

→ The Autoencoder can effectively summarize data points when all relevant info is present in the individual data points themselves.

## Chapter 1

### 1) Vanilla Deep Neural networks

Deep learning aims to automate feature selection, crucial for Computer Vision tasks.

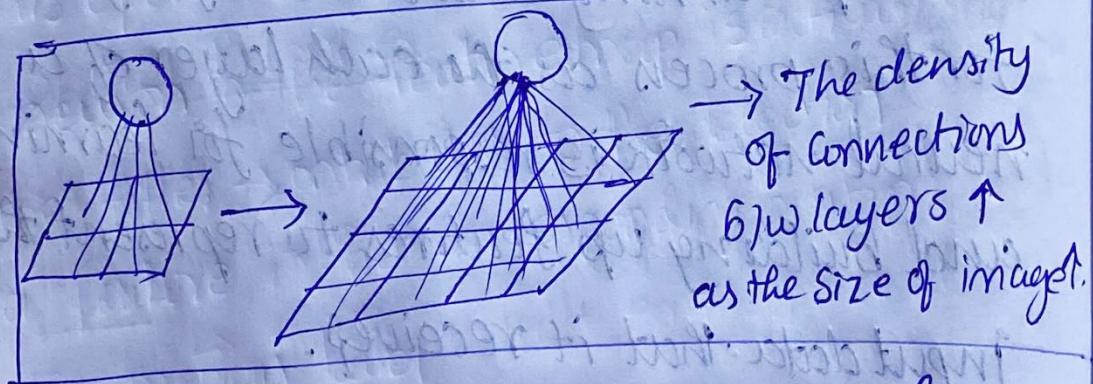
The Deep learning network are prefects for this process as the each layer of a neural network is responsible for learning and building up features to represent the input data that it receives.

The Vanilla Deep neural network are not scalable for image classification as image size increases due to an exponential increase in number of weights.

In the MNIST dataset, the problem appears in manageable with  $28 \times 28$  Pixel images, but Scaling up to larger images rapidly increases the complexity.

For large image, such as  $200 \times 200$  full-color image, The no of weights required for just the input layer becomes high, leading to high risk of overfitting and inefficiency.

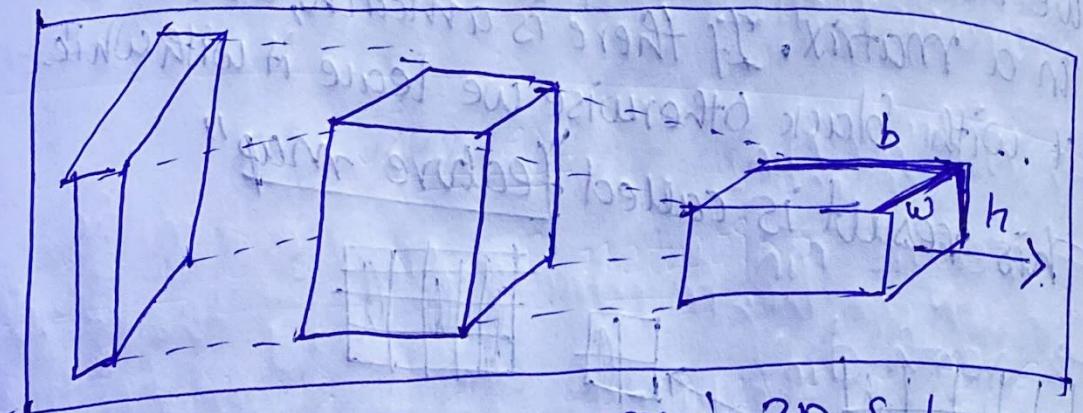
The CNN address this Scalability issue by arranging neurons in 3D and connecting each neuron to only a small, localized region of the previous layer.



This architecture mimics human visual processing and reduce the no of parameters significantly.

CNN manage the complexity by processing info in a 3D; So layers have width, length & depth, so new volumes of processed info layer by layer.

This design of CNN layer avoids the problem of full connectivity and reduce the potential for overfitting while handling larger images more efficiently.



CNN layer arrange neurons in 3D, so layers have width, height; depth.

