

①

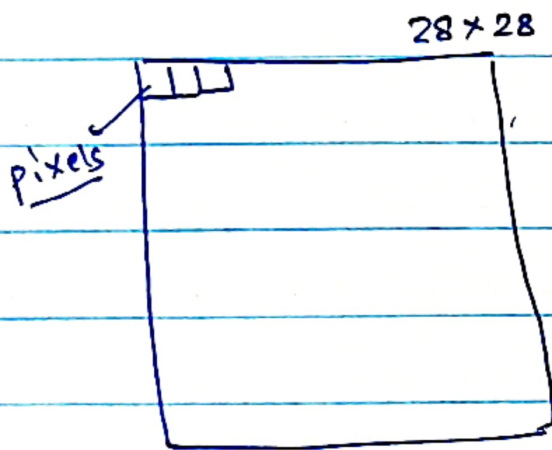
* Images - Pixels.

Why CNN?

* For a simple neural network we have to train a ~~lot~~ lot of parameters. Number of parameters are quite high.

↓
It's very hard to train any network and hence we need to reduce the number of parameters.

Image



→ When flattening happens
↓
Converts this image into single 1D but image was a 2D information. If we pass the 2D information in 1D.
↓
ANN doesn't learn these features.

A similar issue happens in NLP. for eg "I am going to market"

↓
we need to capture the content and not simply convert it to numbers.

- ① Patterns | Edges : when we zoom the images and are able to see some pattern and edges.
- ② Objects : details in an image, person, etc
- ③ Parts of objects : eyes nose ears of the person
- ④ Image : complete information that we process
- ⑤ Scene : multiple images \rightarrow video scene

Edges \rightarrow Pattern \rightarrow Parts of object \rightarrow Object \rightarrow Image \rightarrow Scene

0	0	5	0	0
0	5	18	32	0
0	18	100	64	0
0	100	32	4	0
0	1	2	0	0

If I want to capture this information, i.e. spatial information we cannot capture this if we are performing flattening operation.

