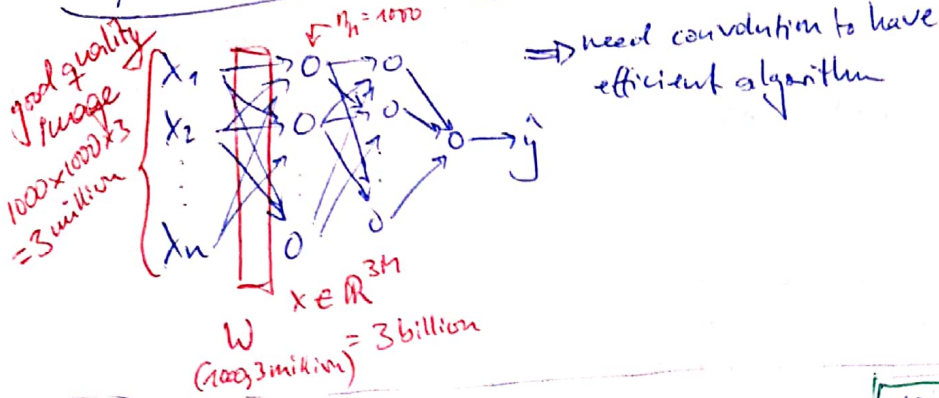


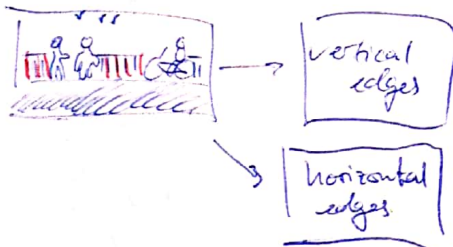
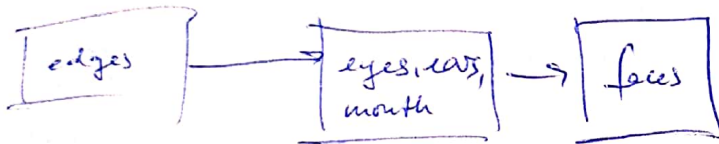
# DL spec // course IV // week 1

## Convolutional Neural Networks

### Computer Vision Example

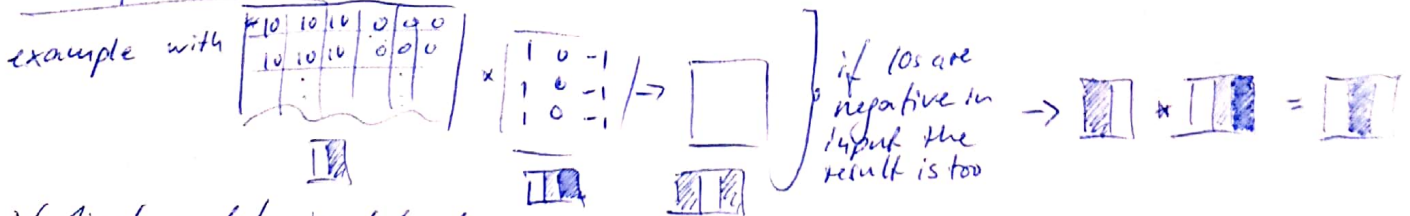


### Edge Detection Example

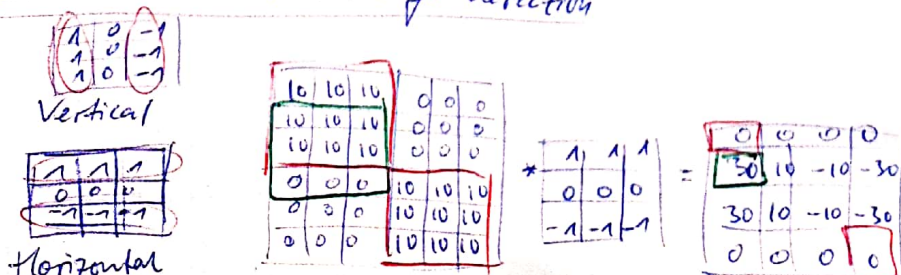


a vertical edge is a  $3 \times 3$  region with bright pixels on left and dark pixels on the right

