

UO-CNNs pt 2Story so far

- scanning, redistr., hierarchical.
- deformations in input \rightarrow max pooling
- 2D scan \rightarrow convnet
- 1D scan \rightarrow time delay NN.

(*) CNNs have neurophysiological roots (inspiration)

- key paper is Hubel & Wiesel paper (1959)

(*) 2 levels of processing:-

- i) S-cells (not robust)
- ii) C-cells. (fires if majority of S-cells fire)

S-cells \rightarrow C-cells

S-C layer (simple-complex module)

- build more complex patterns by composing early neural responses

(*) Neocognition, Fukushima (1980)

- computationalise a remedied version of H-W.
- ~~one~~ off issue \rightarrow position invariance
- grandmother cells fire regardless of position

(*) Neocognition (v2)

- hierarchical, successive layers of S-cells and C-cells.
- S-cells respond to signal in previous layer
- C-cells confirm S-cells response.

(*) Neocognition

S cells \rightarrow S planes (all cells within S-plane - ident resp.)

- each plane is a detector for one specific pattern
(some pattern in a different region of input)



- every one of cells has identical response
(within S-plane)

C-cells - behind

C-cells \rightarrow C-planes

- C-planes smaller than S-planes

- Fukushima's original paper

- S-planes examine elliptical regions of input. } can be used successively
- C-planes examine elliptical - II - of S-planes

- one C-plane per S-plane.

(*) Neocognition (full diagram)

- more complex patterns
- composition/refinement of patterns. (e.g. 2 eyes + nose)

(A1) - Review logic of pattern detection

(#) neocognition

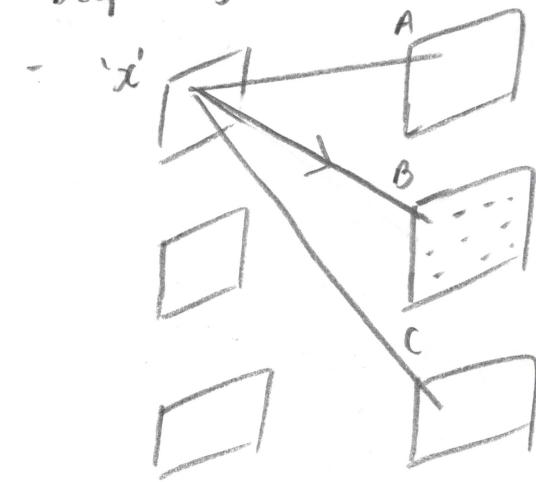
- S-cells } \rightarrow complex orig. form.

- C-cells

- But note BR draws similarity of S-cells with RELU
C-cells with MAX

(*) Neocognition (depth, rec. field)

- deeper layer \rightarrow larger receptive field.



- C-plane contributes to output every S-place in subs. layer
- Which of A, B, C has highest response
- e.g. B.
- use Hebbian rule
- if B largest
- update weights; but due to fixed response; all weights in other parts of place updated

(A2) - rev.
- slides

?)

(*) Learning in neocognition

- unsupervised flavor. \rightarrow clustering
- C-planes do not need to be estimated

(*) Neocognition - fine-tune

- present lots of examples \rightarrow semantic labels 'automatically learned'

(*) Adding supervision

- temporal corr.: Nomura, Atkes, Movius '88
- TDNN: Log, Waibel et al. 1989, 1990.

- Addition of supervision
- final output

(*) (tiny) NN

- CNN: Yann LeCun.

(*) Supervising neocognition

- Add extra layer after final C-layer \rightarrow c1 output.
- BR: fully feedforward MLP with sliced param
- All S-cells within S-plane share per weights.
- Backprop \rightarrow estimate S-cell weights in every plane of every layer
- ~~C-cells not updated.~~

(*) Scanning vs multiple filters / superv.

i) square recept. fields

ii) S-cells RF $\rightarrow K_s \times K_s$

iii) C-cells RF $\rightarrow L_c \times L_c$

(*) Supervising NC

- include term to include exp.

$$u_{S,L,N}(i,j) = \sigma \left(\sum_{p=1}^{K_s} \sum_{R=1}^{K_s} \sum_{t=1}^m w_{S,L,N}(p,R,t) u_{C,L-1,F}(i+t-1, j+R-1) \right)$$

$$u_{C,L,F} = \max_{i \in [1, L], j \in [1, F]} (u_{S,L,N}(i,j))$$

(*) weights $w_{s,i,n}$ are position independent (between C and S -plane)?
and do not change based on which portion of image you are analysing
(x) weights - corresponding region of image
 w

- imposing weights on top of image
- compute pixel by pixel pairwise product
- summing overall planes
- yielding combined response at one particular point
- pass through activation (sigmoid') \rightarrow s-cell output
- C-cell looks at rectangular patch of S-cells; we apply max

(B3) - tie this intuition together with the equations

- (*) instead of having many copies of neuron looking at every region of the image; it is possible instead to have one neuron scanning
(*) Scanning with a neuron \rightarrow 'convolution'

(*) CNN

- AT&T lab \rightarrow Turing Award
- revolutionary \rightarrow USPS.

