobileNets use small filter sizes (3×3) to reduce the no. of computations required by the Conv operation.
- iv) Depthwise separable Convolution :- MobileNets use a depthwise separable Convolution operation which splits a standard Conv operation into two parts - Depthwise Convolution & Pointwise Convolution which helps us to reduce the no. of computations while still achieving good accuracy.

3)

- a) Tasks needed to be achieved during object detection:
 - i) object localization - Determine the spatial coordinates or bounding box using regression
 - ii) object classification - To detect classify an object according to its class in the bounding box.
(i.e. class labels)

b) detected object (2,2) (6,6)
 ground truth (3,3) (7,7)

$$IoU = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$

$$|A| + |B| - |A \cap B|$$

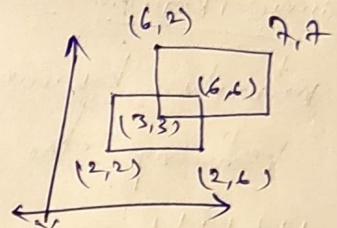
|A|
 ↑
 detected
object

|B|
 ↑
 ground
truth.

$$|A| = 5 \times 5 = 25$$

$$|B| = 5 \times 5 = 25$$

$$\begin{aligned} |A \cap B| &= (6-3)(6-3) \\ &= 4 \times 4 = 16 \end{aligned}$$



$$\frac{16}{50-16} = \frac{16}{34} \frac{8}{17}$$

Jaccard distance = $1 - IoU$

$$= 1 - \frac{8}{17} = 0.529$$

c) To compute Average precision we need to first calculate the precision & recall of each recall value. In this case, the unique recall values are 0.2, 0.4, 0.6, 0.8, 1.

Recall Precision

0.0	1.0
0.2	1.0
0.4	0.6
0.6	0.6
0.8	0.0
1.0	0.0

To compute AP @ 0.5 we need to integrate the precision-recall curve w.r.t to recall at IoU threshold of 0.5. This can be done using the trapezoidal rule.

$$\begin{aligned} AP(0.5) &= (0.2-0)(0+1)/2 + (0.4-0.2)(1+1)/2 + \\ &\quad (0.6-0.4)(0.6+1)/2 + (0.8-0.6)(0.6+0)/2 \\ &= 0.36 \end{aligned}$$

d) Detection box of coordinate :- Normalized between 0 & 1
 boundaries ($x_{\min}, y_{\min}, x_{\max}, y_{\max}$). Since the image is at different scale, relative coordinates is used instead of absolute coordinates.

Scale invariance :- This would make the object detection algorithm scale invariant, which would detect object equally well irrespective of the size.

e) Different terms in the loss function needed for an object detection network

$$\text{loss} = L_{\text{reg}} + L_{\text{cls}}$$

L_{reg} :- loss due to regression of the bounding box
 ex:- MSE

L_{cls} :- loss due to classification of object inside the bounding box
 ex:- Cross entropy

$$L(y_i^{(j)} - \hat{y}_i^{(j)}) = \begin{cases} (y_i^{(j)} - \hat{y}_i^{(j)})^2 & \text{if } y_i^{(j)} = 0 \\ \epsilon (y_i^{(j)} - \hat{y}_i^{(j)})^2 & \text{otherwise} \end{cases}$$

Apart from these one other loss like object loss, box loss etc.

f) Grid cell of 3x3, size of the output tensor - 10 deleted boxes at each cell

$$\text{formula} = 10 \times 5 \times 5 + (k+5)$$

$$= 90 \times (k+5)$$

$$= 90k + 450$$

g) Single Shot vs double shot

- * Single shot detector performs object detection in a single network that looks at image once.
- * Double shot detector look at the image twice - one to propose Object regions which is referred as region proposal and one to refine and clarify regions. This is more accurate but higher computational cost.

h) Object detection loss term :-

center regression

$$\lambda_{\text{coord}} \sum_{i=0}^S \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[\frac{(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2}{a} + \frac{(\text{weight}_i - \hat{\text{weight}}_i)^2 + (\text{height}_i - \hat{\text{height}}_i)^2}{b} \right]$$

weight & height regression loss

$$\lambda_{\text{coord}} \sum_{i=0}^S \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[(\text{cls}_i - \hat{\text{cls}}_i)^2 + (\text{phi}_i - \hat{\text{phi}}_i)^2 \right]$$

object classification loss

$$+ \sum_{i=0}^S \sum_{j=0}^B \mathbb{1}_{ij}^{\text{no obj}} (C_i - \hat{C}_i)^2$$

d - no object classification loss

$$+ \lambda_{\text{no obj}} \sum_{i=0}^S \sum_{j=0}^B \mathbb{1}_{ij}^{\text{no obj}} (L_i - \hat{L}_i)^2$$

e - box classification loss

$$+ \sum_{i=0}^S \sum_{j \in \text{Same}} \mathbb{1}_{ij}^{\text{obj}} (P_i(c) - \hat{P}_i(c))^2$$

i) ROI Pooling

It extracts small feature map of fixed size. It follows a process of region proposal generation, where bounding boxes are generated by the network.

