

Convolution:

$$g(x) = f(x) * h(x) = \int_{-\infty}^{\infty} f(\tau) h(x - \tau) d\tau.$$

'Convolutional Neural

Nets'

- also called CNN convnets.
- CNN chop input image into little patches and then process each patch independently and identically.

Convolutional layer: convolution,

- way of combining two functions (image + filter) to get a new filter
- includes flipping a filter before sliding across the image

Cross-correlation:

- similar to convolution
- filter not flipped

Separable filters in CNN:

◦ instead of using a full square filter (3×3) we can break it into 2 smaller filters

→ vertical filter (3×1) goes up-down ($k \times 1$)

→ Horizontal filter (1×3) goes left-right ($1 \times k$)

◦ fast and cheap

◦ less memory

Downsampling and Upsampling layers:

↓
smaller
shrink

↓
bigger
growing

→ Strided convolution

- skips pixels eg. every 2nd.

→ pooling (Max pooling).

- biggest value chosen from small patches

→ Dilated convolution

(Transposed convolution)

- inserts \neq zeros b/w pixels

◦ apply convolution

→ Simple

Nonlinear filtering layers - Pooling layers:

- reduce size of the data while keeping the important features
- removes extra details (small, fast changing ones) high frequency info

Invariance:

- no matter where the object is in the image, the network identifies it.
- CNN does not care about the position

Equivariance:

- if object moves in the image the output also shifts.

* pooling shrinks the image, keeps important info
it helps the network become less sensitive to invariance.

I know it moved \rightarrow I still know what it is

Normalization layers.

adjusting the values so they are consistent,
no too big or too small values.

Local normalization:

- looks at the small neighbourhood of values and adjusts each value based on the values around it.

Identifying vulnerabilities:

- try finding weak spots.
- identify weaknesses in the learned representation and in our training set.

modulation: change frequency content of an image so it fools the network.

Feature Maps,

- snapshot of what network sees.
- focus on each layer
- image transformed into more meaningful info.

1) image becomes smaller (downsampling)

2) no. of feature maps increased

image size $32 \times 32 \rightarrow 16 \times 16 \rightarrow 8 \times 8$

maps $3 \rightarrow 32 \rightarrow 64 \rightarrow 128$.

→ Early layers

first layer

detect → lines, edges, ~~colors~~ colors.

→ Deeper layers.

detect abstract things < shapes, eyes/wheels, objects.

Sliding filter.

- remove nonlinear parts (like ReLU).
- stack multiple convolution layer to make one big convolution.
- nonlinearities affect each pixel individually
- same process applied to every patch

Treating patches as independent.

- 1)
 - easier to solve
 - faster
 - flexible size
 - locality of visual content.

- 2)
 - Translation invariance in Images.

objects do not change when they move.

- convolution works the same everywhere.
does not care where the obj is

Encoder

Decoder

- shrinks the image (downsampling)
- keep only imp feature.
- uses less memory + faster

- rebuilds the full image
(upsampling).
- recover details

Information Bottleneck.

- tiny connection narrow space b/w encoder and decoder
- only keep essential info
- forced to learn abstract concept.

U-nets: Encoder + Decoder + skip connections.

- encoder when compresses an image some details are lost.
- skip connections pass some of that lost information directly from early layer to later layer.

'Image Pyramids'

Introduction.

→ Translation invariant linear filters.

- works the same way no matter where the obj is in the image

→ Scale invariant processing.

- want our system to recognize obj even if it appear bigger or smaller.

→ To handle different sizes, we need a way to process images at multiple scales (zoomed in and zoomed out).

→ Template matching works if the object size is fixed.

Image pyramids:

- set of copies of same image but at different sizes/resolutions
- original image at the bottom (full size)
- as we go up the pyramid, each level smaller and blurrier.
- image is analyzed at different zoom levels

Gaussian Pyramid:

1) Low pass filtering:

blur the images to remove small details

2) Subsampling

shrink the image to half its size (subsampling by 2).

