

# Convolutional neural network (CNN)

## CNN

CNNs are specialized deep learning architectures designed primarily for image and spatial data. They mimic the process how human brain process the visual informations.

## Structure

- ⇒ Convolutional Layers : Extract features using filter (Mostly edges)
- ⇒ Pooling Layers : Downsample feature maps to reduce dimensionality.
- ⇒ Fully connected Layers : Combine extracted features for classification.

## Why important and advantages

- ⇒ Fewer parameters → less overfitting
- ⇒ Able to capture spatial hierarchy like edges → Shape → Objects.

## Applications

- ⇒ Image classification
- ⇒ Object detection
- ⇒ Medical imaging
- ⇒ Face recognition and many more.

(MS) Answer for question 3

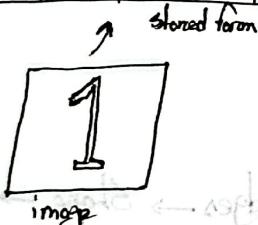
## Convolution operation

Convolution is a mathematical operation that combines two functions. In CNNs, this means sliding a filter (kernel) over the input image

to produce a feature map.

Put filter on the image. Multiply every element that filter covers and take sum. Write on the feature map.

(Input image)					
0	0	0	0	0	0
0	0	0	0.7	0	0
3	0	0	0.4	0.9	0
0	0	0	0.9	0	0
0	0	0	0.9	0	0
0	0	0.5	0.9	0.5	0



\*

-1	-1	-1
0	0	0
1	1	1

=

0	0
0	0

Filter

Feature map

This will be calculated while training

$$\text{output size} = \frac{n-f+2p}{s} + 1$$

|  $n$  = input size  
 $f$  = filter size  
 $p$  = padding  
 $s$  = stride

Visit deep lizard website  
to understand practically

## Padding and Strides

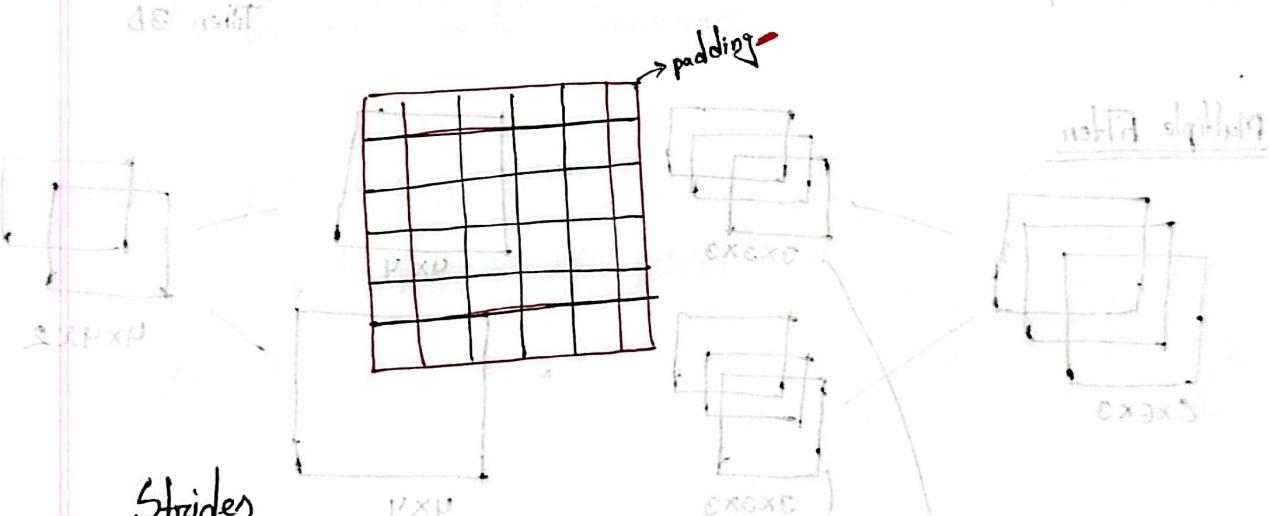
openai 3001 after padding

### Padding

While applying filter on an image its dimension reduces. For this, we may loss information. So, to avoide this reduction, padding is introduced, which adds some extra pixels (usually zeros) around the images edges.