### More edge detection



### Vertical and horizontal edge detection



# Learning to detect edges

$$\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$$

$$\rightarrow \begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$$

**Sobel filter**

→ focus on center row  
→ more robust

$$\begin{bmatrix} 3 & 0 & -3 \\ 10 & 0 & -10 \\ 3 & 0 & -3 \end{bmatrix}$$

**Scharr filter**

just for vertical edges!

filter numbers don't need to be handpicked, can also be learned! using back-prop.

$$\begin{matrix} * & \begin{matrix} w_1 & w_2 & w_3 \\ w_4 & w_5 & w_6 \\ w_7 & w_8 & w_9 \end{matrix} \end{matrix}$$

↳ can also learn different edges than only vertical and horizontal like 45°, 70°, 73°...

→ Neural Networks learning low-level features like edges...

## Padding

$$6 \times 6 * 3 \times 3 = 4 \times 4 \rightarrow \text{because there are only } 4 \times 4 \text{ possible positions to fit the filter in the input matrix}$$

$$n \times n * f \times f = (n-f+1) \times (n-f+1)$$

2 down → each time with convolution the image shrinks sides  
→ pixel at edge is only processed once while a pixel in the center of input image is "used" more often

to prevent, you can pad the image

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

add a border of pixel around the edges with 0-value  
6x6 → 8x8

$$\text{now } 8 \times 8 * 3 \times 3 = 6 \times 6$$

→ managed to preserve initial input size

$P = \text{padding} = 1$  Output:  $n+2p-f+1 \times n+2p-f+1$   
 $p=2$  would be another border around input image

## Valid and Same convolutions

→ no padding

$$\text{"Valid": } n \times n \text{ image} * f \times f \text{ filter} = n-f+1 \times n-f+1 \text{ output}$$

"Same": Pad so that output size is the same as the input size

$$n+2p-f+1 \times n+2p-f+1$$

$$n+2p-f+1 = n \Rightarrow p = \frac{f-1}{2}$$

so for 3x3 filter  $p = \frac{3-1}{2} = 1$

or for 5x5 filter  $p = \frac{5-1}{2} = 2$

↳ f is usually odd

↳ if f was even we'd have unsymmetric padding

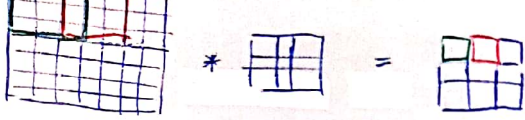
↳ odd dimension filter has a central position → benefit for position of filter



# Strided Convolutions

## Summary of convolutions

jump of 2 pixels with stride 2



$$7 \times 7 * 3 \times 3 \rightarrow 3 \times 3$$

stride 2  $\rightarrow$  take 2 steps

$$n \times n * f \times f$$

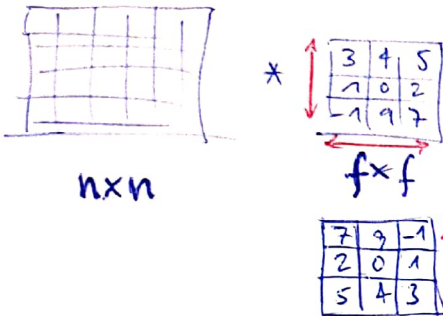
padding stride S  
S=2

$$\left\lfloor \frac{n+2p-f}{S} + 1 \right\rfloor \times \left\lfloor \frac{n+2p-f}{S} + 1 \right\rfloor$$

$\lfloor \cdot \rfloor$  is symbol for "floor", so rounding down to next integer

## Technical note on cross-correlation vs. convolution

Convolution in math-textbooks:



flip horizontally and vertically

actually more interesting for signal processing so that "associativity" is valid

$$(A * B) * C = A * (B * C)$$

$\rightarrow$  doesn't matter for DL

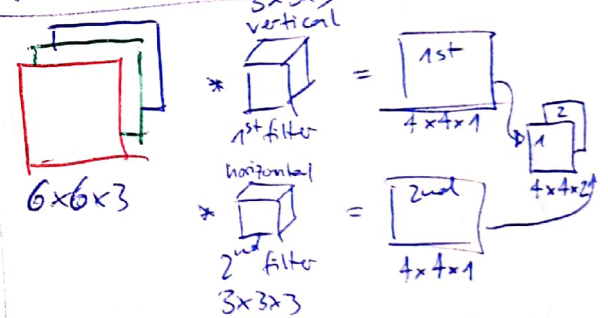
## Convolutions over volume (3D)

convs on RGB image



6x6x3  
height width #channels  
RGB  
channels have to match!  
3x3x3 filter  
height width #channels

