

L12-CNNs pt 10

BR: last lecture on CNNs (ensure you complete CNN material in A1)

(*) There
are differences

between
this S. 2019
and S. 2020

slides

- Make sure
you can
reconcile

(*) note the conventional CNN hierarchical story.

(*) General architec. of NN

- recall downsampling is optional
- brief recap of convolutional layer

- convolution
- activation

(*) (A): A little rusty on multiple maps element of CNNs
+ convolutional maps

- this is essential → review.

(*) CNN vector not.

(A) - understand distinction in pseudocode between using stride = 1
and stride > 1.

- for CNN vector
- can shrink maps via
 - i) convolution with stride > 1
 - ii) downsampling

(*) pooling / downsampling

- max pooling → typically stride > 1
- after convolutional layer (i.e. convolution and activation already applied)
max pooling is applied
- 1) searching for index of largest value (see pseudocode!)
- `maxidx(...)`
-
- 2) store value at that index location
- performed independently for every map (max pooling)
- convolutions combine maps

- (*) Mean pooling → use different op.
- Image is RGB → stack ≈ 3 layers → cube
- (*) each layer of CNN consists of several stacked filters → cuboid
- (*) these filters scan input and produce a collection of planes (conv.)
- (*) over which you perform pooling to reduce no. of planes.
- convolve - pool - c - p -
- (*) Then flatten and pass {C-P-C-P} output through softmax / MLP.
- final

-
- (*) Learning network
 - Parameters to be estimated:-
 - weights of neurons in final MLP
 - weights and biases (aka filters) for every conv. layer

Recap finished

- (*) Parameters just like conventional MLP :-

but i) shared parameters } giant MLP with shared params.
 ii) scanning input.

- (*) Training examples

- (*) Variant of gradient descent; use backprop

Defining loss

- typically cross-entropy or KL divergence between network output, desired output

Problem setup

overall loss is average of divergences over training set

$$\text{loss} = \frac{1}{T} \sum_{i=1}^T \text{div}(y_i, d_i)$$

(*) Training CNNs through g.d.

(*) Need to update weights/biases for every layer; and for every filter

$$w(l, m, j, x, y) = w(l, m, j, x, y) - \eta \frac{d \text{ loss}}{dw(l, m, j, x, y)}$$

② (A) - $l, m, j, x, y \rightarrow$ clarify what these mean (can't recall)

Total train loss: $\text{loss} = \frac{1}{T} \sum_i \text{div}(Y_i, d_i)$

Total deriv: $\frac{d \text{ loss}}{dw(l, m, j, x, y)} = \frac{1}{T} \sum_i \frac{d \text{ div}(Y_i, d_i)}{dw(l, m, j, x, y)}$

(*) Backprop kicks in here

(*) Backprop - final flat layers

(*) Backprop does not differ from MLP up to flattened layer; only reading past flattened layer do things change.

(*) Recall:- derivatives wrt - i) weights, biases (as a consequence)

ii) output of previous layer (recursively)

(*) final layer - $\gamma(x)$

$$\text{derivative of div wrt final layer} = \nabla_{\gamma(l)} \text{div}(\gamma(x), d)$$

(*) Backprop - convol. and pooling

(A) Special treatment of:-

- i) shared computation in conv. layers
- ii) Pooling layers

(*) B.P. - conv layer

(*) B.P. - conv layer

(*) consider what is happening during forward propagation through conv. layer

(*) And, derivative gives info about change of a variable; and effect of that change on final function (divergence)