⇒ Valid padding → No padding, output smaller

⇒ Same padding → Input size = Output size



### Strides

If determines how far the filter moves in each step.

» Stride = 1 : Detailed, slow

» Stride  $> 1$  : Faster but less detailed

## Working with RGB image

As digital images are RGB that can form any color in the picture,

so, in 2D input image cannot do the thing. in CNNs. That's why

we make three channels each for R, G and B.

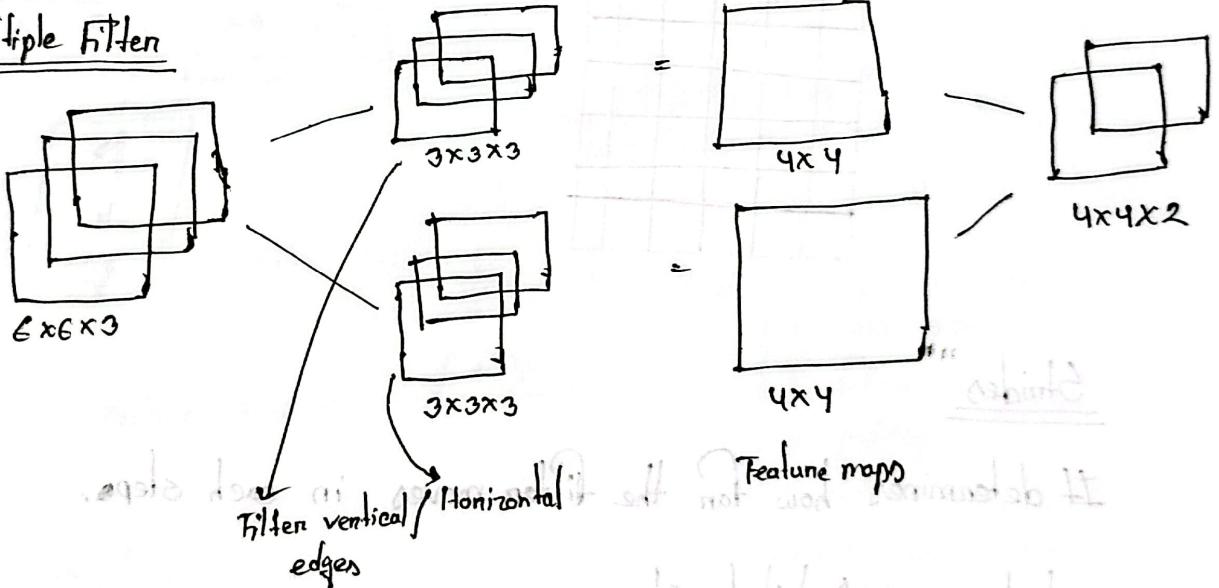
We can also, or we must do the same for filters. Three filter

for three channels.

More formally, the dimension increased from 2D to 3D. Input 3D

Filter 3D

## Multiple Filter



## Problems with convolution

### ① High Computing cost

=> many parameters and multiplication  $\rightarrow$  slow and memory heavy.

### ② Overfitting

### ③ Translation variance

Small shift (translations) in the input image causes different output -  
it doesn't perfectly recognize the object if it moves slightly.

So, we need translation invariance.

S	S
P	P

(S,S)-rule  
(S,P)-rule

$\leftarrow$  padding required

S	L	T	S
S	O	Z	S
I	S	P	I
P	S	T	P

## Pooling Layers

Reduces the spatial size of a feature map  $\rightarrow$  makes computation efficient and prevent overfitting.

