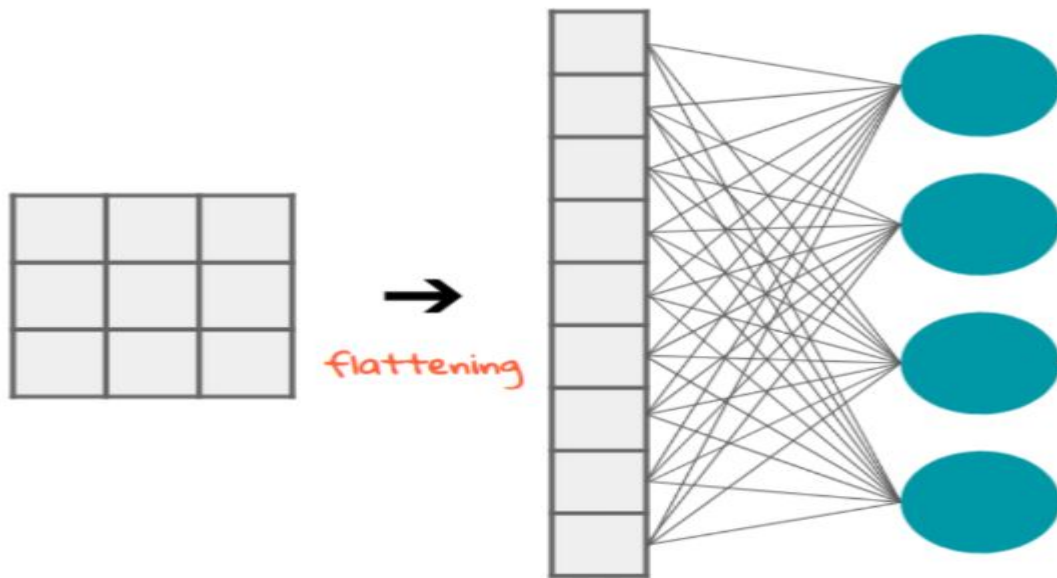# CNN

# Why not use ANN?

- High computational cost
  - → Millions of parameters for even small images
- Severe overfitting
  - → Too many weights, not enough data
- Loss of important spatial information
  - → Treats pixels as independent, ignores 2D structure

When we flatten 2D to 1D input weight increase and training time increase and it causes to loss of spatial arrangement of pixels


flattening

# What is CNN

Convolutional Neural Networks (CNNs), also known as ConvNets, are a specialized type of neural network designed to process data with a known grid-like structure—such as images, videos, and time series.

A CNN typically consists of three main types of layers:

1. Convolutional layers – to extract features

2. Pooling layers – to reduce spatial dimensions

3. Fully connected (FC) layers – to classify the final features

CNNs learn hierarchical features across layers:

- Early layers detect low-level features (e.g., edges, corners, gradients)
- Middle layers combine these into mid-level features (e.g., eyes, noses, wheels, textures)
- Deep layers recognize high-level semantic objects (e.g., faces, cars, dogs)
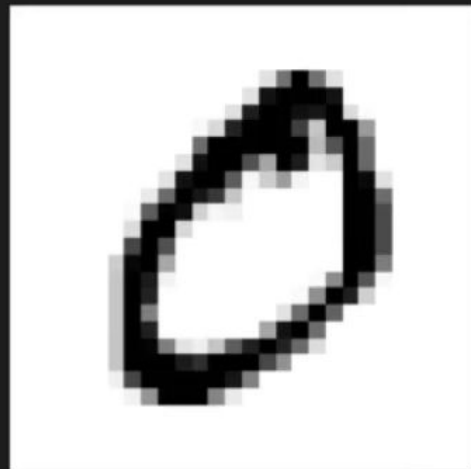
CNNs can process different types of image inputs, most commonly:
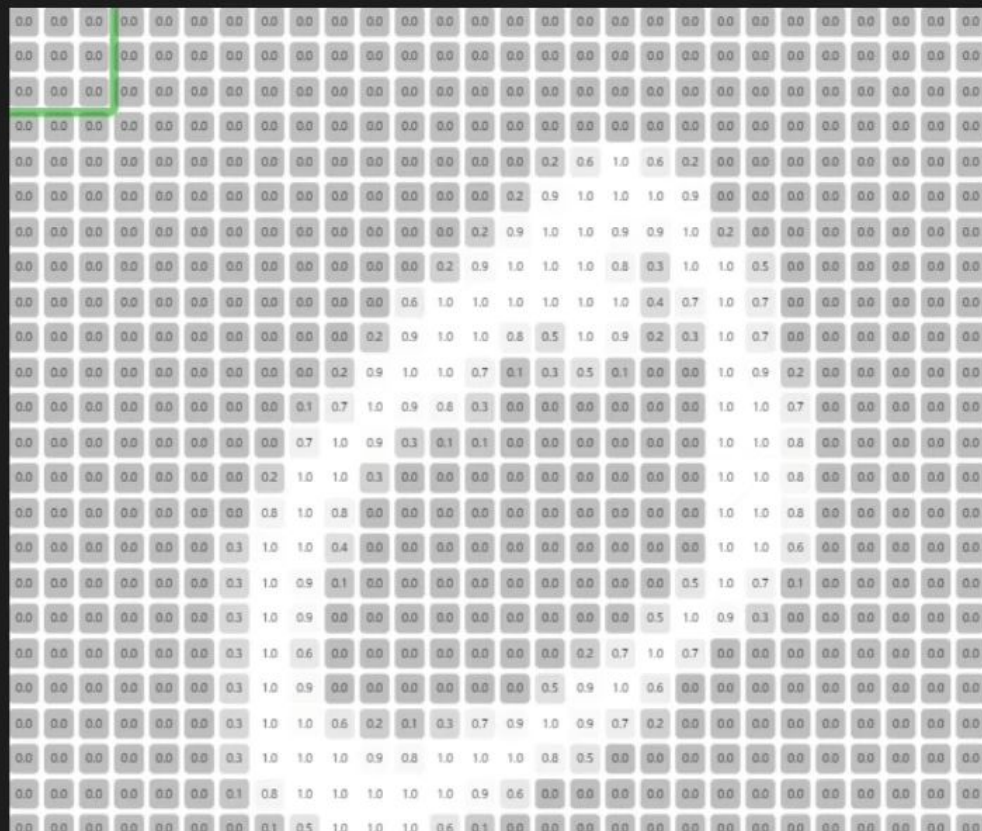
Grayscale images

- Single channel (intensity only)
- Shape: Height × Width × 1
- Example: Medical scans, handwritten digits (MNIST)

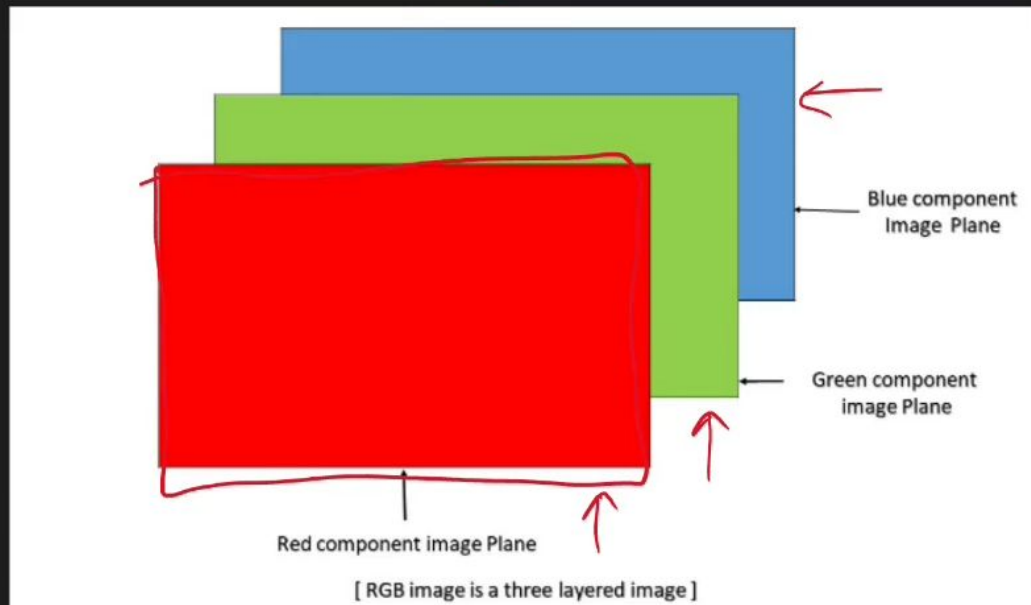RGB images

- Three color channels: Red, Green, Blue
- Shape: Height × Width × 3
- Example: Photos from cameras, web images (CIFAR-10, ImageNet

mnist → 0

( 28 × 28 )