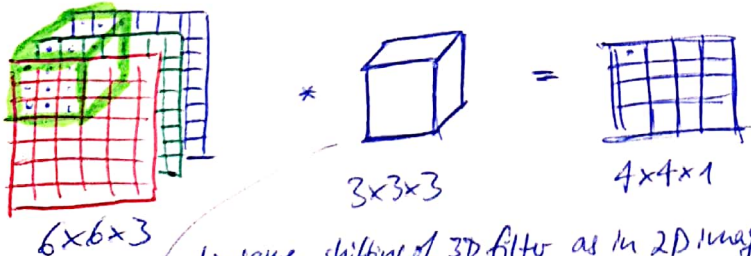
## Multiple filters



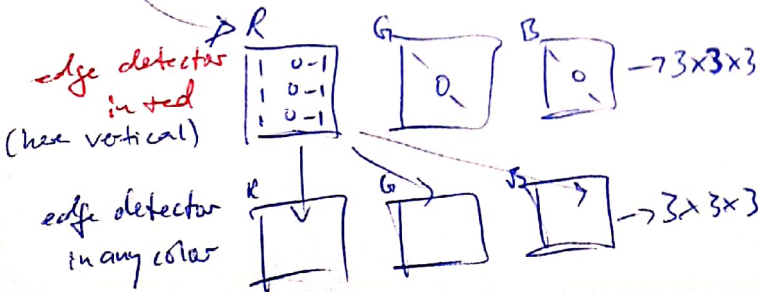
## Summary

$$n \times n \times n_c * f \times f \times n_c \rightarrow \frac{n-f+1 \times n-f+1 \times n_c}{4 \times 4 \times 2}$$

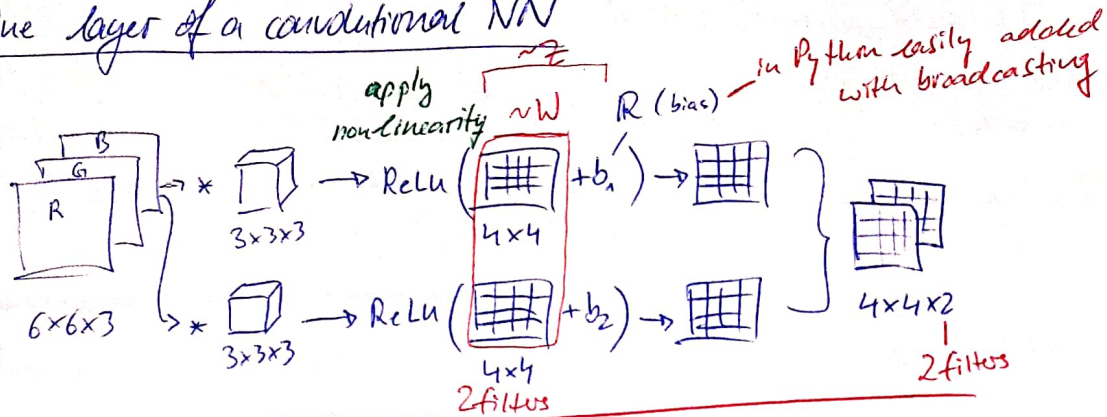
# filters



$\rightarrow$  same shifting of 3D filter as in 2D images...



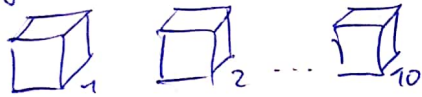
## One layer of a convolutional NN



## One layer of a CNN

### Number of parameters in one layer

If you have 10 filters with 3x3x3... how many params in one layer?



3x3x3

= 27 + bias

→ 28 params

→ 280 parameters

(independent of input image size)

### Summary of notation

If layer  $l$  is a convolution layer.

