

# Convolutional Neural Networks:

## History:

- Mark I Perceptron

- ↳ Connected to a camera that produced 400 pixel image

$$f(x) = \begin{cases} 1, & \text{if } w \cdot x + b > 0 \\ 0 & \text{otherwise} \end{cases}$$

## Update rule:

$$w_i(t+1) = w_i(t) + \alpha (d_j - y_j(t)) x_{j,i}$$

- Adaline & Madaline:

- ↳ Started to stack linear layers into multilayer perceptron
  - ↳ No backprop

- Rumelhart:

- ↳ Backprop was introduced

$$\frac{\partial \mathcal{L}}{\partial w_{ji}} = \frac{\partial \mathcal{L}}{\partial o_j} \times \frac{\partial o_j}{\partial w_{ji}}$$

- Hinton & Salakhutdinov:

- ↳ Showed how to train nn
  - ↳ Not really like today's nn
  - ↳ Required careful initialization

- Hubel & Wiesel:

- ↳ Found that cells had hierarchical organization

- Retinal ganglion cell:

- ↳ Respond to circular regions

- Simple cells:

- ↳ Responsive to oriented edges

- Complex cells:

- ↳ Responsive to light orientation & movement

- Hypercomplex cells:

- ↳ Responsive to movement with an end point (corners)

- Neurocognition:

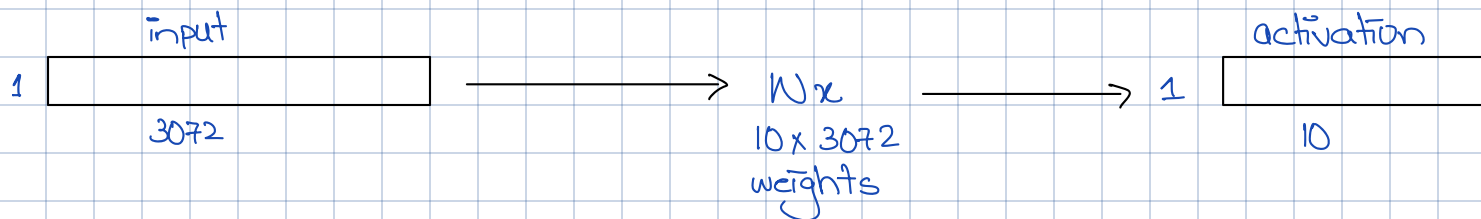
- Simple cells: modifiable parameters

- Complex cells: perform pooling

- Gradient based learning applied to document recognition:
  - ↳ Used backprop & gradients to train CNNs
  - ↳ NOT used on complex data
  - ↳ AlexNet solved that

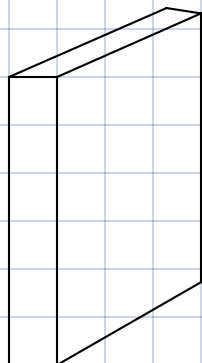
## Fully connected layer:

↳  $32 \times 32 \times 3$  image  $\longrightarrow$  stretch to  $3072 \times 1$

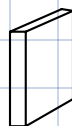


## Convolution layer:

- $32 \times 32 \times 3 \longrightarrow$  preserve spatial structure



$5 \times 5 \times 3$  filter

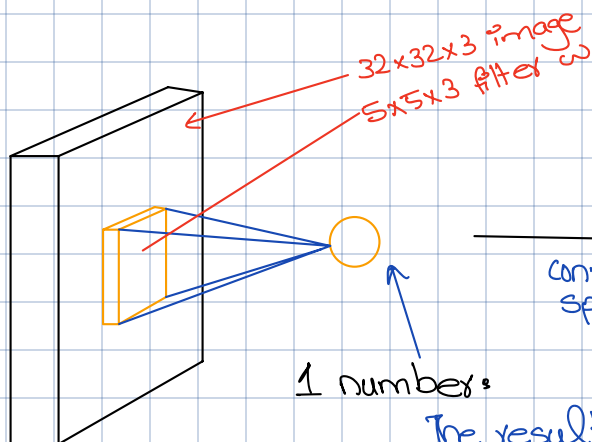


Convolve the filter with the image  
i.e. slide over the image spatially,  
computing dot products