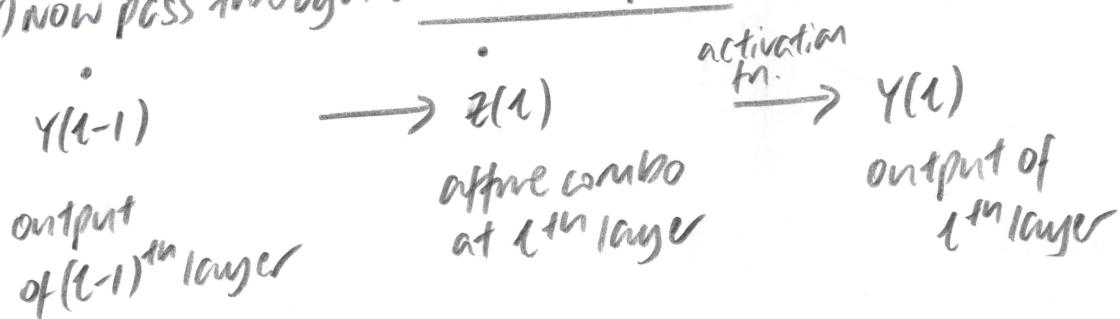
- focus on blue box in Y(1-1)

- focus on blue box in ReLU_1
- (x) for every i^{th} layer filter (convolution); each position in map in $(i-1)^{\text{th}}$ layer affects several positions in the map of i^{th} layer.

(*) In slides; the blue box features in every 'scan' → see slides

(x) True res convolution

- (*) that was convolution
- (*) now pass through activation part of C.L.



- Backprop (one recursive step):-

$$\nabla_{Y(1)} \text{div} \quad \leftarrow \quad \nabla_{Z(1)} \text{div} \quad \leftarrow \quad \nabla_{Y(1)} \text{div}$$

- recall that Δ_{typ} div is given from previous layer

- we need to compute $\nabla_{x(t)} \text{div}$ and $\nabla_{y(t-1)} \text{div}$.

- we need to compute $\nabla_{\theta(t)}$
- not treating bias independently from weights rotationally (absorb)

(*) BP. convlayer

- computing $\mathbb{Z}_{2(n)}$ div :-

$$\frac{d\text{Div}}{dz(t, m, x, y)} = \frac{d\text{Div}}{dy(t, m, x, y)} f'(z(t, m, x, y))$$

- component-wise computation

1-layer
m-map
x-y-position

(I)

(*) BP convlayer

(*) computing $\nabla_{Y(t-1)} \text{div} ; \nabla_{W(t)} \text{div}$

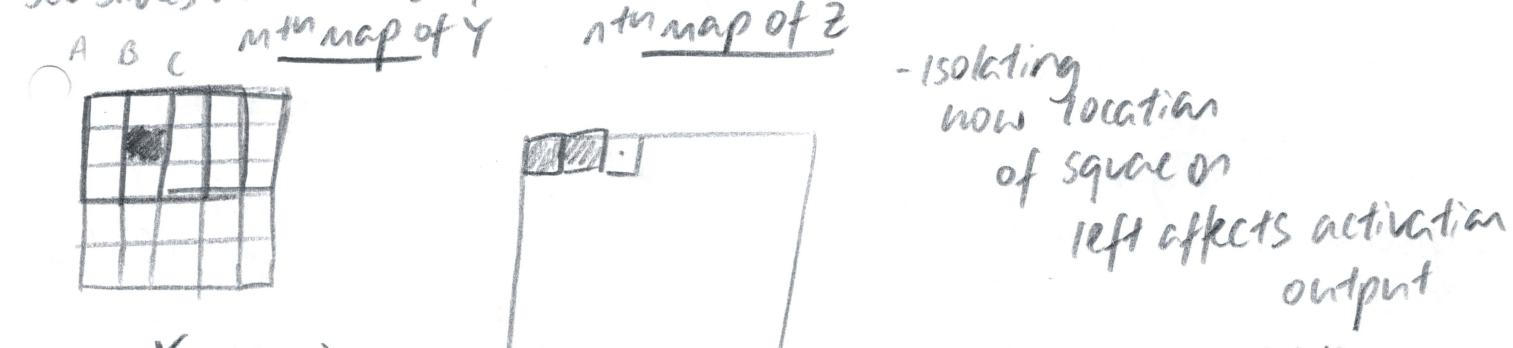
(*) Each $y(t-1, m, x, y)$ affects several $z(t, n, x, y)$ terms.

- All contribute to derivative wrt $y(t-1, m, x, y)$

(*) BP-conv. layer

(*) Each $y(t-1, m, x, y)$ affects several $z(t, n, x, y)$ terms.