R

Blue component Image Plane

Green component image Plane

Red component image Plane

[ RGB image is a three layered image ]

$228 \times 228 \times 3$

# Edge Detection in CNNs

In Convolutional Neural Networks (CNNs), edge detection emerges automatically in the early layers during training. The network uses small learnable filters (also called kernels) — typically 3×3 or 5×5 matrices of weights — that slide across the input image to produce feature maps. Initially, these filters contain random values, but through backpropagation and gradient descent, they adapt to detect patterns that help reduce the overall loss. Remarkably, without any human guidance, the first-layer filters often converge to detect fundamental visual primitives such as vertical edges, horizontal edges, diagonal lines, textures, and color contrasts. This self-organized feature learning is a core reason CNNs outperform traditional methods: instead of relying on handcrafted edge detectors (like Sobel or Canny), the network discovers the most useful filters for the task directly from data.

$255 \times 1 + 255 \times 1 + 2 \dots$

horizontal filter

edge detect

$-1 \times 0 + -1 \times 0 + -1 \times 0$

image

| 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |
| 255 | 255 | 255 | 255 | 255 | 255 |
| 255 | 255 | 255 | 255 | 255 | 255 |
| 255 | 255 | 255 | 255 | 255 | 255 |

$\ast$

| -1 | -1 | -1 |
|---|---|---|
| 0 | 0 | 0 |
| 1 | 1 | 1 |

$=$

| 0 | 0 | 0 | 0 |
|---|---|---|---|
| 255 | 255 | 255 | 255 |
| 255 | 255 | 255 | 255 |
| 0 | 0 | 0 | 0 |

feature map

$\uparrow$ 6×6

filter/kernel

Matrix 3×3

https://deeplizard.com/resource/pavq7noze2

Image * filter=feature Image

(m  x  m )  *(f  x  f) = (m-f+1)(m-f+1)