activation maps

$28 \times 28 \times 1$

arrangement of  
edges, corners,  
lines in the  
image

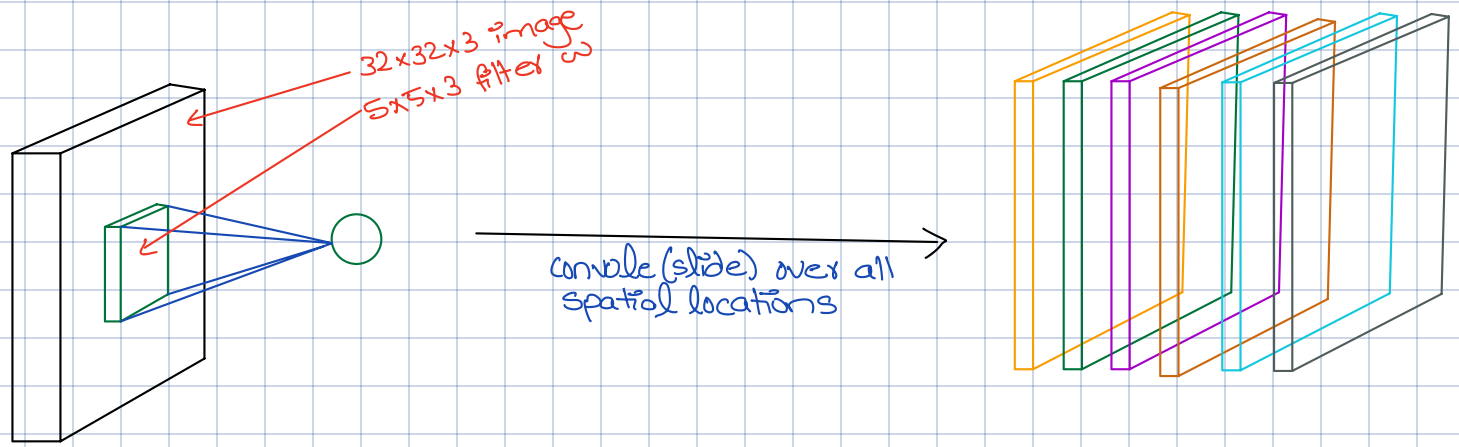


convolve (slide) over all  
spatial locations

1 number:

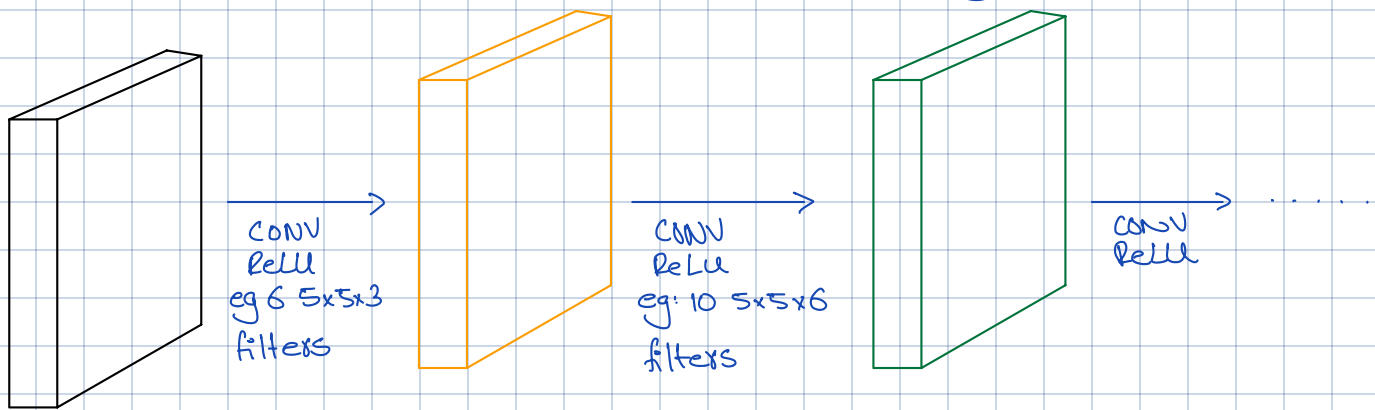
The result of taking a dot product b/w filter  
& small chunk of a image.  
(i.e.  $5 \times 5 \times 3 = 75$  dot product +  $b$ )  
 $w^T x + b$