- see slides with arrows from $Y(t-1)$ to $Z(t)$



- consider 3 filter locations at A, B, C

- A way of visualis. influence

- way of populating sq.



- As square is outside of filter C, there is no influence of square on activation.

- Hence empty square in Z .

(*) The last Z that y will influence \rightarrow any principle?

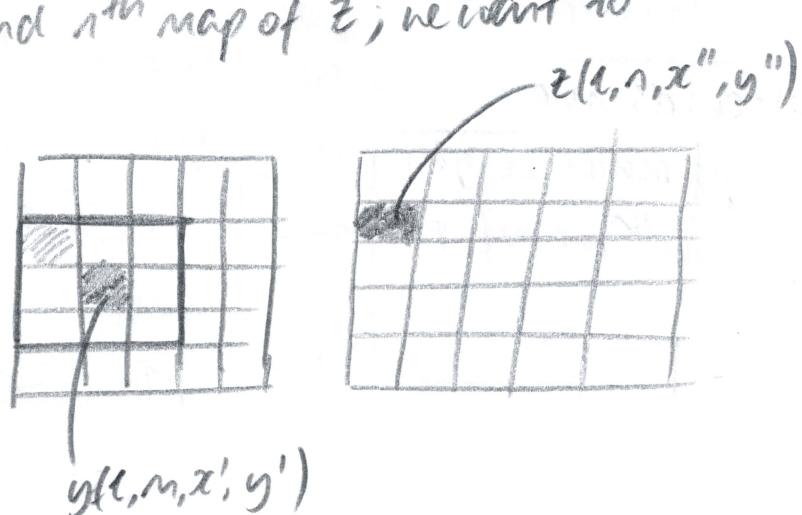
- YES. occurs when dark square in Y is in top-left (front-left) of filter

(*) For a particular m^{th} map of Y and n^{th} map of Z ; we want to assess influence of

$y(t, m, x', y')$ on $z(t, n, x'', y'')$

where (x', y') and (x'', y'')

(*) To do so; must place filter with 'this Z ' at top-left corner



- we want weight that multiplies $y(\ell, m, x', y')$ to $z(\ell, n, x'', y'')$
transforms -
 - each $y(\ell-1, m, x', y')$ affects several $z(\ell, n, x'', y'')$ terms through weight $w(\ell, n, x''-x'', y''-y'')$
 - each $y(\ell-1, m, x', y')$ affects several $z(\ell, n, x'', y'')$ terms through weight $w(\ell, n, x''-x'', y''-y'')$ (all) (on)
 - for all ℓ^{th} layer z-maps $\{z(\ell, 0, x'', y''), z(\ell, 1, x'', y''), \dots, z(\ell, N, x'', y'')\}$
($n=0, 1, \dots, N$) deriv.
 - for all ℓ^{th} layer z-maps $\{z(\ell, 0, x'', y''), z(\ell, 1, x'', y''), \dots, z(\ell, N, x'', y'')\}$
($n=0, 1, \dots, N$) deriv.
 - can you see multiplicity of influence paths now?
- $\sum \frac{d\text{Div}}{dz} \cdot \frac{dz}{dy}$
"z influenced by y"
- $\frac{d\text{Div}}{dy(m, x', y')} = \sum_{n=1}^N \sum_{x'', y''} \frac{\frac{d\text{Div}}{dz(n, x'', y'')}}{dz(n, x'', y'')} \cdot \frac{dz(n, x'', y'')}{dy(m, x', y')}$
 $= \sum_{n=1}^N \sum_{x'', y''} \frac{\frac{d\text{Div}}{dz(n, x'', y'')}}{dz(n, x'', y'')} \cdot w(m, n, x''-x'', y''-y'')$
(II)
- weights ↴
- note:-
 $y(\ell-1, m, x, y) = y(m, x, y)$
 $z(\ell, n, x', y') = z(n, x', y')$
(no ℓ layer m is)