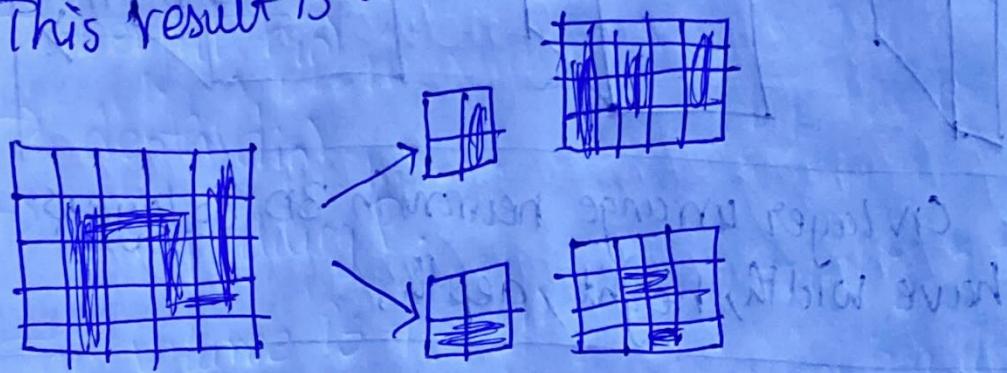
## Filters and Feature Map:

The visual cortex layers structures, where each layer builds upon the features detected by the previous one, inspired the design of CNN layers in neural networks.

The filters is essentially a feature detector that slide across the image to identify features like edges, lines & angles.

The

To detect Vertical lines, we would use the feature detectors on the top, slide it across the image and at each and every step check if we have a match. We keep track of answer in a matrix. If there is a match, we shade it with white otherwise we leave it with black. This result is called "feature map".

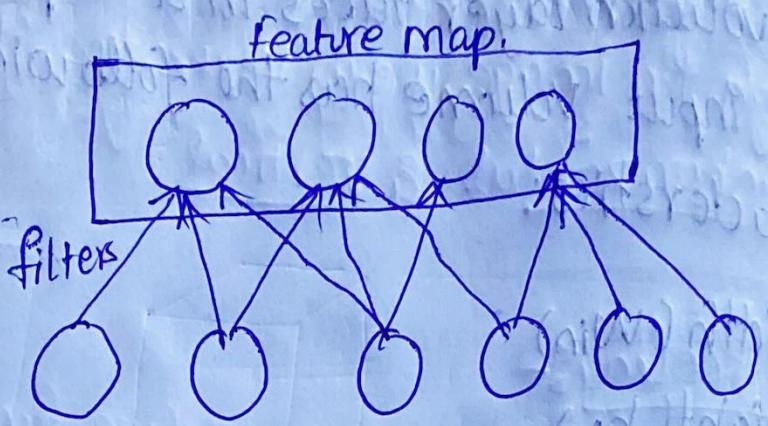


This operation is called Convolution. A neuron in the feature map is activated if the filter contributes to the detection of a feature in the corresponding position in the previous layer.

Expressing the Convolution Operation as neurons in a network

Layers of neurons in a feed forward neural net represents either the original

image or a feature map that get replicated across the entirety of the input.



Let's denote the  $k^{th}$  feature map in layer  $m$  as  $m^k$ .  $W \rightarrow$  corresponding filters weight.  $b^k \rightarrow$  bias. we mathematically Express feature map as:

$$m_{ij}^k = f((W * x)_{ij} + b^k)$$

The fundamental operation of Sliding filters across the input data to produce feature maps, capturing more complex features. They capture patterns like edges or textures.

Filters in CNN have depth, allowing them to process multi-channel input, such as colour images.

The CNN converts one volume of values into another volume of values.

