

CNN Architecture

- modern CNN originate from LeNet
- Typical CNN have 2 stages
- stage 1: Conv + Pool layers. \leftarrow Feature transformer / finder
- stage 2: Fully connected layers \leftarrow non-linear image classifier.
 \leftarrow must flatten 3D \Rightarrow 1D.

Pooling

- for downsampling / upsampling
|
output smaller image from larger image
eg. $512 \times 512 \Rightarrow 128 \times 128$. (downsample by 3).
- Max vs Avg pooling
|
from filter, choose max / avg values.

1	2
5	6

\leftarrow max = 6
avg = 3.5

1	2	3	4		max		avg
5	6	7	8		6	8	3.5 5.5
16	15	14	13	\rightarrow	16	14	or 13.5 11.5
12	11	10	9				

Purpose:

- 1) less data to process
- 2) translation invariance (location does not matter)

Found	Not Found	$\xrightarrow{\text{max}}$	Found.
Not Found	Not Found		

Pool size & step

- Typically square.
- stride = steps to slide
- if stride < size, overlap occurs. (not common)

Conv-Pool Design

conv-pool	\rightarrow	conv-pool	\rightarrow	conv-pool
<u>simple lines & edges</u>		<u>more complex features (eyes, lips)</u>		<u>entire face, beard etc.</u>

- After each conv-pool, image shrinks
- filter sizes stays the same, typically 3×3 , 5×5 , 7×7
- \therefore After each conv-pool layer, filter search for increasingly large patterns.
- As image shrinks, No. of feature maps \uparrow
 \leftarrow spatial info \downarrow
No. features \uparrow

CNN Hyperparams

- Filter size & stride
 - No. feature maps
- Pool size & stride } \exists existing conventions to follow.
- 1) follow conv-pool design
2) Increase No. feature maps $32 \rightarrow 64 \rightarrow 128$.

Global max pooling 2D

- To allow for input image of different shape to be flattened.
eg: Input 1: $32 \times 32 \times C_2$
Input 2: $64 \times 64 \times C_2$ } Flatten does not work.
- Global max pooling always get output of $1 \times 1 \times C_2$
↳ Takes max value of each feature map.

* Image cannot be too small eg 2×2 .

Final activation f₂:

Regression = none

Binary = sigmoid

multi = softmax