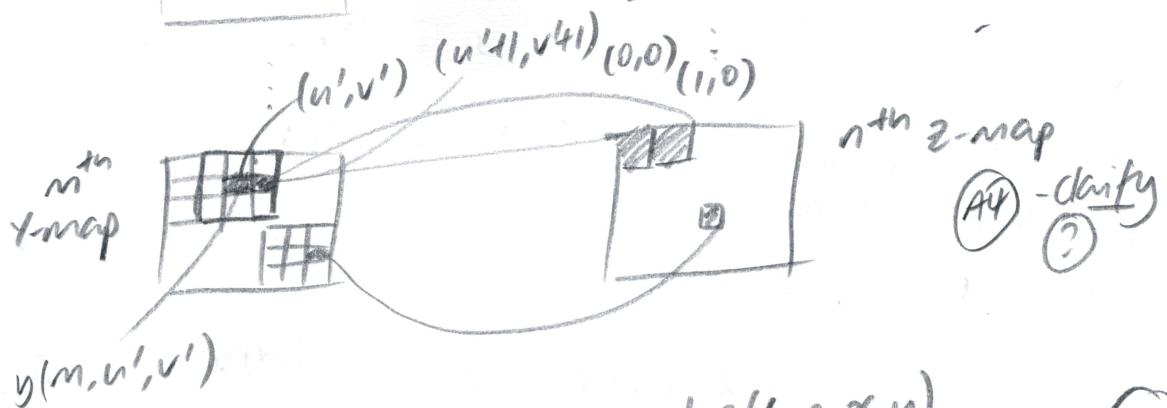
- BR uses (x, y) and (x', y') to refer to pos. of y-map and z-map
- I used (x, y) and (x', y')
- I used (x', y') and (x'', y'')
(to differentiate between argument and instantiation)

-
- (*) Towards finding derivative of weights (BP conv. layer)
- A particular weight $w(m, n, x, y)$ governs the rel. between the m^{th} map of Y to the n^{th} map of Z

$w(m, n, u', v')$



(*) - slides missing
- see SW22 slides



nth z-map

(A4) - clarify
?

- each weight $w_e(m, n, x', y')$ also affects several $z(l, n, x, y)$ (A3)

for every l

Affects items in only one z-map

All entries in map contrib. to the derivative of the divergence not $w_e(m, n, x', y')$ (3?)

$$\frac{d\text{div}}{dw_e(m, n, x, y)} = \sum_{x'', y''} \frac{d\text{div}}{dz(l, n, x'', y'')} \frac{dz(l, n, x'', y'')}{dw_e(m, n, x, y)}$$

$$= \sum_{x'', y''} \frac{d\text{div}}{dz(l, n, x'', y'')} y(l-1, m, x''+x, y''+y)$$

(*) De my!
Simple than maths

) - compute mulp within each loop

(*) CNN forward

- key principle:-

$$z = wY \quad \frac{df}{dw} := \frac{df}{dz} Y$$

f(z)

$$\frac{af}{dY} := \frac{df}{dz} w$$

(*) Go through layers, filters;
each filter scans input.
position filter at a
partic. location, compute
pairwise product of
filter pos and y and
add them up.

- use this for CNN backward layer l

(*) Backward layer l

- (*) Pseudocode puts (I), (II), (III) together for CNN backprop
- (*) only 3 lines changed in backprop compared to forward prop
- (*) Multiple way of casting as tensors
-
- (*) Backward with strides
- messy algebraically; but trivial in code. (?)
 - (AS): make sure you included distinction between pseudocode.
-
- (*) How a single $y(t-1, m, x, y)$ influences $z(t, n, x', y')$
- examine how each (x, y) position in m^{th} y -map (one map) influences various locations in an arbitrary n^{th} z -map
 - previously, for CNN backprop, we examined how m^{th} y -map influenced multiple N z -maps
- (*) (A6): Review interactive, color coded slides
- (*) Have to reverse direction of influence to compute the deriv. of that x, y component of Y .

(*) Each z is sum of component-wise product of filter elements and elements of the region of Y it is placed on.