Consider a 6 filters  $5 \times 5 \times 3$ , because we need to work with multiple filters, as each filter is looking for a template

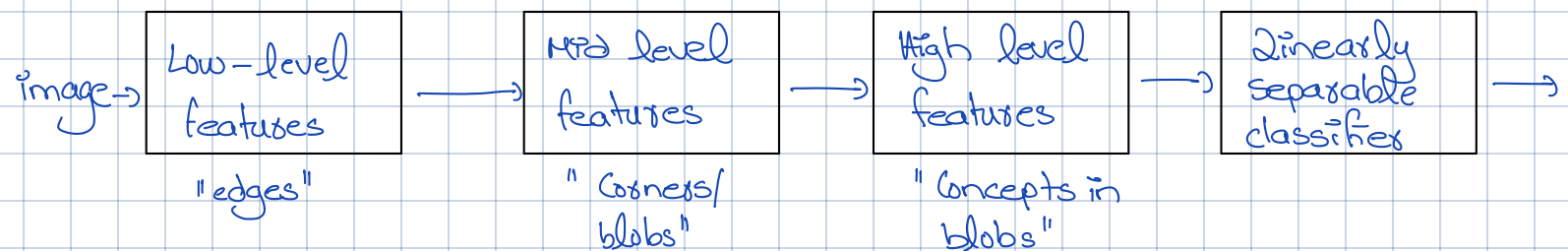


• We can stack these to get a new image of size 28x28x6!

ConvNet is a sequence of these convolution layers



In ConvNet, you end up learning this hierarching of filters,



• We call layer convolution because it is related to convolution of two signals

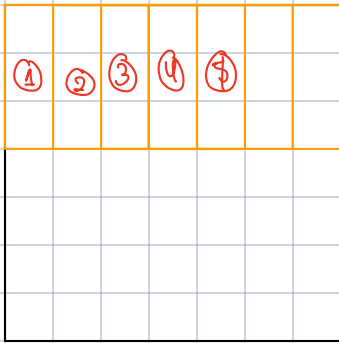
$$f[x_1, x_2] \cdot g[x, y] = \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f[n_1, n_2] \cdot g[x-n_1, y-n_2]$$

A closer look at spatial dimensions:

- How we get  $28 \times 28 \times 1$  activation map

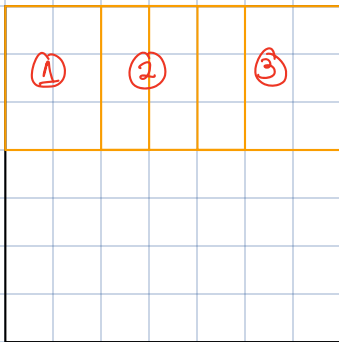
Consider  $7 \times 7$  input,  $3 \times 3$  filter & stride 1

↳ interval of slide



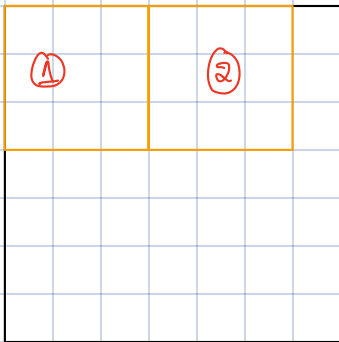
⇒  $5 \times 5$  output

Consider  $7 \times 7$  input,  $3 \times 3$  filter & stride 2



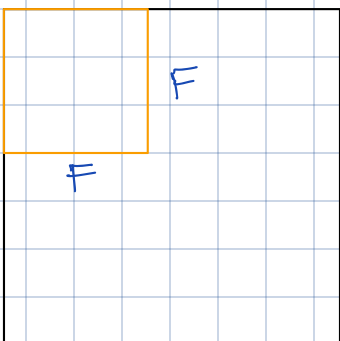
⇒  $3 \times 3$  output!