### Types

Max pooling \*

Min pooling

Avg pooling

Global pooling

### Effects

$\Rightarrow$  Reduces dimensionality

$\Rightarrow$  increase in variance to small shift

$\Rightarrow$  generalize

map			
5	3		
3	1	1	3
2	5	0	2
1	4	2	1
7	4	2	4
			9

size = (2,2)

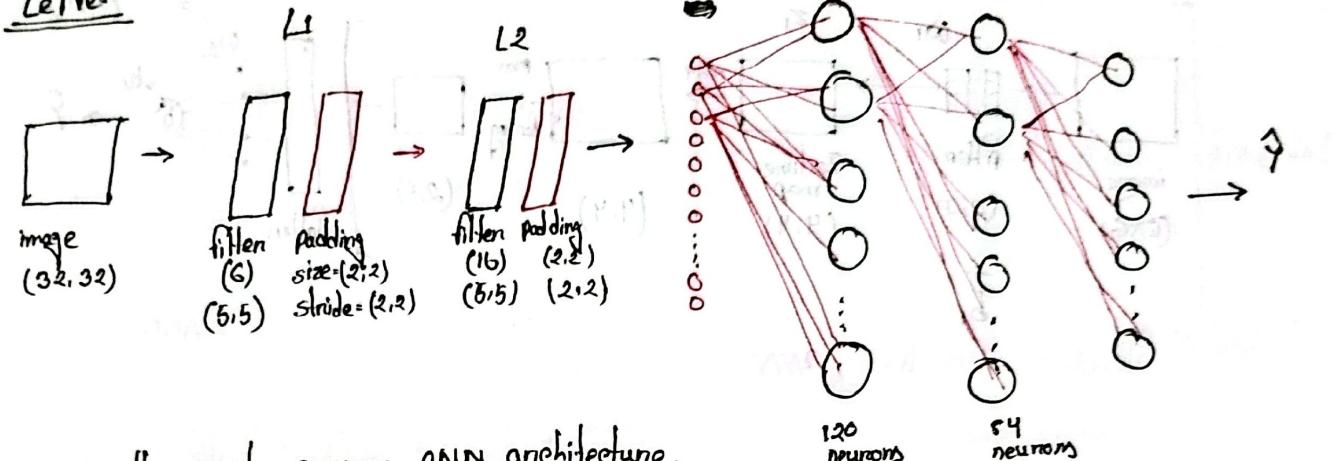
stride = (2,2)

type  $\rightarrow$  max pooling  $\rightarrow$

5	3
7	4

## CNN architecture

### LeNet



This is the most common CNN architecture.

### Demo application on mnist

[model = Sequential ()]

[model.add(Conv2D(6, kernel\_size=(5,5), padding='valid', activation='tanh', input\_dim=(32,32,1)))]  
[model.add(AveragePooling2D(pool\_size=(2,2), strides=2, padding='valid'))]

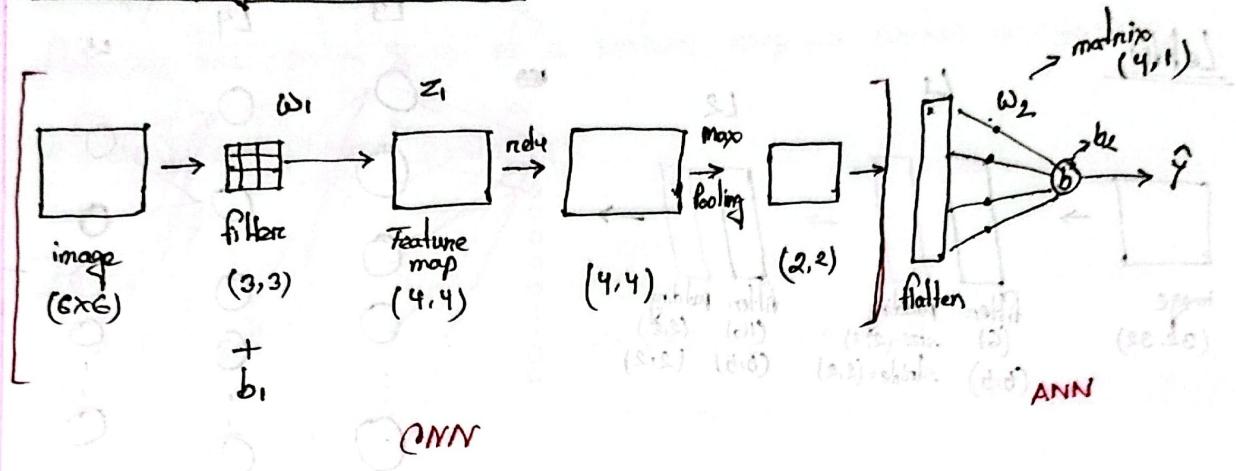
{ model.add(Conv2D(16, kernel\_size=(5,5), padding='valid', activation='tanh')) }  
{ model.add(AveragePooling2D(pool\_size=(2,2), strides=2, padding='valid')) }

