

## Convolution Neural Network

Page No.	
Date	

Q] What is Convolution Neural Network CNN  
Discuss applications of CNN.

→ Convolution Neural Network (CNN)  
is type of deep learning model specially designed for processing data with grid-like structure such as image.

- CNN represents the input data in the form of multidimensional data it works well with large number of labeled data.
- The goal of CNN is reduce the image so that it would be easier to process without losing features that are valuable for accurate prediction.
- Let's take the example of human brain as it contains we use the images we classify the image by its color, shape, size and other message it is upto to convey. Similarly for machine image is just the array of pixels. Machine tries to learn the unique pattern in the image & computer tries to find these pattern & information from the image.
- Most of the companies have adopted the CNN for image recognition.
- Machine can be trained by giving tons of images to learn & find patterns.

- (iii) Application of CNN in industry
- .(i) In monitoring environment
  - 1) Image classification
  - 2) Face Recognition
  - 3) Medical image
  - 4) Self-driving cars
  - 5) Object detection (YOLO, R-CNN)
  - 6) Video analysis
  - 7) Gesture Recognition

Q

Explain convolution operation in CNN with example. Explain discrete & circular convolution operation. Take input  $5 \times 5$  input data,  $3 \times 3$  kernel data & calculate convoluted features.

- Convolution operation focus on extracting the feature from input. Convolution operation allows the network to detect horizontal & vertical edges of an image & based on those edges high-level features are build.
- Convolution operation uses three parameters i) Input image ii) feature detector iii) Feature map

• Convolution operation uses input matrix & filter i.e., also known as Kernel. Input matrix are pixel value of greyscale image and filter is small matrix that detect the edges of image.

• Input image is converted to binary 1 & 0 format.

1	0	0	1	1
0	1	1	0	0
0	0	1	1	0
0	0	1	0	1
0	0	1	0	1

1	0	1
0	1	0
1	0	1

Input is padded with Kernel

Step by Step O/P

$$(I \otimes K) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} = 1 + 0 + 0 + 0 + 1 + 0 = 4$$

$$I[0,2] = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} = 3$$

$$\begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} = 3$$

final output feature map (3x3)

$$\begin{bmatrix} 4 & 3 & 3' \\ 3 & 3 & 2 \\ 9 & 2 & 85 \end{bmatrix}$$

### 1] Discrete convolution →

- It is method of combining two discrete signals to produce a third signal.
- In CNN it is used to apply kernels (filters) to for feature extraction.
- Usually zero-padding is used if kernel goes outside the boundary.

$$y[n] = \sum_{m=-\infty}^{\infty} x[m] \cdot h[n-m]$$

- Output depends on padding and stride.

### 2] Circular convolution →

- Works like discrete convolution but wraps around instead of filling with 0.
- Output size is always same as input.

$$\left[ \text{DiscreteConv}[m \bmod N] \right]$$

$$\sum_{m=0}^{N-1}$$

Q) What are pooling layers in CNN? Explain its types with example

- pooling is technique in CNN to reduce the spatial size of the representation. Its main purpose is to decrease the amount of parameters and computation in the network.
- pooling layer takes the output of previous layer & compute summary of neighbourhood around specific location.
- The main purpose of pooling is to-
  - reduce computation by lowering number of values
  - Control overfitting by reducing parameters
  - Reduce memory requirement for storage.
- The pooling layer will always reduce the size by a factor of 2. Suppose you have image of  $(6 \times 6)$  pixels it will be pooled  $(3 \times 3)$  pixels.
- Pooling operation is also called as subsampling.
- Pooling layer is also called as downsampling. It reduces the number of parameters in the input.