Consider  $7 \times 7$  input,  $3 \times 3$  filter & stride 3



⇒ Doesn't fit  
So, we don't do  
this

Formula for checking how strides will work,



Output size:

$$\frac{(N-F)}{\text{stride}} + 1$$

eg:  $N=7, F=3$

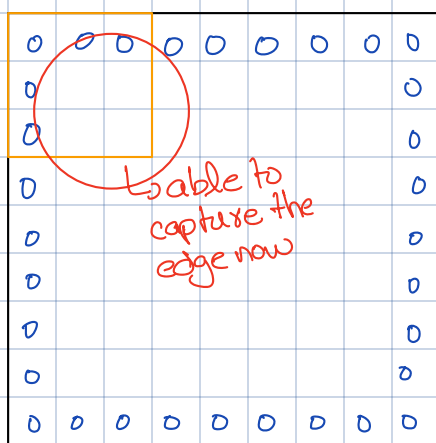
stride 1:  $(7-3) / 1 + 1 = 5$

stride 2:  $(7-3) / 2 + 1 = 3$

stride 3:  $(7-3) / 3 + 1 = 2.33$

↓  
oops!

In practice, it is common to zero pad the corners so that filters can fully overlap the image



Now,

$$\frac{N - F + 2 \times P}{\text{stride}} + 1$$

padding ↑

If we have multiple convolution layers without zero padding with 5x5 filters, the volume will shrink very quickly.

Example,

Input: 32

Filters: 10 of 5x5

Output size:  $\frac{32 + 2 \times 2 - 5}{1} + 1 = 32$

An activation map is a 28x28 sheet of neuron outputs:

1. Each is connected to a small region in the input
2. All of them share parameters

$32 \times 32 \times 10$

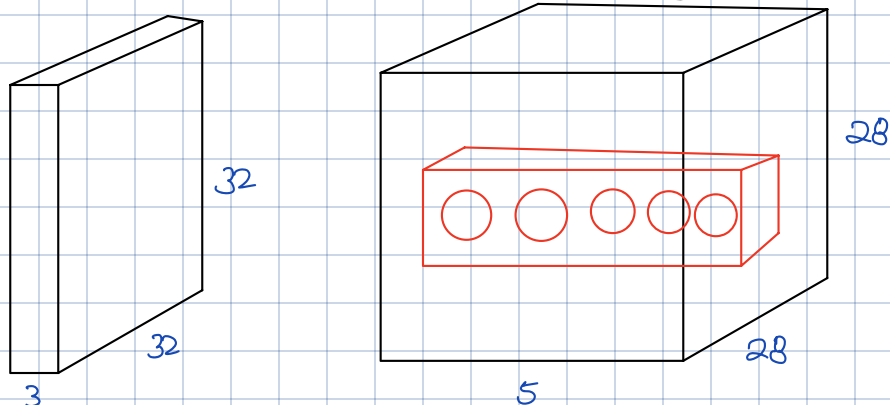
Number of params:  $\overset{\text{Filter}}{\uparrow} 5 \times 5 \times 3 + \overset{b}{\uparrow} 1 = 76$   
 $76 \times 10 = 760$

Hyperparameters:

- 1) Number of filters  $k$  ↑ Powers of 2
- 2) Their spatial extent  $F$
- 3) The strides  $S$
- 4) The amount of zero padding  $P$

kernel = filters = receptive field

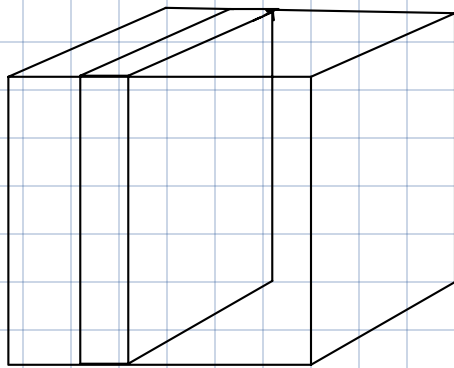
The brain/Neuron view of CONV layer:



The five neurons are all looking at the same region but for different things

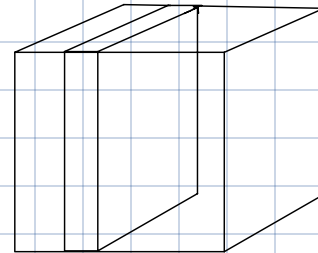
## Pooling layer:

- ↳ Makes the representation smaller and more manageable
- ↳ Operates over each activation map independently