model.flatten

[model.add(Flatten())]

[model.add(Dense(120, activation='tanh'))]  
[model.add(Dense(84, activation='tanh'))]  
[model.add(Dense(10, activation='softmax'))]

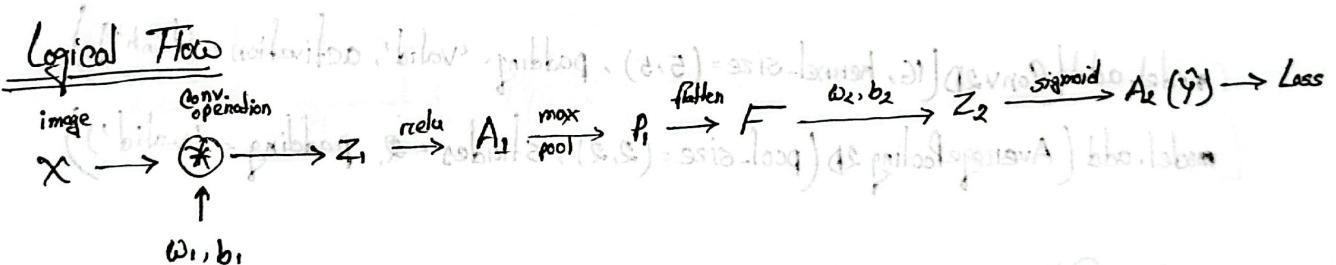
## Backpropagation in CNN



## Trainable parameters

$$\omega_1 = (3, 3), \quad \omega_2 = (4, 1) = 15$$

$$b_1 = (1, 1), \quad b_2 = (1, 1)$$



## Forward propagation

$$Z_1 = \text{conv}(x, \omega_1 + b_1)$$

$$A_1 = \text{relu}(Z_1)$$

$$P_1 = \text{maxpool}(A_1)$$

$$\begin{aligned}
 F &= \text{flatten}(P_1) \\
 Z_2 &= \omega_2^T F + b_2 \\
 A_2 &= \sigma(Z_2)
 \end{aligned}$$

From the ANN part, If we want to calculate  $\omega_2$  and  $b_2$ , process will be,

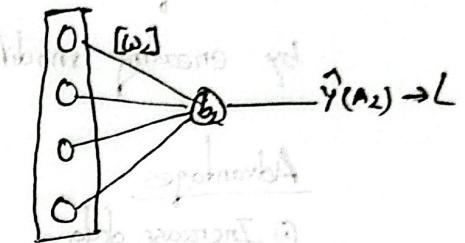
$$\omega_2 = \omega_2 - \eta \boxed{\frac{\partial L}{\partial \omega_2}}$$

$$b_2 = b_2 - \eta \boxed{\frac{\partial L}{\partial b_2}}$$

$$\frac{\partial L}{\partial \omega_2} = \frac{\partial L}{\partial A_2} \times \frac{\partial A_2}{\partial Z_2} \times \frac{\partial Z_2}{\partial \omega_2}$$