Image processing by brain

$V_1, V_2, V_3, V_3A, V_4, V_5$ \rightarrow these layers are a network

we see edges - pattern - part of objects \rightarrow
Object - image

Colors / ① RGB

② CMYK

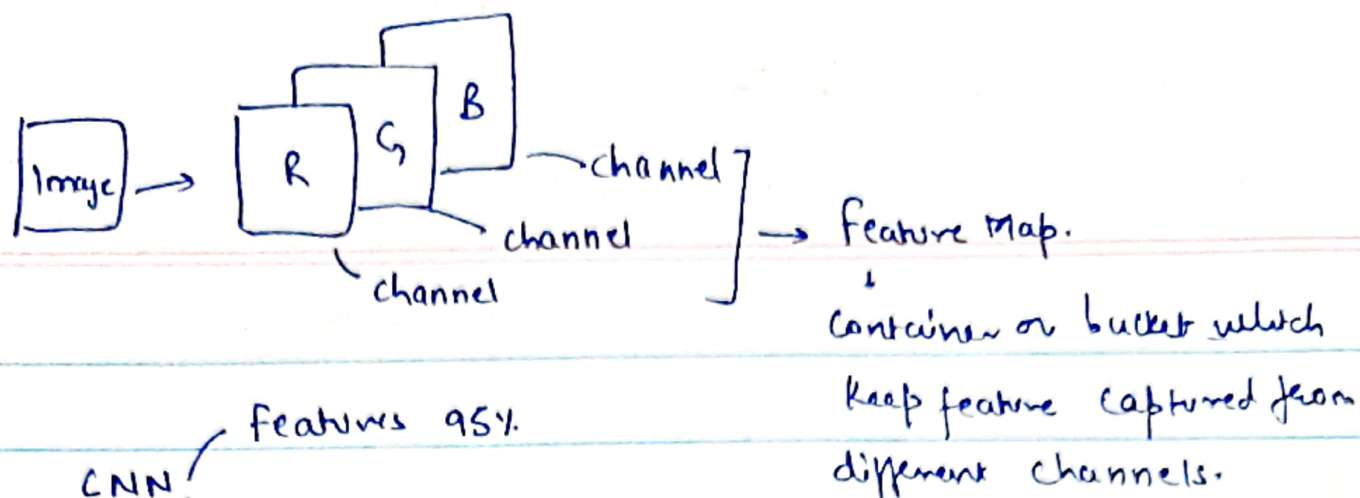
└ News paper: Printing

good for
human
reading

\rightarrow This color space is good for reading

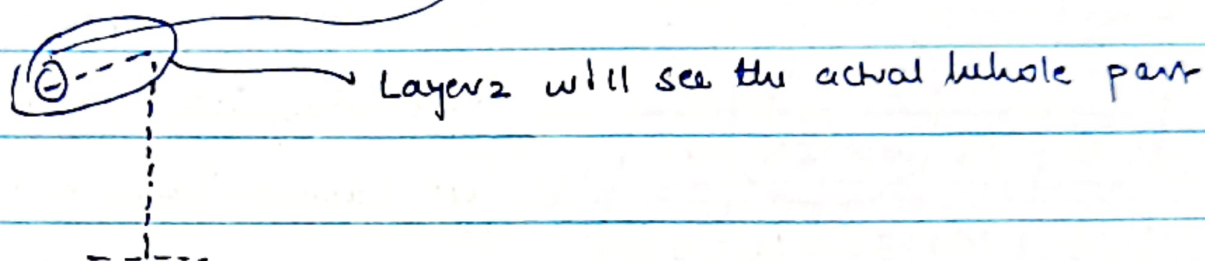
α - Transparency

③



CNN { Features 95%
Colours 5% }

Early layers of CNN — ① Feature Extraction - Similar features.



Kernels

- ① 3×3 : Most optimised kernel till now
 - ② 1×1 : point wise convolution - not used for feature extraction but rather for dimensionality reduction
 - ③ 5×5
 - ④ 7×7
 - ⑤ 11×11
- as we can see we are just writing odd numbers - why?
- There is a symmetry issue. \rightarrow This is a experimental idea

$7 \times 7 \xrightarrow{(3 \times 3)} 5 \times 5 \xrightarrow{(3 \times 3)} 3 \times 3 \xrightarrow{(3 \times 3)} 1 \times 1$

More your kernel looks at the pixels, more its going to extract information from the image.

Most optimised filter: 3×3 - b/c it has least number of parameters (4)
↳ NVIDIA has optimised the GPU to respond to 3×3
Very fast

Kernel may not be a square matrix but also something like 3×1 or 1×3 matrix but experimentally this has been proved.

Always define the number of filters to the tune of 2^n . Similar things goes for images 512×512 , etc.

1) Number of filters = number of feature maps.

Parameter sharing:- same weights can be shared for two similar objects. for eg. parameters will be same for both the eyes.

Basic Components of CNN

1. Filters/ Kernels/ Feature Extractors.
2. Channels.
3. Feature Maps
4. Stride
5. Padding

Basic Components
of Convolution
Operations

6. Receptive Field

7. Output Dimension Size

8. Max Pooling

9. Network design

10. Example of MNIST

Receptive Field

Let's there be an image (I) which is 7×7

$$I_{7 \times 7} \xrightarrow[k_1 (3 \times 3)]{\downarrow} I_{5 \times 5} \xrightarrow[k_2 (3 \times 3)]{\quad} I_{3 \times 3}$$

filters are odd to maintain
symmetry around origin

1. local receptive fields $\rightarrow 7 \times 7 \ K_1$

2. global receptive fields $\rightarrow 5 \times 5 \ K_2$

$\hookrightarrow K_2 \ 5 \times 7 + 7 \times 7$

In local receptive field Kernel can only see the preceding image but not beyond that. Global receptive field could see beyond the preceding image. Eg. K_2 in a local receptive field can only see 5×5 image but K_2 in global receptive field can see both 5×5 and 7×7 image. ②

Designing a Network

Isize \rightarrow Osize :

$200 \times 200 \xrightarrow{K(3 \times 3)} (1 \times 1)$

① Try not to use many layers. Not more than 15

② Generally layers. b/w 5-15

③ Rare cases 15-30

My entire image has been seen by the network

close to 100 convolution layers.

There will be huge number of parameters
↑
higher training time and cumbersome

↓
but there is no fixed rule if image has higher resolution then one needs more layers.

You can increase the size of kernel to reduce the image size quickly and have less layers.

Max Pooling Layer : It reduces the size quickly and reduce the number of layers. (3)

By default it reduce the image by 2

eg 50×50 maxpooling 25×25

Convention:- Do not use two maxpooling layer simultaneously.
So, to reduce the dimension of image quickly than Max Pooling is our best friend. Filters/Mask are needed to save location of the highest feature.

Input Image \rightarrow Output Image

① No stride No padding

② Only padding

③ Only stride

$$n_{\text{output}} = (n_{\text{in}} - f + 1) \quad \text{--- ①}$$

\downarrow \downarrow \downarrow
Output size Input size Filter size

In deep learning we mostly use zero padding and tend to avoid the reflective padding \rightarrow discontinuity of information

Two types of padding \rightarrow Valid - Normal, no padding (Image Reduced)
 \rightarrow Same - Input = Output, No Redn.

4

Cases: ~~only stride~~ $\text{padding} + \text{stride}$ \rightarrow easy computation
 \rightarrow skip pixels.

Hyperparameter:

- Kernel size - 3×3
- Number of Convolution Layers.
- stride - 2 $\begin{matrix} 5 \times 5 \\ 7 \times 7 \end{matrix}$ \hookrightarrow Image Reduction
- Padding - Same, valid
- Epochs.
- Dense layer / Number of Neurons $\rightarrow 64, 128$
- Activation function - ReLU
- Loss function \rightarrow Categorical
- Optimiser \rightarrow Adam, In CV (SGD)
- Image Resolution