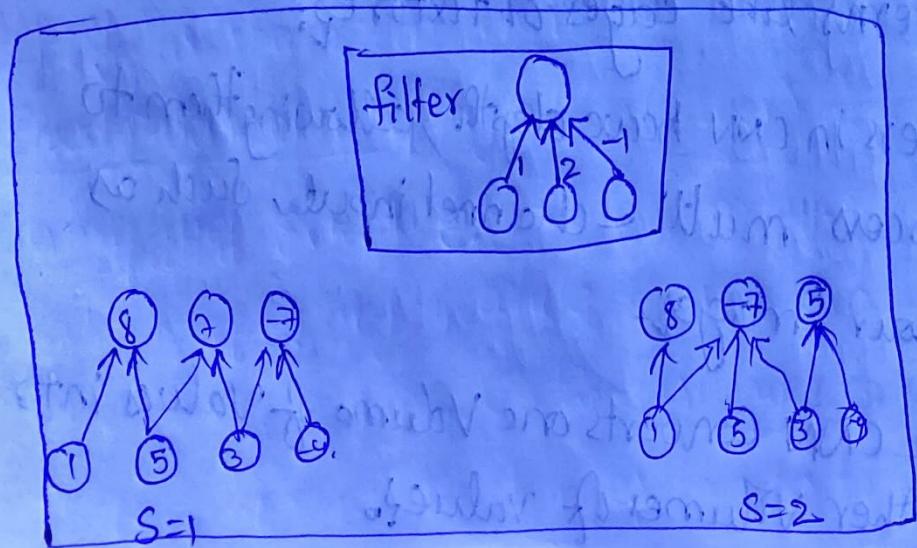
# Full Description of Convolution layer

A Convolution layer takes in a Input Volume.  
This input volume has the following  
Characteristics:

- Width (Win)
- height (hin)
- depth (din)
- Zero padding P.

This volume is processed by a total of  $K$  filters,  
This filters has many no. of hyper parameters:

- Spatial Extent e. → which is equal to the filters height & width.
- Stride S → distance b/w consecutive applications of the filters on the input volume.
- bias b, → which is added to each component of the convolution.



The Result in the output volume has :

→ function  $f$ ,

$$\rightarrow \text{Width } W_{\text{out}} = \left\lceil \frac{W_{\text{in}} - e + 2p}{s} \right\rceil + 1$$

$$\rightarrow \text{height } h_{\text{out}} = \left\lceil \frac{h_{\text{in}} - e + 2p}{s} \right\rceil + 1$$

→ depth  $d_{\text{out}} = k$

$$\begin{aligned} w &= 5 \\ h &= 5 \\ d &= 3 \\ \text{Zero padding} &= 1 \\ w &= 3 \\ h &= 3 \\ d &= 2 \end{aligned}$$

Padding:

If the input is  $n \times n$  and filter size is  $f \times f$   
the output size will be  $(n-f+1) \times (n-f+1)$

There are two disadvantages:

→ Every time we apply a convolution Operation, the  
Image size shrinks

→ Pixels present in the corner of the image are used  
only a few number of times during convolutional  
as compared to center pixels.

To overcome the issue we can pad the image with  
additional borders. We add one pixel around the  
edges. This means input size will be  $9 \times 9$  instead  
of  $7 \times 7$ .

input :  $n \times n$

padding :  $p$

filter size :  $f \times f$

output :  $(n+2p-f+1) \times$   
 $(n+2p-f+1)$

## Striding

No of steps taken for convolution.  
Suppose we choose a stride of 2. So while convoluting through the image, we will take two steps: both in horizontal & vertical.

## Max Pooling:-

It returns the maximum value from the position of the image covered by the kernel.

~~we can describe~~

The pooling layer is responsible for Reducing the spatial size of the convolved feature.

This is to decrease the computational power required to process the data, through dimensionality reduction, and hence speed up the computation.

## two parameters:

1) Spatial Extent "C".