- * Generates a RPN which represents the location of possible objects.
- * Then for each & every subregion - pooling operation is applied reducing the spatial dimensions, preserving the channel dimensions.

We do this because to provide a fixed size representation of variable-sized image, which makes the object detector able to handle objects of different sizes & aspect ratios. This helps in scale variant object detection & also helps in learning the features of the object for different aspect ratio.

i) NMS (Non Maximum Suppression)

RPN :- basically propose the regions where the object has to be detected

Here non maximum suppression is to select coordinate with highest classification & delete candidate overlapping with it.
(IOU > Threshold)

K) RCNN Loss function

$$\text{Loss} = L_{\text{Class}} + L_{\text{reg}} + L_{\text{seg}}$$

$$L_{\text{Class}} = L_{\text{class}}^{\text{obj}} + L_{\text{class}}^{\text{RPN}}$$

$$L_{\text{reg}} = L_{\text{reg}}^{\text{RPN}} + N_{\text{reg}}^{\text{obj}}$$

Q. 4.

a) Semantic Segmentation

In semantic segmentation it classifies each pixel into 1 or more classes. The goal is to group each and every pixel ~~to~~ into ~~one~~ specific class regardless of the object.

Instance Segmentation

In instance segmentation, it assigns unique label to each instance of the object and it provides accurate object boundaries with object localization.

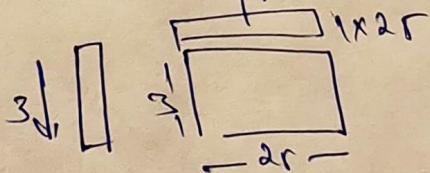
b) regular convolution

$$O/P \leftarrow \frac{n - filter + 2(pad) + 1}{\text{Stride}}$$

$$\frac{5-3}{1} + 1 = 2 + 1$$

$$= 3 \times 3$$

$$O/P \leftarrow 3 \times 3 \quad I/P \leftarrow 25 \times 1$$



$$3 \times 1 = 3 \times 25 \times 25 \times 1$$

$$4 \text{ of } 3 \times 3$$

$$\therefore \text{Size of the map} = 3 \times 25$$

of 3×25

$$\begin{bmatrix} abc \\ def \\ ghi \end{bmatrix}$$

Size of the filter matrix = 3×25

~~25x25~~

- (a) To perform Convolution. On a vectorized image & the filter, we need to slide the filter over the image. For image with dimensions 5×5 & filter size of 3×3 , there are $(5-3)+1 = 3$ possible positions to place the filter both horizontally & vertically. And do no. of operations that can be performed = $25 \times 3 = 75$ operations.

Each operation involves Convolving with 3×3 filter with the part of the region of 3×3 vectorized image that overlaps with the filter. This will produce the convolution result of 9×1 .

The final output will be 3×3 matrix where each pixel/element represents the result of the operation. (Assuming that padding or stride is used).

- (b) To compute the filter of the transpose convolution matrix, we need to reverse the convolution. Instead of sliding the filter over the image, we now slide the output of the convolution over a larger matrix to obtain the original image. The output of convolution will be 3×3 and hence 3×3 by filling the gaps with zeros. Each element of the filter represents the subarea of the original image. We first fill in the corresponding position in transpose convolution matrix with this output to get the transpose convolution. There will be 9 possible positions in 5×5 matrix and so the transpose filter will be 5×5 with 9 channels. We need to flatten this now which results in (25×9) .