$f^{[l]}$  = filter size

$p^{[l]}$  = padding

$s^{[l]}$  = stride

$n_c^{[l]}$  = number of filters

Each filter is:  $f^{[l]} \times f^{[l]} \times n_c^{[l-1]}$

Activations:  $a^{[l]} \rightarrow n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

vectorized:  $A^{[l]} \rightarrow m \times n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

Weights:  $f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$

bias:  $n_c^{[l]} - (1, 1, 1, n_c^{[l]})$   $L$  # filters in layer  $l$

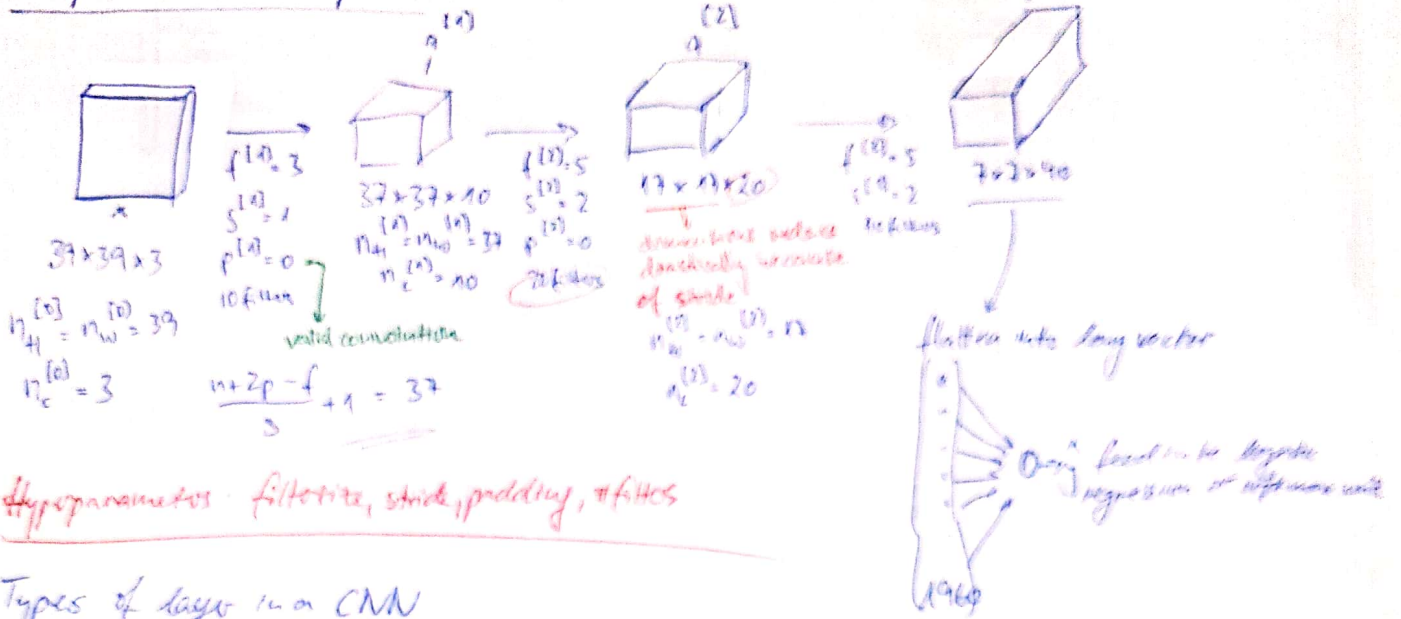
Input:  $n_h^{[l-1]} \times n_w^{[l-1]} \times n_c^{[l-1]}$   
height width  
of previous layer

Output:  $n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

$$n_{h/w}^{[l]} = \left\lfloor \frac{n_{h/w}^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$$



## Simple CNN example



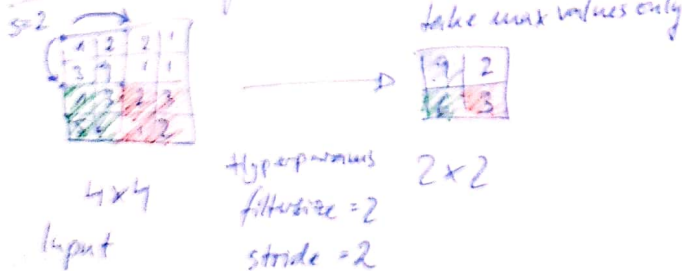
Hyperparameters: filter size, stride, padding, # filters

## Types of layers in a CNN

- Convolution (CONV)
- Pooling (POOL)
- Fully connected (FC)

