# Convolution Neural Networks

- Grid like data
  e.g. ⌐ 2D images.
       └ 1D time series

- Instead of fully connected layers
    (to all neurons)
  Focus on small regions ⌐→ Features
                 └ Using filters
                   (Sliding over the data)

- Helps focus on immediate neighbors
      instead of noise from all the features

- Same filter across all input so kind of same
                                      weight

- Creating automatic features
    (like maybe eyes, nose etc. in an image)

## CNN

e.g.  A sample input →  I =

| 10 | 11 | 3 | 8 |
| 2 | 16 | 12 | 19 |
| 7 | 1 | 18 | 11 |
| 13 | 4 | 3 | 14 |

$$I = H \times W \times C$$

H×W×C
4×4×1

Each value corresponds to a pixel (in case of an image) — height — width — channels (e.g. R,G,B)

○ Usually we do a bit of preprocessing

  - Normalize the matrix  - So values b/w 0 & 1 - So weights are not too funky & we don't run into problems training
  - Add a padding - Adding empty rows/columns to make sure the output doesn't shrink too much
          - Usually we add it, if we want our 'features' down the line to be of specific size
            Some networks add padding in every layer to preserve resolution

(1.1) Convolution

For simplicity we will walk through our sample without normalization OR padding for now

**What's Happening**

$$I_{H \times W \times C} \xrightarrow[\text{with } K_{(f_h \cdot f_w \times C)} \text{ filters}]{\text{Convolution Layer}} S_{H_{out}, W_{out}, C}$$

**Mathematically**

Nested loops for each position in filter (every row, col, channel) → Multiplication

$$S_{(i,j)} = \sum_{m=0}^{f_h - 1} \sum_{n=0}^{f_w - 1} \sum_{c=0}^{C-1} I(i \cdot s + m, j \cdot s + n, c) \cdot K(m,n,c)$$

Output of convolution layer — element on $i^{th}$ row & $j^{th}$ column

Input Matrix H×W×C

Stride (sliding window)

Filter ($f_h \cdot f_w$, C)
{ (m),(n),(c) under row, col, channel }

## Inside the Layer

Let    $K = $



$2 \times 2 \times 1$
$f_w \times f_h \times C$

$S = 1$
Stride
(Sliding window)

1 means window moves by 1 index

$p = 0$
No padding
(Downsampling)

### How it works



$p = 0$

9 times

$H_{out} = \frac{H - f_h + 2p}{s} + 1 = \frac{4-2}{1} + 1 = 3$

$W_{out} = \frac{W - f_h + 2p}{s} + 1 = \frac{4-2}{1} + 1 = 3$

$S(0,0) = \begin{bmatrix} 10 & 11 \\ 2 & 15 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} = 10 - 15 \Rightarrow -5$

$= \begin{bmatrix} I(0,0) & I(0,1) \\ I(1,0) & I(1,1) \end{bmatrix} \begin{bmatrix} K(0,0) & K(0,1) \\ K(1,0) & K(1,1) \end{bmatrix}$

$S(0,1) = \begin{bmatrix} 11 & 3 \\ 15 & 12 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} = 11 - 12 \Rightarrow -1$

$= \begin{bmatrix} I(0,1) & I(0,2) \\ I(1,1) & I(1,2) \end{bmatrix} \begin{bmatrix} K(0,0) & K(0,1) \\ K(1,0) & K(1,1) \end{bmatrix}$

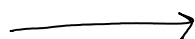$S(i,j) = \sum_m \sum_n I(i+m, j+n) \; K(m,n)$

## Overall

$I_{H \times W \times C}$    $\xrightarrow[\text{with } K_{(f_h, f_w \times C)} \text{ filters}]{\text{Convolution Layer}}$    $S_{H_{out}, W_{out}, C}$



$I_{(4 \times 4 \times 1)}$

$K(2,2,1), p=0, s=1$

$S(3,3,1)$

(1.2) ## Activation

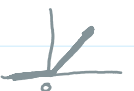This is pretty standard (happens in other NN)

Here we applied on the convolution output ($S$) above

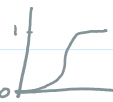There are various functions that can be used

Here we'll use ReLU, one of the popular activation function for sparsity ^(introducing some)

e.g.
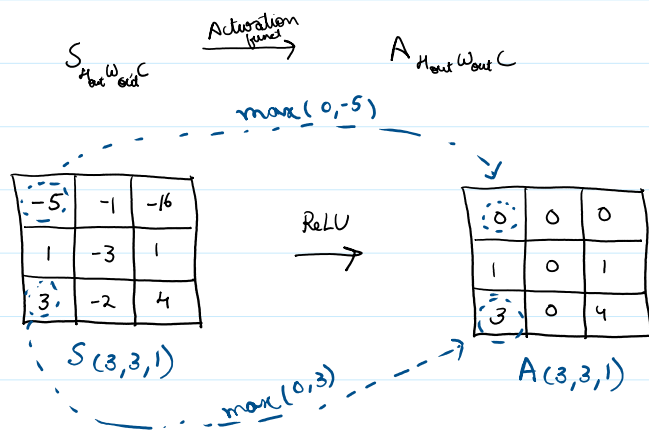
ReLU    $f(x) = \max(0, x)$

Sigmoid    $f(x) = \frac{1}{1 + e^{-x}}$

tanh    $f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

...we'll a use "ReLU", one of the popular activation function for sparsing

$$A(i,j) = ReLU(S(i,j))$$

- The output element in $i^{th}$ row & $j^{th}$ column
- $max(0,x)$
  - Only keep +ve values
  - Helps with gradient potential optimizations downstream
- Convolution output

| | | |
|---|---|---|
| tanh | $f(x) = \dfrac{e^x - e^{-x}}{e^x + e^{-x}}$ | |
| Softmax | $f(x) = \dfrac{e^x}{\sum_i e^{x_i}}$ | |
| LeakyReLU | $f(x) = \begin{cases} x & x > 0 \\ \alpha x & x < 0 \end{cases}$ | |
| Swish | $f(x) = \dfrac{x}{1 + e^{-x}}$ | |

$$S_{H_{out} W_{out} C} \xrightarrow{\text{Activation func}} A_{H_{out} W_{out} C}$$

