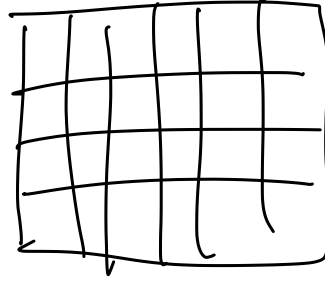
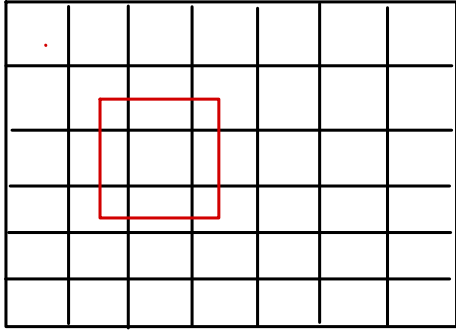


Recall



## Typical workflow

- Take in a set of  $n$  images

- Do this  $K$  times  
(new channels)
- For each image, apply same convolutional filter (either slightly reduces size, or retains size)
  - Apply ReLU (or other activation) at each pixel (same size)
  - Max pool (reduce size)

These steps are iterated to form a "deep" CNN.

Eventually, with all the pooling, the images become small (e.g.  $5 \times 5$  or  $4 \times 4$ ), at which point we

"Flatten" them:

1	2	3
4	5	6
7	8	9



1  
2  
3  
4  
5  
6  
7  
8  
9

and then do a final one (or a few) regular fully connected layers. If the response is quantitative, the final output is numeric; if it is qualitative, the final output is a vector of the ~~the~~ classes, and a softmax (or similar) function is applied.

Remark Images are usually color, not grayscale.

This is often encoded as 3 channels (AKA #s)

$(R, G, B) = (\text{red}, \text{green}, \text{blue})$  for each pixel, with values in  $\{0, 1, 2, \dots, 255\}$ .

$(0, 0, 0) = \text{black}$ ,  $(255, 255, 255) = \text{white}$

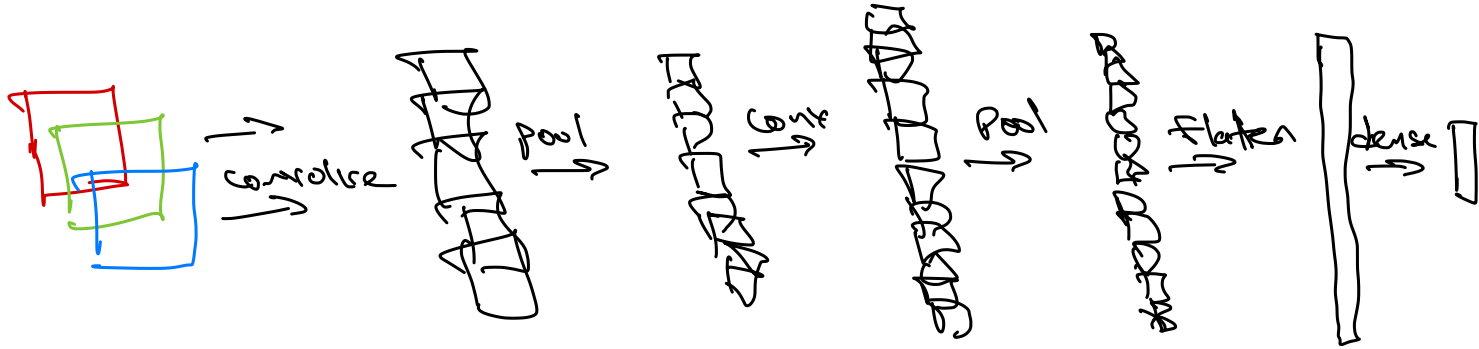
$(0, 0, 255) = \text{Blue}$   $(5, 71, 42) = \text{evergreen}$

$(0, 128, 128) = \text{teal}$  etc...

So each original image has 3 "channels", i.e. there are really 3x original images at the input stage.

A CNN applies a different filter to each channel, and adds them together at the 1st stage (~~before~~ ReLU or pooling).

## A typical representation



**Ex** Classification (10 classes) of  $32 \times 32$  color images.  
Define + count parameters of a "deep" CNN:

- Input arrays of dimension  $(32, 32, 3)$   
Spatial image  $\rightarrow$  # of channels = (R, G, B)

- Do a  $3 \times 3$  convolutional layer 32 times

output of stage:  $(30, 30, 32)$

   
image size      # of channels


Count parameters:  $(3 \text{ input channels}) \times (9 \text{ params/conv})$   
 $\times (32 \text{ convs}) + (32 \text{ biases}) = 896 \text{ parameters}$

- ReLU previous step, no change
- Max-pool, output dimension  $(15, 15, 32)$
- Do another  $3 \times 3$  conv layer, 64 times,  
 output dim:  $(13, 13, 64)$

   
image size      64 channels

$$\text{Count param: } (32 \text{ input channels}) \times (9 \text{ params/conv}) \\ \times (64 \text{ convs}) + (64 \text{ biases}) = 18496$$

- ReLU, no change

- Max pool, output dim:  $(6, 6, 64)$   


- Final convolutional layer of  $3 \times 3$ , 64 times

output dim:  $(4, 4, 64)$

$$\text{param count: } (64 \text{ input channels}) \times (9 \text{ params/conv}) \\ \times (64 \text{ output channel}) + (64 \text{ biases}) = 36,928$$

- ReLU, no change

- Flatten into vector of length  $4 \times 4 \times 64 = 1024$

- Dense layer with 64 nodes,

output: vector of length 64

$$\text{param count: } (1024 \text{ inputs}) \times (64 \text{ params}) + (64 \text{ biases}) \\ = 65600$$

- ReLU, no change

- Final stage dense layer w/ 10 outputs

dim: 10

$$\text{param: } (64 \text{ hidden nodes}) \times (10 \text{ outputs}) + \\ (10 \text{ biases}) = 650$$

Total # of params: 122,570

- Softmax prior stage, classifier is class with highest prob.
- 

Note For "natural images" data augmentation is often used

- rotation
- stretching
- zooming
- cropping
- flipping