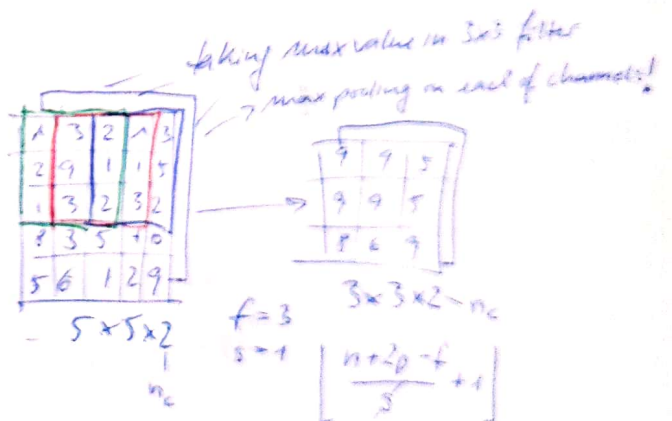
## Pooling layers

### Max Pooling

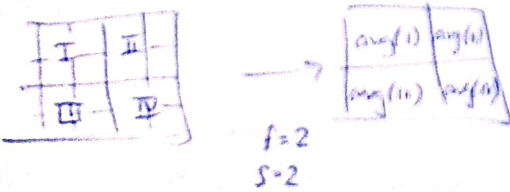


Hyperparameters  
 filter size = 2  
 stride = 2

→ no params to learn! Fix computation! (with GD...)



### Average Pooling (can be used deeper in a CNN)



### Summary of pooling

Hyperparameters *no params to learn*

$f$  filter size  $f=2, s=2$   
 $s$  stride  $f=3, s=2$

Max or average pooling (usually no padding)

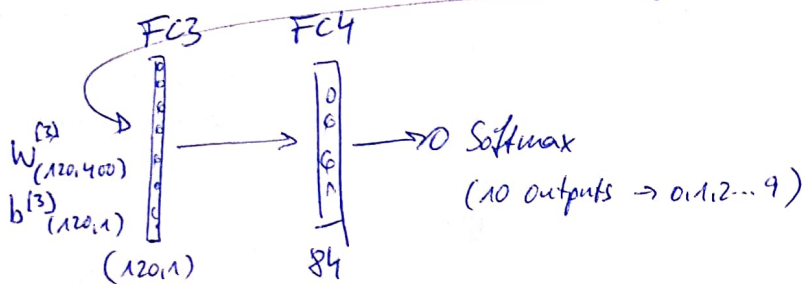
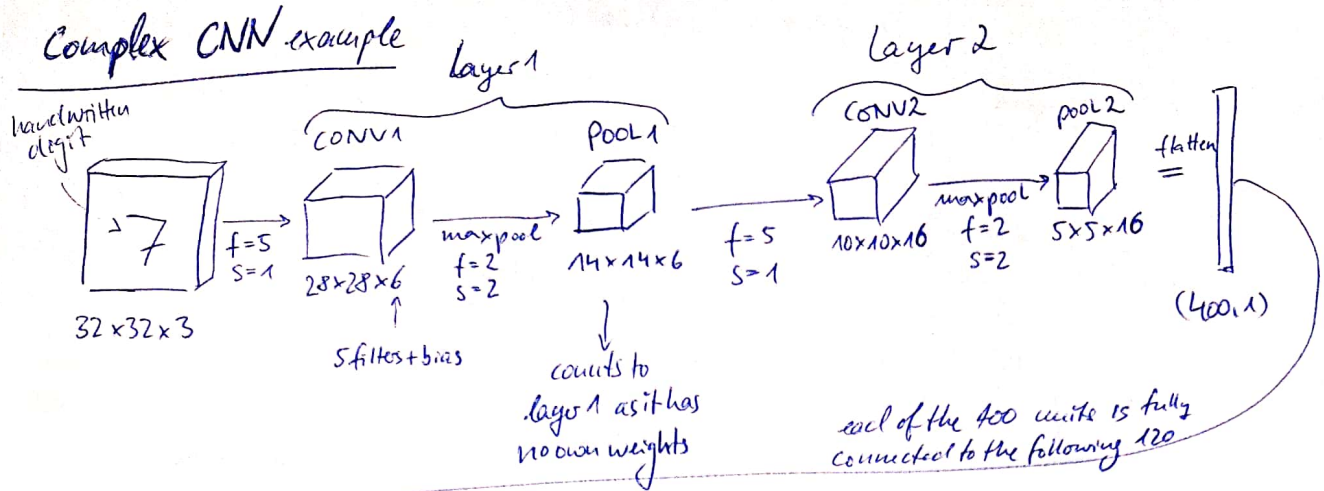
$n_H \times n_W \times n_C$