$(3\times3)\times3$

3channel

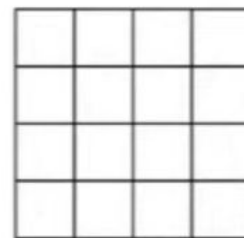

6 x 6 x 3          *          3 x 3 x 3          =          4 x 4
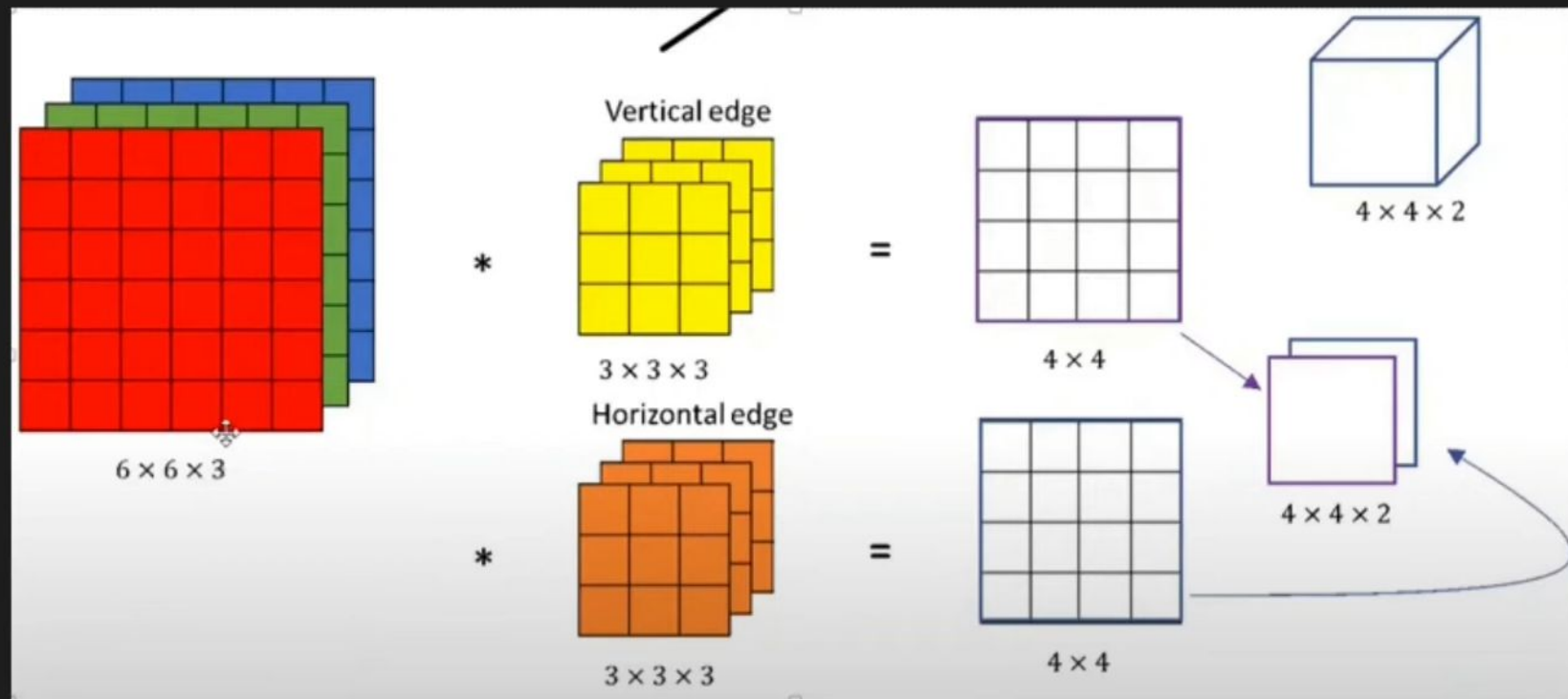
# Padding

Cause size of image less than original and numbers aka pixel value and border doesn't have that much importance than others

So we increase the size of margin

5*5 to=7*7 which give 5*5 feature map

2 types of padding

      Valid:no padding,

      Same padding to get same number of feature image as input

https://colab.research.google.com/drive/1HBMLctcBnhvV6Rj62Zc8eAXERQw54I2H?usp=sharing

Input matrix:

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 60 | 113 | 56 | 139 | 85 | 0 |
| 0 | 73 | 121 | 54 | 84 | 128 | 0 |
| 0 | 131 | 99 | 70 | 129 | 127 | 0 |
| 0 | 80 | 57 | 115 | 69 | 134 | 0 |
| 0 | 104 | 126 | 123 | 95 | 130 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Kernel

| 0 | -1 | 0 |
|---|----|---|
| -1 | 5 | -1 |
| 0 | -1 | 0 |

Output:

| 114 |  |  |  |  |
|-----|--|--|--|--|
|  |  |  |  |  |
|  |  |  |  |  |
|  |  |  |  |  |
|  |  |  |  |  |

# Stride

Is the moment space or speed/shift of filter over image

When stride decrease feature map size decrease

((n-f)/s)+1=output image

Stride are used when only high lvl features are required

# Pooling

**Pooling (typically max pooling) is used to:**

Reduce Spatial Dimensions

→ Shrinks the width and height of feature maps

Introduce Translation Invariance

→ Helps the network recognize features even if they shift slightly in the image

→ Example: A cat's eye is still detected whether it's 1 pixel left or right

Retain the Most Important Information

→ Max pooling keeps the strongest activation (e.g., the clearest edge response)

→ Discards less relevant details, acting as a form of controlled downsampling

Prevent Overfitting

# Types of pooling

https://medium.com/@abhishekjainindore24/pooling-and-their-types-in-cnn-4a4b8a7a4611

https://colab.research.google.com/drive/1F4F6Q9O-hPvCDeOWcqMUa5BuBOvuOBWc?usp=sharing