(*) (A7): Influence pattern is mirror image of weights

(*) (A8): Review formulae here; index correspondences

$$(67): - \frac{\partial \text{Div}}{\partial y(t-1, m, x, y)} = \sum_n \sum_{x', y'} \frac{\partial \text{Div}}{\partial z(t, n, x-x', y-y')} w_{t, n, x', y'} \quad \text{we}(m, n, x', y')$$

(*) (A9): write some examples if this is unclear during review. mirror image

(*) computing derivative.

(*) Each deriv wrt y at (x, y) is obtained by flipping the filter left-right, top-bottom and computing the inner product

with respect to the square patch of $\frac{\partial \text{Div}}{\partial z}$ ending at (x, y) .

- this would be for any (x, y) .

Q8: Draw diagrams during review (see slides)

Q9: $\frac{\partial \text{Div}}{\partial z(l, n, x, y)}$ is obtained by taking filter $w_l(m, n, *, *)$

• then we flip filter left-right, top-to-bottom; and convolve with ✓

$$\frac{\partial \text{Div}}{\partial z(l, n, x', y')}$$

Q10: except we are not just interested in one map - interested in many

(*) Moving to derivatives over multiple maps

Q11: In reality, derivative at each (x, y) location is obtained from all z-maps

- filter will have one place for each map

- each of these will be mirror-imaged (is this kernel flipping?)

Q12: see slides, colored interactive to go up another level of abstraction
from this.

BR: maths is tedious, consuming; but rely on visual intuition here
to minimise the difficulty.

(*) observe: to compute derivative at top-left corner, filter begins outside
the $\frac{\partial \text{Div}}{\partial z(l, n, x', y')}$ plane

- more precisely; convolution of flipped filter and $\frac{\partial \text{Div}}{\partial z(l, n, x', y')}$

only results multiplication of bottom right corner

of flipped filter and top-left corner of $\frac{\partial \text{Div}}{\partial z(l, n, x', y')}$ at

beginning of 'scanning' from top-left

(*) note flipped filter will be 'hanging off edge' (see slides)

(*) backprop → have to allow filter to fall off edge
→ require zero-padding

(*) computing deriv.

- convolution of $\frac{\partial \text{Div}}{\partial z(l, n, x, y)}$ by the inverted filter

- after 0 padding with $(l-1)$ zeros on every side
 $\sim (?)$

- can rely on forward code - amend.

(*) derivatives for a single layer l

- vector not.

A10 - review pseudocode

(*) for derivatives w.r.t weights (as opp. to y)

- interpretation needs tweaking...

(*) pooling and downsampling (stride > 1 typically)

(*) max pooling

- At each location:-

- i) find position of input where value is maximum
- ii) copy value at that position

- recall defn of deriv; change of output not change of prem
BR: If 6 is 'copied' value; will modif. of 1,3,5 change output?

- NO. (as they do not 'influence' output)

"derivative with respect to 1,3,5" - 0

(A11) - review max pooling equations



(~) (?) A11 - tighter logic
loose

(*) derivative of max pooling

$$(iii): \frac{d\text{div}}{dy(\ell, m, k, l)} = \begin{cases} \frac{d\text{div}}{du(\ell, m, i, j)} & \text{if } (k, l) = p(\ell, m, i, j) \\ 0 & \text{otherwise} \end{cases}$$

MAX

pooling op: $p(\ell, m, i, j) = \underset{\substack{k \in \{(i-1)d+1, (i-1)d+K_{\text{pool}}\} \\ n \in \{(j-1)d+1, (j-1)d+K_{\text{pool}}\}}}{\operatorname{argmax}} y(\ell, m, k, n)$ (A12) review

$$u(\ell, m, i, j) = y(\ell, m, p(\ell, m, i, j))$$

(*) max pooling at layer ℓ

- pseudocode (A12)

(*) Mean pooling:-

mean pooling op: $u(\ell, m, i, j) = \frac{1}{K_{\text{pool}}^2} \sum_{\substack{k \in \{(i-1)d+1, (i-1)d+K_{\text{pool}}\} \\ n \in \{(j-1)d+1, (j-1)d+K_{\text{pool}}\}}} y(\ell, m, k, n)$