(*) Story so far

- good summary
- visual cortex \rightarrow neocognition \rightarrow external sup. \rightarrow CNNs

(*) CNNs

- started as a computational model of a biological phenomenon ①
- (ii) A cybernetic development \rightarrow instantiation of a mechanism

(*) General architecture of CNN

- successive convolutional, downsampling layers

(*) convolutional - neurons that scan inputs for patterns.

- (*) Downsampling → perform max operations on groups of outputs from convolutional layers.
- may occur in sequence; but typically alternate.
- (*) followed by MLP with one or more layers (final)
- (*) doesn't usually make sense for downsample to be followed by downsample.
 - ↳ one large downsample op.
- (*) downsamples (max) reduce size of input (map).
- (*) MLP operating on flattened version of the outputs of the final layer of neurons.

(*) G.A of CNN - Learnability

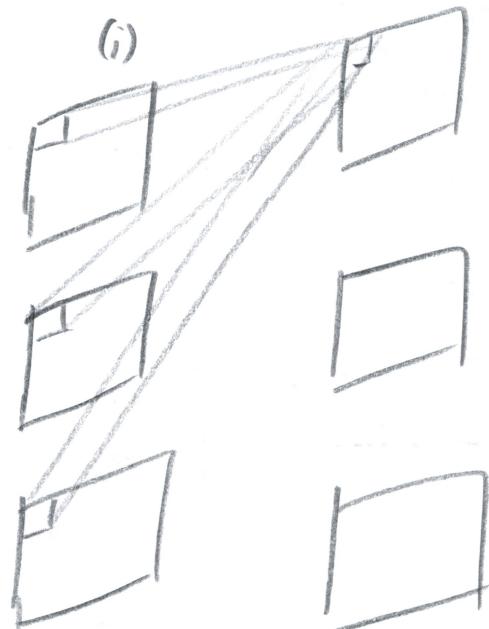
- convolution layers and MLP are learnable
i.e. params estimated from training for target class. task
- downsampling layers
- no learnable parameters.

(*) Convolutional layer (CL)

- CL - series of maps
- analogous to NC S-planes
- create/act. maps.

(*) CL

- within each map - pair of operations
- i) Affine combination
 - A linear map obtained by convolution over maps in the previous layer.
 - each affine map assoc. with a learnable filter



ii) Activation that operates
on convolution

⑧ - see diagram
(excellent)

- CL \rightarrow affine + activation (similar to MLP)

(*) If you want to compute affine combo at any one location, it is done
so over all maps from previous layer ⑨

(*) All maps in previous layer contribute to each convolution

(*) What is convolution?

- consider conv. of single map (brown patch)

⑨ (*) Examine how brown patch contributes to one point
in the corresponding map

(*) Weights are binary in pedagogical example
(1/0)

- green \rightarrow patch from previous layer.

- orange \rightarrow place filter on patch

input map ; filter ; convolved feature

↓ ↓
values values

(*) Take affine combo of filter values with input map values for
a particular region

(affine combo)

(*) See slides for a visual demo ⑩ ⑪.

(*) stride > 1 possible \rightarrow smaller output map.

(*) What really happens.

(*) See slides \rightarrow small gray patch is a result of computation
over multiple input maps simultaneously.
intuitive.

- 3×3 convolved feature \rightarrow 9 weights, 1 bias.

⑫ - match visual display to formulae - crucial

(*) 1st map, 2nd map

- weights for 1st, 2nd map very different.

② - weights, activ. formulae

(*) A different view

- stacked arrangement of planes

- generalise to cubes / cuboids

③ : key formula to anchor intuition visually:-

$$z(s, i, j) = \sum_p \sum_{k=1}^K \sum_{l=1}^L w(s, p, k, l) y(p, i+l-1, j+k-1) + b(s)$$

(I)

(II)

(I) - pixel by pixel product; then summing over \rightarrow (has specific weights)

(II) - for every patch (over planes)

- computation of convolutional map at any location sums the convolutional outputs at all planes

(*) filter (cube) - has different values at every single location
- weights at one plane not identical to weights on all other planes

(*) CNN - vector not.

pseudocode

- below corrs.

- weight $w(l, j)$ is a 3D $D_{l-1} \times K_l \times K_l$ tensor (ass. sq RFs)

- prod. in blue \rightarrow tensor mat prod. with scalar output

$y(0)$ - image

layers op. on vector at (x, y)

FOR $l = 1:L$

FOR $j = 1:D_l$

FOR $x = 1:W_{l-1}-K_l+1$

FOR $y = 1:H_{l-1}-K_l+1$

segment = $y(l-1), :, x:x+K_l-1, y:y+K_l-1)$ - 3D tensor

$z(l, j, x, y) = w(l, j). \text{segment}$

tensor mat prod.

$$y(l, j, x, y) = \text{activation}(z(l, j, x, y))$$

$$y = \text{softmax}((Y(l, :, :, :)))$$

(*) engineering consid.