There are two types of pooling

- i) Max pooling
- ii) Avg pooling

A) max pooling

- Takes maximum value in the window
- Keeps strongest activations

1 3

2 4

$$O/P = 4$$

b) Average pooling

- Takes average of values in the window
- produce smoother output

1 3

2 4

$$(1+2+3+4)/4 = 2.5$$

c) Explain the following hyperparameters for the convolution layer.

i] filter size

ii] output depth

iii] stride

iv] zero-padding

→ i] filter size

- filter size is the height x width of the matrix sliding over the input image.
- common filter sizes are:

-  $3 \times 3$ ,  $7 \times 7$ ,  $5 \times 5$

• Two types of filters are

- 1) small filter → fine detail (eg edges)
- 2) large filter → broad patterns (eg texture)

- Each filter has values (weights) learned during training.
- Multiple filters detect different shapes.
- Example:

$3 \times 3$  filter

$$\begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}$$

### ii) Output depth map created.

- Number of Feature maps output by a convolutional layer:

- Each feature map = result of 1 filter
- More filters  $\rightarrow$  more features detected.
- Different filter detects diff color, corner, edges.
- Output depth increase with number of filters.
- more depth = more computation.

eg  $\rightarrow 28 \times 28 \times 3 \rightarrow$  RGB image.

filter : 8 filters ( $3 \times 3 \times 3$  each)

O/P  $\rightarrow 26 \times 26 \times 8$  (assume no padding, stride)

### iii) Stride.

~~of stride~~ = number of pixels filter moves horizontally and vertically over the pixels of the input during convolution.

- Stride is component for compression of images & video data. ex if NN stride = 2 the pixel will move by 1 unit at time.

- Stride is smaller if we want fine/high details in our output.
- Stride is higher if we want macro level details.

eg  $\rightarrow$  Input =  $6 \times 6$   
 filter =  $3 \times 3$  no padding  
 stride 1  $\Rightarrow 4 \times 4$  outputs.  
 stride 2  $\Rightarrow 2 \times 2$

### 1) Zero-padding

- padding refers to adding extra pixels around the boundary to an image before applying the convolution.
- without padding the output gets zeroed after each convolution.
- zero padding helps making the output dimension & kernel size dependent.
- Three common zero padding are -

#### 1) Valid convolution $\rightarrow$

- No zero-padding used
- filters only processes areas fully inside the image
- output size is smaller than input

Input =  $5 \times 5$ , filter =  $3 \times 3$ , o/p =  $3 \times 3$

#### 2) Same Convolution

- padding added so output size = ipsize.
- kernel size is K if ip is padded by K-1 zeros

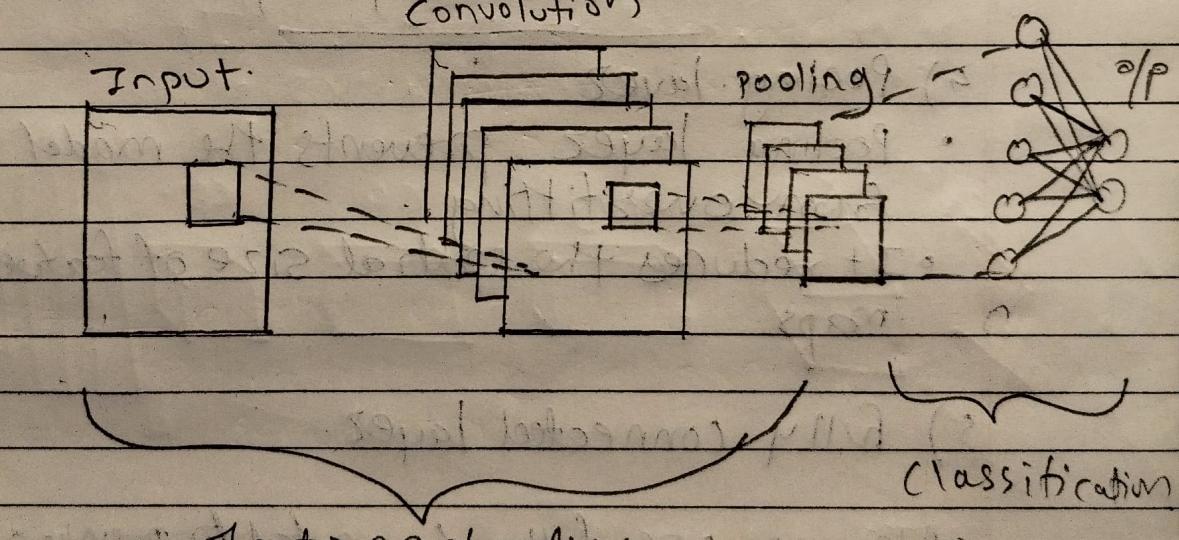
eg  $I/P = 5 \times 5$  filter, filter  $= 3 \times 3 \rightarrow$  Add 1 layer  
of zeros  $\rightarrow O/P = 5 \times 5$