224x224x64

Pool →



112x112x64

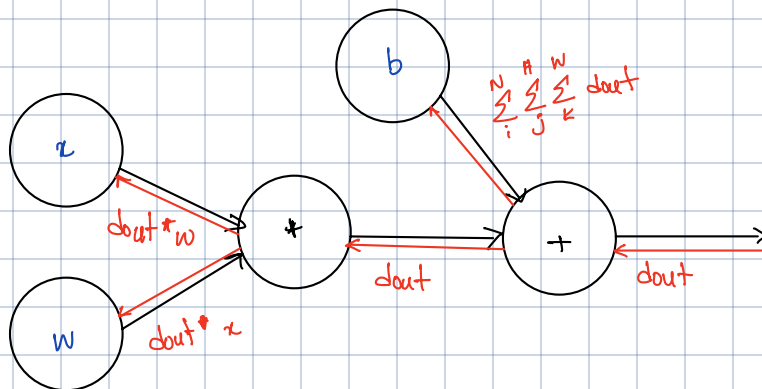
## Max pooling:

|   |   |   |   |
|---|---|---|---|
| 1 | 1 | 2 | 4 |
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 2 | 3 | 4 |

max pool with 2x2 filters  
stride 2 →

|   |   |
|---|---|
| 6 | 8 |
| 3 | 4 |

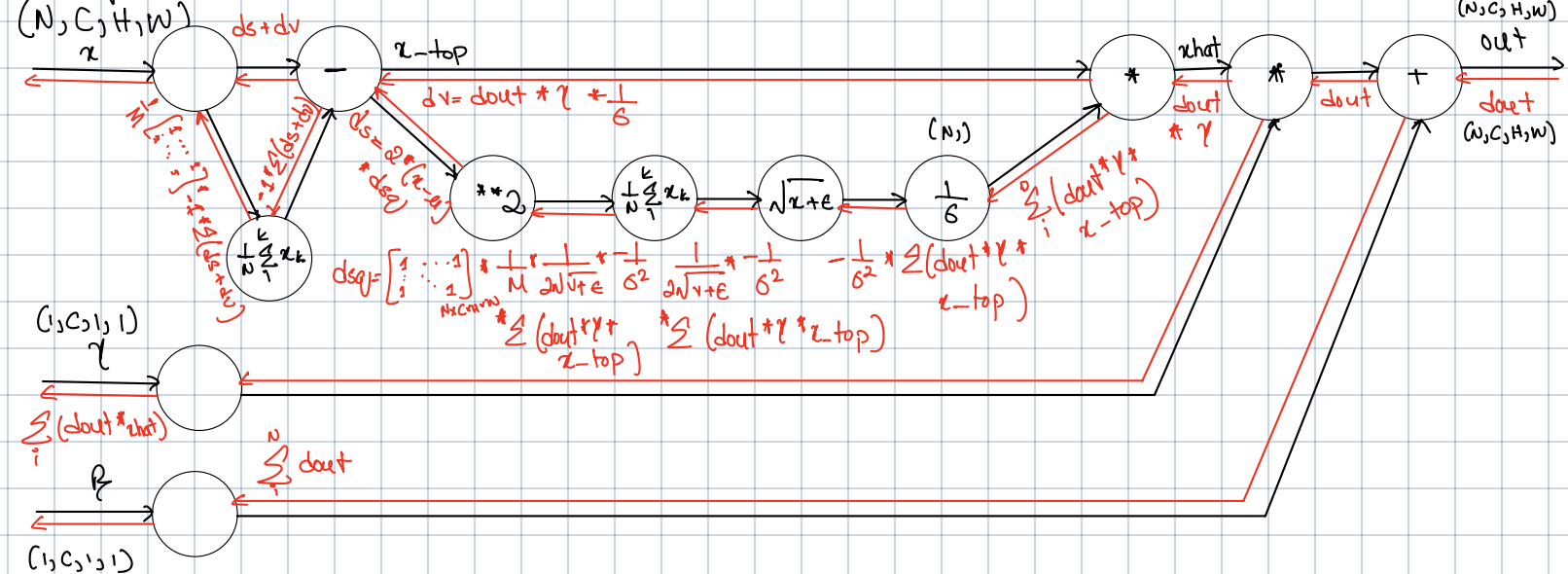
## Assignment helper:



Group norm:

$$U = N^* H^* W$$

$(N, C, H, W)$



$$dx = ds + dv - \frac{1}{N} \begin{bmatrix} 1 & \dots & 1 \\ \vdots & & \vdots \\ 1 & \dots & 1 \end{bmatrix} * \sum (ds + dv)$$