3) Repeat the process.

previous image \rightarrow blurring \rightarrow shrinking

Laplacian Pyramid:

◦ shows what's different b/w two levels.

◦ focuses on details (edges and textures).

1) Difference b/w two levels.

take small blur image resize to the size of larger one and subtract it.

shows what's lost b/w the two images/levels.

2) Repeat for each level.

Do the same process for every pair and form a pyramid

of detailed images called laplacian pyramid.

3) Low-pass residual.

the smallest gaussian image is kept at the top of the pyramid

Image Blending:

◦ Blending in gradual and natural way can be done by using laplacian and gaussian pyramids

1) Build laplacian pyramid for both images.

Image A \rightarrow build laplacian pyramid \rightarrow capture details of A.

Image B \rightarrow " \rightarrow "

3) Create a mask and its Gaussian Pyramid.

- make a mask and decide where to blend

eg left half = 1 (Image A) Right half = 0 (Image B)

- create gaussian pyramid of the mask.

4) Blend the laplacian pyramids

- + Image A Laplacian * mask.

- + Image B Laplacian * 1 - mask.

Add.

$$l_{out} = l^A + m + l^B * (1-m)$$

5) Reconstruct the final image

take blended laplacian pyramid and collapse it.

final image found

Sampling and Aliasing

Intro to Sampling.

- converting something ~~smooth~~ smooth and continuous into a set of points or pixels that a computer can understand.

Two goals of sampling.

→ More samples = more details ↑ memory + time to process

→ Fewer samples = faster and lighter lose details.

Aliasing.

- sample slowly and computer misunderstands the original signal fast wave (18Hz) but sample it slowly and looks like a slow wave, this is wrong and is called aliasing error).

Sampling Theorem: Nyquist Theorem.

- states how fast we need to sample

sample at twice the highest frequency

sampling frequency $W_s = 2\pi/T_s \geq W_{\text{max}}$

sampling period $T_s \leq T_{\min}/2$, $T_{\min} = 2\pi/W_{\max}$

Anti-Aliasing:

- shrink an image ($412 \times 512 \rightarrow 26 \times 32$) and make it big again
- image looks blurry details lost
- anti-aliasing smoothes out these problems

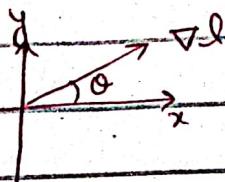
'Image Derivatives'

- Image derivative tells us how intensity changes in an image.
- detects edges and transitions.

Derivative - discrete approximation

$$\nabla l \approx \begin{pmatrix} \frac{\partial l}{\partial x} \\ \frac{\partial l}{\partial y} \end{pmatrix}$$

$$\nabla l = \left(\frac{\partial l}{\partial x}, \frac{\partial l}{\partial y} \right). \quad \frac{\partial l}{\partial x} = \text{change in intensity horizontally}$$
$$= l(x,y) - l(x-1,y)$$



$$\frac{\partial l}{\partial y} = \text{change in intensity vertically}$$
$$= l(x,y) - l(x,y-1)$$

* derivative can't be used on pixels, difference used to approximate it

Derivative as a convolution:

The derivative operator is linear and translation invariant.
Therefore, it can be written as convolution.

Filter is slides across the image to compute derivative

$$I(x,y) \xrightarrow{\left[\begin{array}{c} \frac{\partial}{\partial x}, \frac{\partial}{\partial y} \end{array} \right]} (I_x(x,y), I_y(x,y))$$

Derivative - discrete approximation

- 2 filters to compute 1D derivative

$$1) d_0 = [1, -1]$$

$$\Delta I = I[n] - I[n-1]$$

subtracts previous value from current value.

fast but not smooth

$$2) d_1 = [1, 0, -1]/2$$

called central difference

$$\Delta I = \frac{I[n+1] - I[n-1]}{2}$$

neighbor on left + right, i.e.

balanced + smoother