(*) Derivative of mean pooling:

- derivative of mean pooling is distn over the pool:-

$$\begin{matrix} k \in \dots \\ n \in \dots \end{matrix} \quad (\text{from above}); \quad \frac{dy(\ell, m, k, n)}{du(\ell, m, k, n)} = \frac{1}{K_{\text{pool}}^2}$$

(*) mean pooling layer at layer ℓ + derivative

- pseudocode (A13)

- type 0 - $dy(\ell, j, p, n+i, m+j) += (1/K_{\text{pool}}^2) du(l, j, x, y)$

(*) learning network

- Backprop for CNN complete (after review!!)
- Now embed in gradient descent

(*) implicit assumption

- Can subsequent maps increase in size?

(*) 1D figure

- Scanning with stride of 18
- Arrow \Rightarrow contains 6 weights (2 input neurons \rightarrow 3 output neurons)
- distributed parameter network
- individual neurons have smaller strides (covered in H/W)
- No. of circles within bar \rightarrow no. of channels
- Scanning \Rightarrow no. of bars stays the same or goes down
(tree structure)
- stride > 1 \Rightarrow no. of bars reduces

(ai2) - logic - make
see water tight

59:00 \rightarrow

with layers of increased resolution

- new bars' \rightarrow middle gets only one contributions from lower bars
left and right gets multiple -||- — — ||—

(*) transposed convolution

(ai3) - seriously confused

- see this for 2-D

(*) 2D-expanding convolution

- output size typically an integer multiple of input
- +1 if filter width is odd.

$$z(l, i, j) = \sum_m \sum_k w(l, m, i-k, j-l) I(m, k, l)$$

b is stride \rightarrow scaling factor between size of z -map and I -map.

(*) Recasting CNN as shared param MLP; expanding layers have weight matrices that are taller than wide.

(*) P no. of weights matrices $\Rightarrow P$ no. output channels

\uparrow stride \Rightarrow \uparrow size of individual maps

(*) use scale factors to better understand relation between scale factor between 2 maps and stride

(*) Invariance:-

- shift invariance accounted for.
- How about rotation/scale/reflection invariance?

(*) Shift invariance - different perspective

(A14) rotation invariance

(*) Transform invariance

Each filter produces a set of transformed maps

(*) network becomes invariant to all transforms considered

- affects backprop

(*) Transform invariance

- (A15) - review pseudocode (enumerate over transforms;
compute a scan with transformed filter)

(*) BP with transform invariance

- derivatives flow back through transforms to update individual filters.

- as long as you have correspondence between original and transformed filter, all fine.
 - computations not covered in some amount of depth.
-
- ### (A) exact location of object
- implicitly → by recording where max is triggered in final layer
 - explicitly → incorporate additional output layer
 - predicts corners of bounding box of object (8 coordinates)
- Ⓐ: why 8?
- can be used for pose estimation
- training: compute 2 divergences over 2 output layers (conventional, bound box)
- final divergence is weighted sum of 2 divergences
 - (e.g. X-entropy of class. layer + L_2 loss of bounding box)

(B) Model Variations

- depth-wise convolutions → MobilNet (HW2)

(B) depth-wise convolutions

- alternate view of convolution Ⓛ (standard) Ⓛ - review
 - not clear
- | depth-wise convol. as a switch in order
- | only perform a 'local convolution' at each map?
- | - layerwise convolution over all maps and then add
 - (think this just a slight variation on how to view the operation from a matrix/tensor multiplication perspective).
- | - still some op

⑦ Ⓛ: collapsing? → review convolution op.

(*) convolutional vs depth-wise convol.

- Review this slide
- Shows efficiency (in terms of no. parameters required) for depth wise convol.
- note use of 30 vs 20 filters in terms of no. parameters

(*) what do filters actually learn/test?

- receptive fields
- subsequent neuron weights after 1st (more difficult)
- ImageNet

(*) Training Issues

- Convergence issues ↳ review slides, for history of ImageNet
- Data aug.