$$= (A_2 - \hat{y}) F^T$$

$$\frac{\partial L}{\partial b_2} = \cancel{(A_2 - \hat{y})} \quad \frac{\partial L}{\partial A_2} \times \frac{\partial A_2}{\partial Z_2} \times \frac{\partial Z_2}{\partial b_2}$$



$$L = -y_i \log(\hat{y}_i) - (1-y_i) \log(1-\hat{y}_i)$$

$y$  = actual output

$$= (A_2 - \hat{y})$$

For the CNN part,

$$\omega_1 = \omega_1 - \eta \boxed{\frac{\partial L}{\partial \omega_1}}$$

$$b_1 = b_1 - \eta \boxed{\frac{\partial L}{\partial b_1}}$$

$$\frac{\partial L}{\partial \omega_1} = \frac{\partial L}{\partial A_2} \times \frac{\partial A_2}{\partial Z_2} \times \frac{\partial Z_2}{\partial F} \times \frac{\partial F}{\partial P_1} \times \frac{\partial P_1}{\partial A_1} \times \frac{\partial A_1}{\partial Z_1} \times \frac{\partial Z_1}{\partial \omega_1}$$

$$= \text{conv}(x, \frac{\partial L}{\partial Z_1})$$

( $Z_1$  is a matrix)

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial A_2} \times \frac{\partial A_2}{\partial Z_2} \frac{\partial Z_2}{\partial F} \frac{\partial F}{\partial P_1} \frac{\partial P_1}{\partial A_1} \frac{\partial A_1}{\partial Z_1} \frac{\partial Z_1}{\partial b_1}$$

$$= \text{sum} \left( \frac{\partial L}{\partial Z_1} \right)$$

with batch size N

## Data Augmentation

Data augmentation is the process of artificially expanding a training dataset by creating modified version of the existing images.

### Advantages

- ① Increase data
- ② Reduce overfitting

from keras ~~layers~~ import layers, models

data-augmentation = ~~model~~ Sequential([

layers.RandomFlip("horizontal"),

layers.RandomRotation(0.2),

layers.RandomZoom(0.2),

layers.RandomContrast(0.1)).

model = ~~model~~ Sequential()

model.add(data-augmentation)

model.add(layers.Rescaling(1./255)) // normalize

// Create model now

Pretrained Models

Pretrained models are neural networks already trained on large benchmarks dataset (like Imagenet) then reused to other tasks.

(CIFAR-100 + aplo. bign. celeb - pt. celeb, ImageNet + cifar) → Base model  
celebA  
cifar10

Adv

- ⇒ Faster training
- ⇒ Less data required
- ⇒ Higher accuracy
- ⇒ Avoid overfitting

Common pretrained CNN architectures

<u>Model</u>	<u>Parameter</u>	<u>Use case</u>
VGG 16 / VGG 19	~138M	Easy to Fine-tune, slower
ResNet50/101	~25M	Very Stable and accurate
InceptionV3	~23M	High performance, Fast
MobileNet	~4M - 7M	For edge device

```
from keras.applications import VGG16, ResNet50, InceptionV3  
from keras import layers, models
```

choose model of base with (base, add) bracket

base-model = VGG16 (weights = 'imagenet', include\_top = False, input\_shape = (224, 224, 3))  
ResNet50  
InceptionV3

```
base-model.trainable = False # Freeze convolutional base
```

```
model = models.Sequential([
```

base-model,  
layers.Flatten(),  
layers.Dense(128, activation = 'relu'),  
layers.Dense(1, activation = 'sigmoid')

```
model = models.Sequential([
```

→ data-augmentation,  
layers.Rescaling(1/255)

base-model,

:  
:      VGG  
:      ResNet  
:      Inception

])

Best practice

VGG → Flatten()

ResNet / Inception → GlobalAveragePooling2D