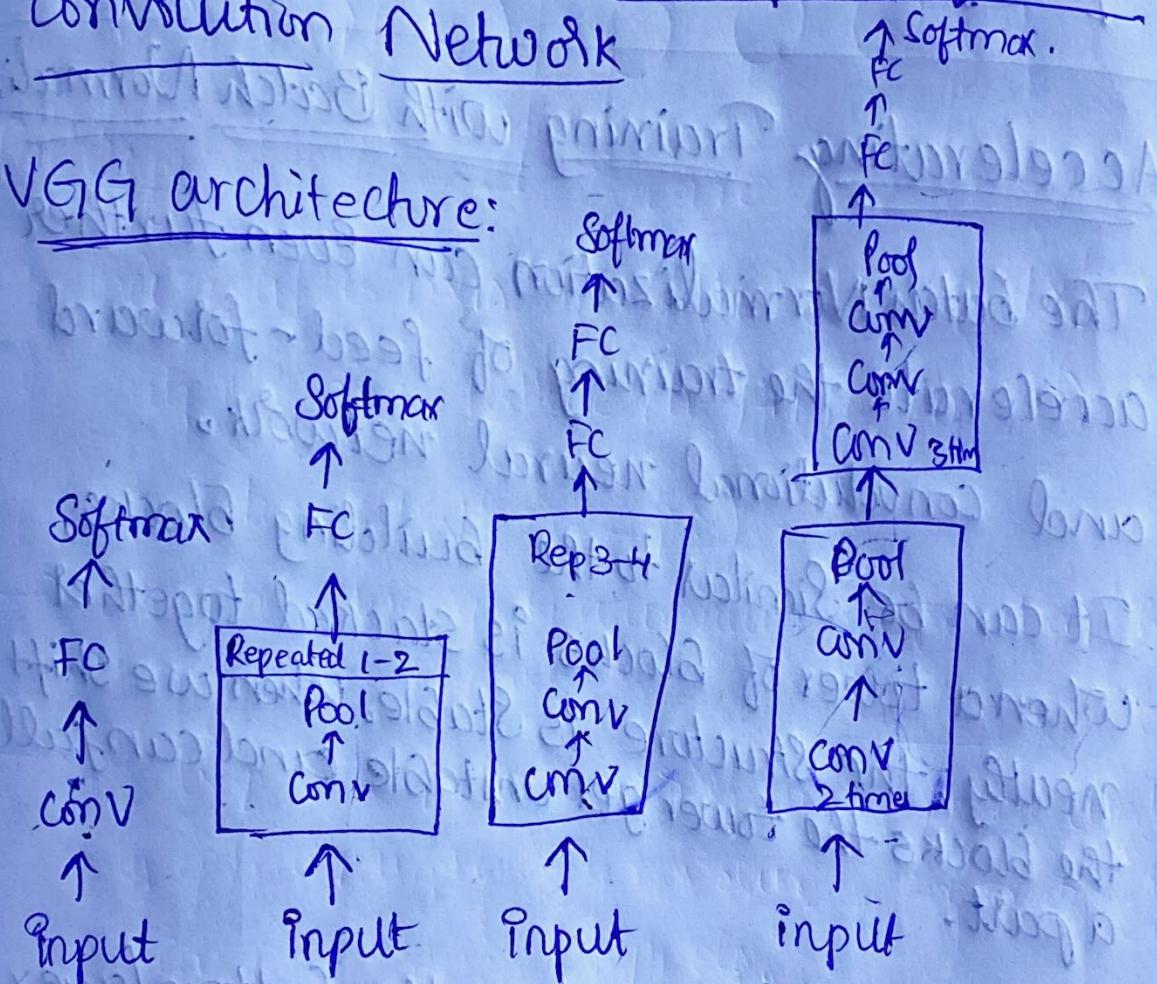
2) Stride "S".

$$\text{Width } W_{\text{out}} = \left\lceil \frac{W_{\text{in}} - e}{s} \right\rceil + 1$$

$$\text{Height } H_{\text{out}} = \left\lceil \frac{H_{\text{in}} - e}{s} \right\rceil + 1$$

## Full Architectural Description of Convolution Network

### VGG architecture:



$\text{FC} \rightarrow \text{Flatten}$

The VGG architecture is a deep convolution neural network known for its simplicity and its deep stack of convolution layers with small receptive fields and large strides.

Small filters.

Input layer: The input layer would be size  $28 \times 28 \times 3$ , assuming the images are RGB.

Link

## Accelerating Training with Batch Normalization

The batch Normalization can even further accelerate the training of feed-forward and convolutional neural network.

It can be similar to the building blocks; when a tower of blocks is stacked together neatly, the structure is stable. When we shift the blocks the tower get unstable and can fall apart.

A similar phenomenon where the bottom layer (first layer) and top layer (last layer) happens, when we give data to the bottom layer the data keep changing while passing from bottom layer to upper layers. So the last layer has to keep adapting to the changes in order to give

## UNIT 3

⑫ Collect output so the more the layers we have the more problem we get.

Normalization of image inputs helps out the training process by making it more robust to variations. Batch Normalization takes this a step further by normalizing inputs to every layer in our neural network.

We modify our network to include the below operations for batch Normalization:

- 1) Grab the vectors of bias which is incoming to a layer before they pass through the nonlinearity.
- 2) Normalize each component of the vector of bias across all examples of the minibatch  
~~by subtracting mean.~~
- 3) Given Normalized input  $\hat{x}$ , use an affine transform to restore representation power with 2 vectors of parameters:  $\gamma \hat{x} + \beta$ .

Batch Normalization also allows us to significantly  $\uparrow$  the learning rate.