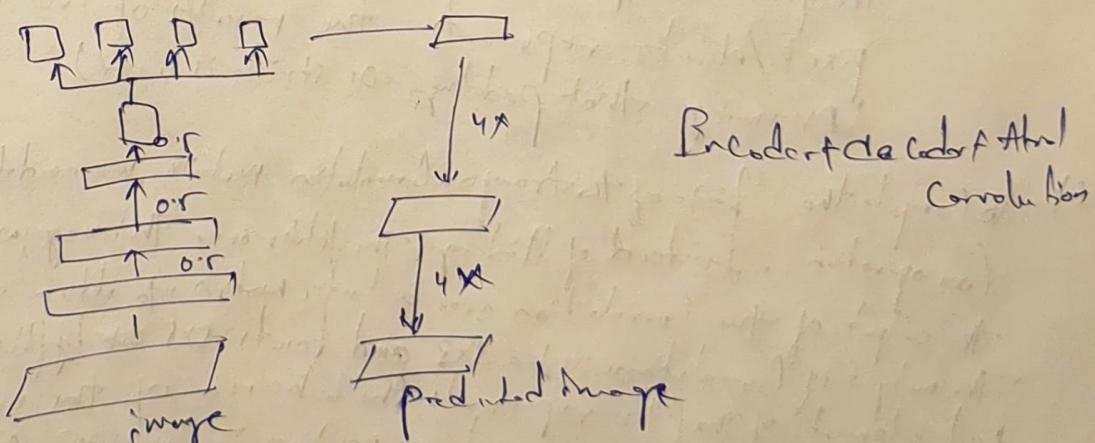
a) Skip Connections

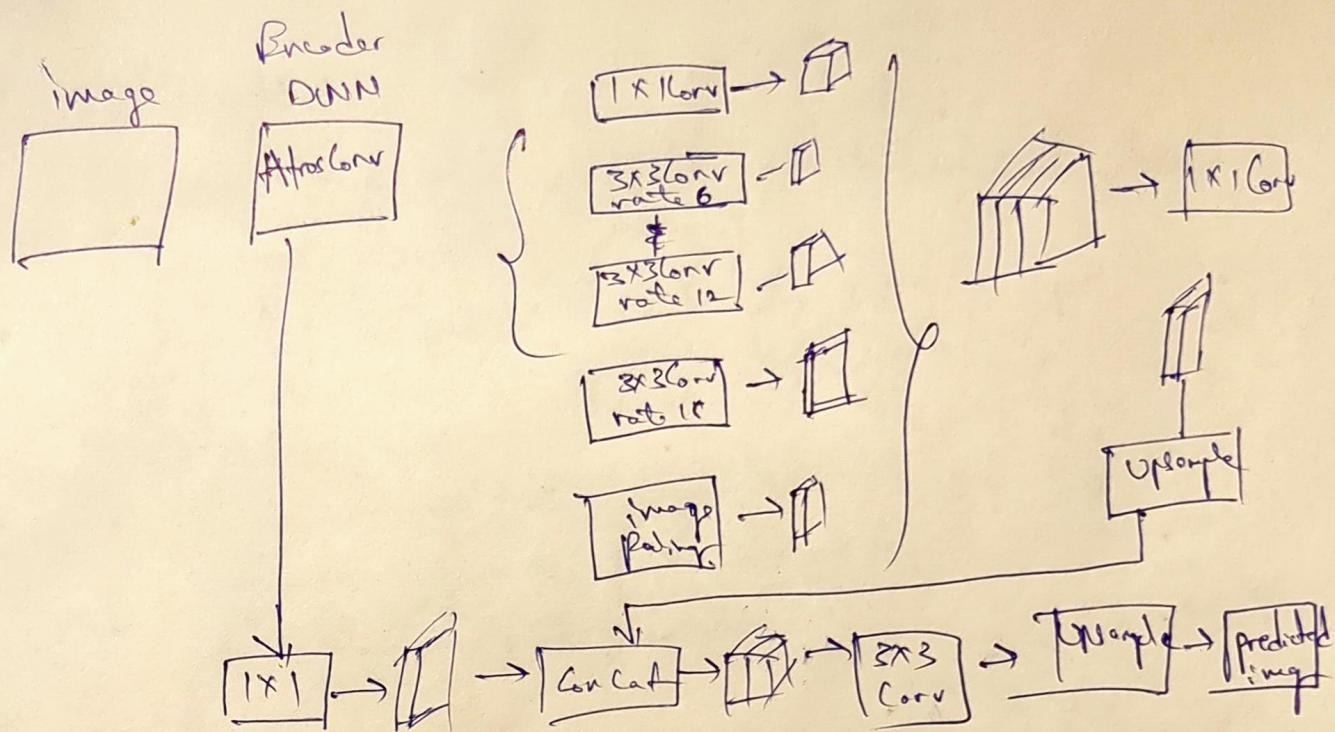
- i) Connection b/w encoder & decoder to its corresponding layer (dimensions should be same.)
- ii) Upsampled output is concatenated with cropped input - loss due to loss of information (pixels) in every convolutional layer.
- iii) Used to address vanishing gradient & improve the information flow & feature reuse which helps in accurate image segmentation.

b) DeepLab Architecture

* Encoder & Decoder

- * Use Atoms spatial pyramid pooling with skip connections.
- * Use depthwise separable convolutions





Multiscale Contextual information by applying other Convolution at multiple scale , the decoder defines the Segmentation result along object boundaries.

f) Convolution matrices - Semantic Segmentation results.

The possible evaluation matrices.

1. IOU

2. pixel accuracy :- no. of pixels that have been segmented into its corresponding class

3. Inference of Classification

* precision

* recall

* f1 score .