# Pooling

| 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |
| 255 | 255 | 255 | 255 | 255 | 255 |
| 255 | 255 | 255 | 255 | 255 | 255 |
| 255 | 255 | 255 | 255 | 255 | 255 |

**\***

| -1 | -1 | -1 |
|----|----|----|
| 0 | 0 | 0 |
| 1 | 1 | 1 |

**=**

relu

feature map (non-linear)

Max pooling
Min pooling
Avg pooling
L2 pooling
Global pooling

| 3 | 1 | 1 | 3 |
|---|---|---|---|
| 2 | 5 | 0 | 2 |

2) Translation invariance

Regularization

ANN — Weight $\downarrow$ bias

$\quad\Box\rightarrow$ loss function

we add a penalty term
in cost func

Cost $= L + ($ penalty $)$

$L_2$
$L_1$

when penalty value will come near

10. 2070

---

10. 2070

Intuition behind Regularization

learn rate

$$W_{now} = W_{old} \cdot \eta \left(\dfrac{\gamma L}{\partial w_d}\right) \quad -\text{partial derivative of loss fun}$$

$$= W_0 - \eta\left(\dfrac{\partial L}{\gamma w_i} + \lambda w_i\right)$$

so weight decrease

so model am simple, miss genl da

# 1) Dropout

we randomly drop node in layer

so architecture change

$2\%$ increase allunely

so depender on nodes decrease

s.o weigh dulubyllon increase

# Batch Normalizatios

It malle learning faster and stable

say Normalic (learning can so faster and slase)

hidon layer outpul → inpucs layor

we normalize every rcd individ

$$z_{11} = wx + b$$
$$y = sig(z_{11})$$

in batch Normal

2 method

① $z_{11} \to z_{11} \xrightarrow{Normalle} g(M_{11}^{N}) = a_{11}$

① $z_{11} \to g(z_{11}) \to a_{11} \to a_{n1}^{N}$

① $z_{11} \to g(z_{11}) \to a_{11} \to a_{n1}^{N}$

$$\frac{z_{11} - \mu w}{2} = z_{11}^{N}$$

$$z_{11} = \frac{w}{2} = z^{N}$$

$$M_{1} = \frac{1}{m} \sum_{i=1}^{m} z_{11}^{i} \quad m = batch \; size$$

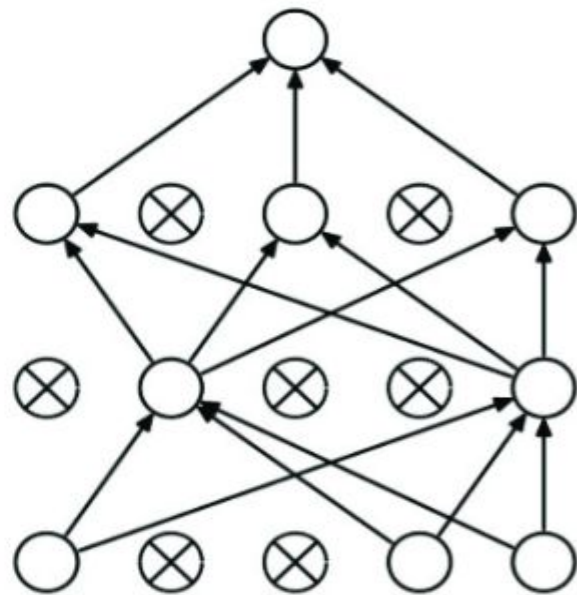$$\alpha = \sqrt{\frac{1}{m} \sum_{i=1}^{?} (z_{11}^{i} - MB)}$$

https://colab.research.google.com/drive/1PObj5KrXLDDmHjoJ1x0bVmxAFbif5s7q?usp=sharing

https://colab.research.google.com/drive/1KyMLdV1yB0qVdS-1huxKMN9xVKhrfxGL?usp=sharing

https://colab.research.google.com/drive/1473vOd0ICPbRW-co_Rm-_TBXgeajkJZ_?usp=sharing#scrollTo=xMDnrLgHHAcP

(a) Standard Neural Network (b) Neural Net with Dropout