$max(0,-5)$

| -5 | -1 | -16 |
|----|----|-----|
| 1 | -3 | 1 |
| 3 | -2 | 4 |

$S(3,3,1)$

$\xrightarrow{ReLU}$

| 0 | 0 | 0 |
|---|---|---|
| 1 | 0 | 1 |
| 3 | 0 | 4 |

$A(3,3,1)$

$max(0,3)$

② **Pooling**

Point is to reduce the size further (to control overfitting & improve computation)

Coalescing different features - a sliding window over the output of convolution + Activation

e.g      Max pooling  - max of the window  (keep only the strongest feature)

Avg. pooling - avg. of the window  ( combines features into a smaller more smoother matrix)

In this case, we'll use max pooling

$$P(i,j) = \max_{0 \le m, n < k} A(i.s+m, j.s+m)$$

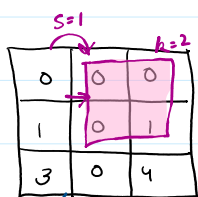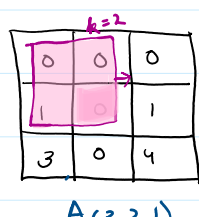- Output element in $i^{th}$ row & $j^{th}$ column
- Pooling window
- Activation output
- stride

e.g      $k = 2$        $s = 1$        Usually no padding
         Pooling        Stride
         window

$k=2$

| 0 | 0 | 0 |
|---|---|---|
| 1 | | 1 |
| 3 | 0 | 4 |

$A(3,3,1)$

$s=1$   $k=2$

| 0 | 0 | 0 |
|---|---|---|
| 1 | 0 | 1 |
| 3 | 0 | 4 |

4 times

$$H_{out} = \frac{H - k_h + 2p}{s} + 1 = \frac{3-2+0}{1} + 1 = 2$$

$$W_{out} = \frac{W - k_w + 2p}{s} + 1 = \frac{3-2+0}{1} + 1 = 2$$

$$A(3,3,1) \quad\quad A(3,3,1)$$

$$H_{out} = \frac{H - k_h + 2p}{s} + 1 = \frac{3-2+0}{1} + 1 = 2$$

$$W_{out} = \frac{W - k_w + 2p}{s} + 1 = \frac{3-2+0}{1} + 1 = 2$$

$$P \; 2\times 2$$

$$P(0,0) = \max\begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}_{2\times 2} = 1 \quad\quad P(0,1) = \max\begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}_{2\times 2} = 1$$

$$P(0,0) = \max\begin{pmatrix} A_{00} & A_{01} \\ A_{10} & A_{11} \end{pmatrix} \quad\quad P(0,1) = \max\begin{pmatrix} A_{01} & A_{02} \\ A_{11} & A_{12} \end{pmatrix}$$

$$P(i,j) = \max_{0 < m, n < k} A(i+m, j+n)$$

$$A_{(H_{out}, W_{out}, C)} \xrightarrow{\text{Pooling Layer}} P_{(H_{out}', W_{out}', C)}$$



$$A(3,3,1) \quad\quad k=2, \; s=1 \quad\quad P(2,2)$$

## ③ Flatten

Just as it sounds

Takes this condensed feature matrix (after convolution, activation, pooling) & puts into 1d vector

$$P_{(H, W, C)} \xrightarrow{\text{Flattening Layer}} F_{(1, \; H \cdot W \cdot C)}$$



$$[1, \; 1, \; 3, \; 4]$$

## ④ Fully Connected

Again a standard layer (last layer) followed by an activation function for final output

Fully Connected

Again a standard layer (last layer) followed by an activation function for final output

$y = \omega x + b$

e.g a softmax function

Learnt params

$\Rightarrow$ Final Output $= \sigma(\omega x + b)$

$$y_j = \sigma\left(\sum_i^N \omega_{ij} x_i + b_j\right)$$

In our case

$$\text{Final output} = \sigma\left(\omega \begin{bmatrix} 1 \\ 1 \\ 3 \\ 4 \end{bmatrix} + b\right)$$

If our fully connected layer has lets say 2 neurons, $\omega$ will be $4 \times 2$
e.g for a binary classification & b " " $2 \times 1$
here output of each neuron
describes the probab. of that class

$$y = \sigma\left(\begin{bmatrix} \omega_{01} & \omega_{11} \\ \vdots & \vdots \\ \omega_{03} & \omega_{13} \end{bmatrix}_{2\times 4} \begin{bmatrix} 1 \\ 1 \\ 3 \\ 4 \end{bmatrix}_{4\times 1} + \begin{bmatrix} b_1 & b_2 \end{bmatrix}_{2\times 1}\right)$$

$$= \sigma\left(\begin{bmatrix} \omega_1 x_1 + b_1 \\ \omega_2 x_2 + b_2 \end{bmatrix}_{2\times 1}\right) = \begin{bmatrix} p_1 \\ p_2 \end{bmatrix} \rightarrow \text{Neuron1} \rightarrow \text{Neuron2}$$