

## What is a Convolution

- Operation that involves adding & multiplying
  - Input image \* Filter (kernel) = Output Image
- $\uparrow$   
 convolve  
 operation

- useful for performing edge detection

$$G_u = \begin{pmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{pmatrix} \quad G_y = \begin{pmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{pmatrix}$$

- convolution is like image modifier

- ↳ Feature transformation

- Different filters will result in different outputs

- ↳ Gaussian filter = Blurs image

Edge detection filter = Edge detection

## Convolution Operation

0	10	10	0
20	30	30	20
10	20	20	10
0	5	5	0

$$\begin{array}{r} 1 \ 0 \\ 0 \ 2 \\ \hline \end{array}$$

1)  $(0 \times 1) + (10 \times 0) + (20 \times 0) + (30 \times 2) = 60$   
 2)  $(10 \times 1) + (10 \times 0) + (30 \times 0) + (30 \times 2) = 70$   
 3)  $(10 \times 1) + (0 \times 0) + (30 \times 0) + (20 \times 2) = 50$

\* Convolution is commutative

$$A_j * W = W * A_j$$

Output matrix		
60	70	50
60	70	50
20	30	20

Convolution output

if input length =  $N$  then output length =  $(N - K) + 1$   
kernel length =  $K$

- \* Images are typically not square  $\leftarrow$  some neural nets use square for convenience.
- \* Filters are almost always square

## Convolution Equation

$$(A * w)_{ij} = \sum_{i'=0}^{K-1} \sum_{j'=0}^{K-1} A(i+i', j+j') w(i', j')$$

↑                      ↑  
kernel           kernel  
row                col

- \* Filters in CNNs are learned using gradient descent

1. Filters assigned with random values
2. Input convolve filters to produce feature maps
3. Loss calculated & compared with labels
4. Backprop to calculate gradient
5. Filters updated using gradient descent

Findling

- Pad input array with 0s to make convolved output same size as input
- "valid": output =  $N - K + 1$  (filter can only touch valid input)
- "same": output =  $N$  (output same size as input, done with padding)
- "full": output =  $N + K - 1$  (ensure non-zero output) (Not typically used)

## Dot Product & correlation

↖ Angle b/w a & b

- $a \cdot b = \sum_{i=1}^N a_i b_i = |a| |b| \cos \theta_{ab}$
- Also called cosine similarity / cosine distance
- $\cos(0) = 1$  ↗ same dir.
- $\cos(90) = 0$  ↗ orthogonal
- $\cos(180) = -1$  ↗ opp dir.

$$\cos \theta_{ab} = \frac{a \cdot b}{|a| |b|} : \text{how similar are } a \text{ \& } b.$$

max = 1

min = -1

0 : orthogonal

$$\text{magnitude of vector } |a| = \sqrt{\sum_{i=1}^N a_i^2}$$

$$\text{Pearson correlation} = \frac{\sum_{i=1}^N (a_i - \bar{a})(b_i - \bar{b})}{\sqrt{\sum_{i=1}^N (a_i - \bar{a})^2} \sqrt{\sum_{i=1}^N (b_i - \bar{b})^2}} \quad \leftarrow \text{very similar to dot product, just subtract mean.}$$

Dot product is like a correlation measure of 2 vectors.

↳ related to filter in CNN finding correlated features

## Matrix operation & convolution

$$a = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{bmatrix} \quad W = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \quad a * W = \begin{bmatrix} a_1 w_1 + a_2 w_2 \\ a_2 w_1 + a_3 w_2 \\ a_3 w_1 + a_4 w_2 \end{bmatrix} = \begin{pmatrix} w_1 & w_2 & 0 & 0 \\ 0 & w_1 & w_2 & 0 \\ 0 & 0 & w_1 & w_2 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{pmatrix}$$

↑  
Repeating filters  
simulate convolutions  
(Takes up space)

### Consider opposite case:

- using convolutions to replace matrix multiplications.
- main idea behind parameter / weight sharing

$$\text{Activation} = W^T x$$

↳ instead of full weight matrix,  
use same 2 weights over.

### Translation invariance

- convolution operation recognize patterns / features regardless of location within input space
- filter slides over input data, ensuring same filter applied across entire input

## CNN Input vs output

- For colour images, shape =  $H \times W \times 3$
- Kernel size =  $K \times K$
- output =  $H \times W$  for same size operation.

\* Discrepancy: Input = 3D, output = 2D. Need same shape to stack convolution layers.

Soln: For each feature output, stack it.

$$\begin{array}{l} \text{eg: Nose filter} = H \times W \\ \text{Eye filter} = H \times W \end{array} \quad \left. \vphantom{\begin{array}{l} \text{eg: Nose filter} = H \times W \\ \text{Eye filter} = H \times W \end{array}} \right\} \text{stacked} = H \times W \times 2$$

$$\therefore \text{output} = H \times W \times \text{No. of filters}$$

## CNN in Deep Neural Network

$$B = A * w, \quad A = H \times W \times C_1$$

↑

filter =  $C_1 \times K \times K \times C_2$  ← No. of features. per convolution layer.

$$H \times W \times C_2, \quad \text{Assuming same shape convolution}$$

↑

No. of features. in final convolution layer.

\* 3<sup>rd</sup> axis no longer colour but known as feature maps. or channels

$$\begin{array}{l} \text{conv layer: } \sigma(\underbrace{W * x}_{H \times W \times c} + b) \quad \text{size } c : \text{Different shape but} \\ \text{Dense layer: } \sigma(\underbrace{W^T \cdot x}_{\substack{\text{size } c \\ M}} + \underbrace{b}_{\substack{\text{size } c \\ M}}) \quad \text{rules of broadcasting in} \\ \text{Numpy allow such operation} \end{array}$$

## Cost savings (conv vs matrix multiplication)

### Convolution:

Input :  $32 \times 32 \times 3$

Kernel:  $3 \times 5 \times 5 \times 64 = 4800$  params.

output:  $28 \times 28 \times 64$   
( $32-5+1$ )

### Matrix Multiplication: (flatten)

Input =  $32 \times 32 \times 3 = 3072$

output =  $28 \times 28 \times 64 = 50176$

Weight:  $3072 \times 50176 = 154\,000\,000$  params.

↳ 32000 x more than conv operation.

Supplies from translation variance.

Need to learn every possible location