$$\left\lfloor \frac{n_H - f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n_W - f}{s} + 1 \right\rfloor \times n_C$$

because pooling applies to each channel independently

# of input channels is equal to output channels

# Complex CNN example



look for hyperparameters in literature

$n_H, n_W \downarrow$   
 $n_C \uparrow$

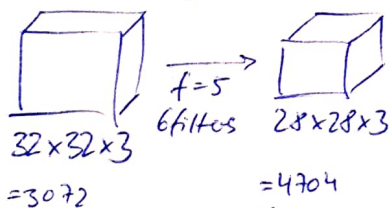
usually CONV-POOL-CONV-POOL-FC-FC-FC-SOFTMAX

- POOL layers have no parameters
- CONV layers' params increase with depth
- FC have huge #params!
- Activation size will decrease with depth but should not drop quickly → performance issues

## Why Convolutions?

2 main advantages over fully connected layers

- parameter sharing
- sparsity of connections



as NN  
the connections would lead to a weight matrix of 14 million parameters!  
→ can easily become too large

$$\text{filters } 5 \times 5 \rightarrow 25 + 5 \text{ bias} = 26 \text{ params}$$

$$(5 \times 5 \times 3 + 1) = 456$$

Parameter Sharing: A feature detector (vertical edges) that is useful in one part of the image is probably useful in another part of the image  
→ apply same filter to different regions

Sparsity of connections: In each layer, each output value depends only on a small number of inputs

CNN → translation invariance

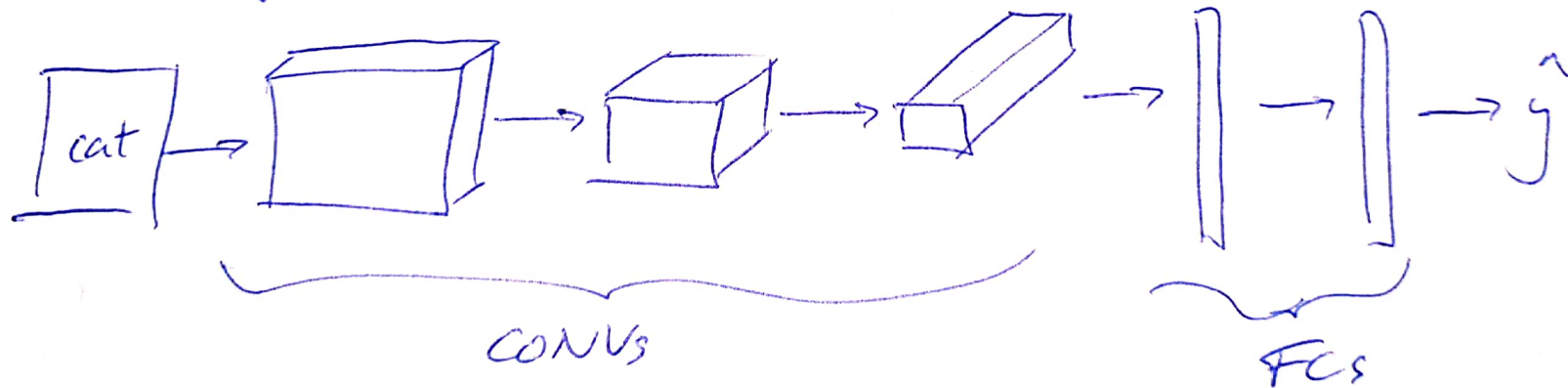
- ↳ a picture of a cat shifted a few pixels to the right still looks like a cat
- ↳ the fact that we apply the same filter adds to robustness

DL spec // course IV // week 1

Putting it together

Training set  $(x^{(m)}, y^{(m)})$

$W, b$  randomly initialize



$$\text{Cost } J = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Use gradient descent (or RMS, adam) to optimize parameters to reduce  $J$