- stride of conv. can be greater than 1 pixel.
- size of output of convol. op depends on implementation
- no need for it to be same size as input

(**) essentially examine constraints of problem

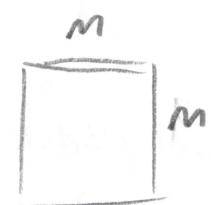
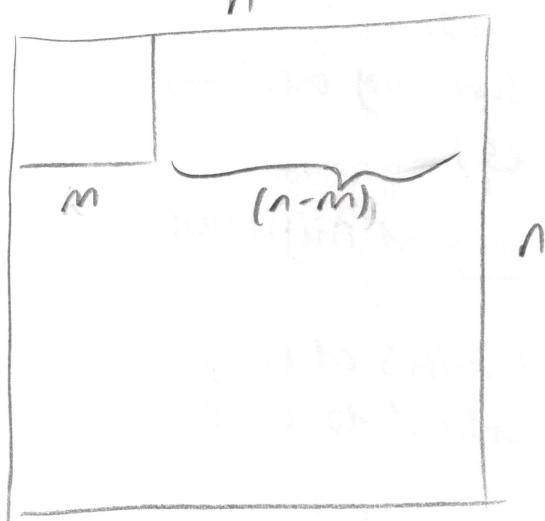
↳ does filter 'fall off edge'?

(*) convolution size

(*) In general for an $(n \times n)$ image / input map

$(m \times m)$ filter

output will be:-



$$\text{output size: } 1 + \left\lfloor \frac{(n-m)}{s} \right\rfloor - \text{no. of pixels } m \text{ output in each dim.}$$

→ reduction in output size

- even if $s=1$

need to maintain size and keep track of proportions

(**) often; this is not acceptable

⑤ Hence allow filter to 'fall off the edge' → systematic way.

- assume stride $s=1$, how many rows are 'lost' when going from left to right?
 $N - (N - M + 1) = (M - 1)$ rows
- Hence $\frac{M-1}{2}$ - pad this many columns of 0s to left and right

(**) review calc.

- output size = input size - regardless of stride

- after 0 padding

(*) padding (A)

(*) review skipped slides

(*) Activation

HW2 → ReLU activ.

(*) other component

(*) review

- downsampling/pooling (follows convolution)

- maxpooling → selects largest from a pool of elements

- pooling performed by 'scanning input' (?)

(*) looks at groups of pixels from output map (?)

(*) one downsampling corresp. to every output map

- one C for every S

- C is looking at one S; not across all Ss

(*) Max strides by more than one pixel typically. (to account for jitter).

(*) Pooling-size of output

(*) See slides!)

○ interactive. (similar principle)

- NxN picture compassed by a P x P pooling filter with stride 0
results in an output map of side:- $\left\lceil \frac{N-P}{0} \right\rceil + 1$.

- typically do not zero-pad

(*) All-mean pooling

(*) Purpose of max is to use jitter and variance

- if petal location moves a little, should still detect petal (?)
(shift-invar and jitter distinction?)

- compute mean instead.

- or softmax
- or p-norm.

(*) other options

- pooling may be a learned filter

(?) quick.
 review

(*) All convolutional net

- downampling - achieved with conv. layer with stride > 1
- softmax instead of max \rightarrow standard convolution

(BR): - these always \rightarrow convolutions all the way through

(*) fully convolutional network (no pooling)

pseudocode

(AQ) - review.

(*) setting everything together (engineering / actual task)

- image scaling \rightarrow GPU contingent / informed.

(*) CNN -

- parameters to choose: K_1, l, S

1. No. of filters $\rightarrow K_1$

2. size of filters $\rightarrow l \times l \times 3 + \text{bias}$

3. stride of convolution $\rightarrow S$

3 colors

 $K_1(3l^2 + 1)$

total no. param

need to be est. as part of training.

design choice (not param.)

(A10) - review equations

'RGB' CNN

- ✓ i) convolution
- ii) pooling

. max pooling

keep track of which position had highest value

(use max res picked up from)

- i.e. reg argmax .

required for backprop

(*) note CNN jargon - kernels / channels

(*) A11 - CNN's

- diagram, equations review
- 1st pooling 2nd conv. layer

(*) size of the layers.

A12 - review

- summarises no maps, sizes of 1/10 maps. no.
- (*) general, no. of conv. filters increases with layers no...
- going through layers \rightarrow maps get smaller and smaller; maps \uparrow , no. filters/layer \uparrow .
- to retain info; want more of them
- e.g. original 3 color, 1024 x 1024 pixels
 - \hookrightarrow 3 million values
 - retain all info.
 - no. of channels \times output of each channel has to be 3 million
 - else info loss

(*) information loss in general a good thing acc. to BR

(*) parameters to choose (design choices)

- lots of design choices, but not damning

(*) Training

- don't be put off
- large one giant MLP
- shared params networks
- pool derivatives

(*) only difference is structure of network

- next class in depth

(*) convolution does not change anything

(A) learning AN

(B). review

(C) Note connec/refresh ^{per instance}
training loss - divergence - etc.