### 3) Full convolution

- padding added so filter can process all possible overlaps including partial ones.
- $O/P > I/P$

eg  $5 \times 5$ ; filter  $= 3 \times 3$  Add 2 layers of zeros  
 $O/P = 7 \times 7$

### Q] Explain Basic architecture of CNN



Q] Input layer: This layer accepts the raw data such as image dimensions  $\rightarrow$  height  $\times$  width  $\times$  depth.  
eg  $\rightarrow 12 \times 12 \times 4$

## 2) Convolution layer.

- Core layer performing convolution using filters / kernels.
- Detects local features like edges, corners & textures
- Filters are learned during training.

## 3) Activation Function Layer

- This layer applies activation  $f^n$  to each element in the output of the convolution layer.
- Common  $f^n$ : ReLU, sigmoid, Tanh.
- Helps CNN learn complex patterns.

## 4) Pooling layer

- Pooling layer prevents the model from overfitting.
- It reduces the spatial size of feature maps.

## 5) Fully connected layer.

- Neuron are fully connected to previous layer's o/p
- Converts extracted features into a final classification output.

to extract meaningful features & words [P]

Function of hidden layers in CNN.

1] Learning Features.

2] Dimensionality Reduction.

- Pooling reduces feature map size.

3] Translation Invariance.

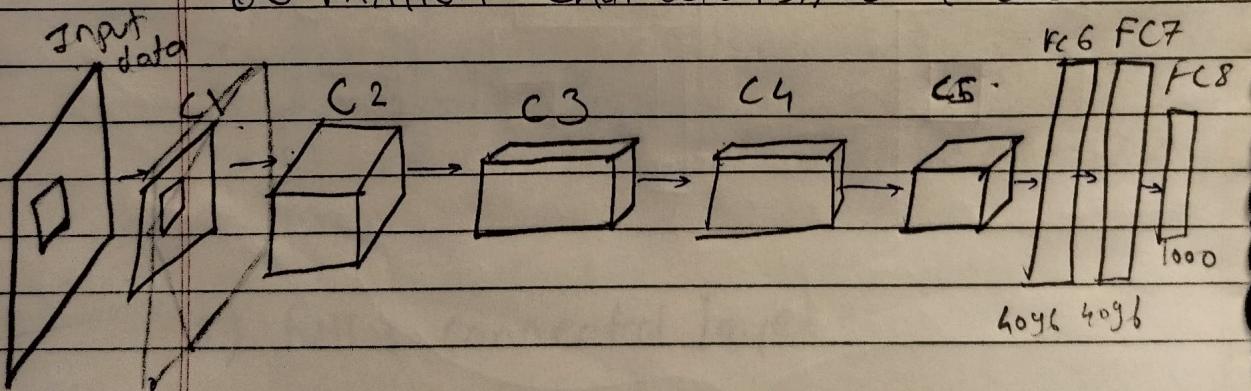
- Pooling ensures detection even if feature slightly shift position.

4] Hierarchical Representation.

- Layer-by-Layer CNN builds a hierarchy.
- low-level  $\rightarrow$  mid-level  $\rightarrow$  high level Features.

Q] Draw & explain architecture of Alexnet

- • Alexnet is Deep convolution NN for image classification
- Alexnet composed of 5 Convolution layers with combination of max-pooling layers, 3 Fully connected layers, & 2 dropout layers
- The Activation function used in all layers is ReLU, & Activation f<sup>n</sup> used in output layer is softmax,
- Alexnet architecture was created with large-scale image datasets in mind, and it produced It has 60 million characteristics of all



1) Input layer → accepts  $227 \times 227 \times 3$  image.

2) Conv layers → There are five convolution layers in total. These five layers use convolution kernels to generate feature maps.

- 3) Max-pooling layers → Max pooling layer is applied bet<sup>n</sup> every two convolutional layers
- 4) Fully connected layers → After the convolutional & pooling layers the tensors are flattened and fed into three fully connected layers
- 5) Dropout layers: Two dropout layers are used to reduce overfitting which was a major issue due to model have 60M parameters
- 6) Output layer → Last layer is softmax layer which computes the loss f<sup>n</sup> for learning & provides the classification output.