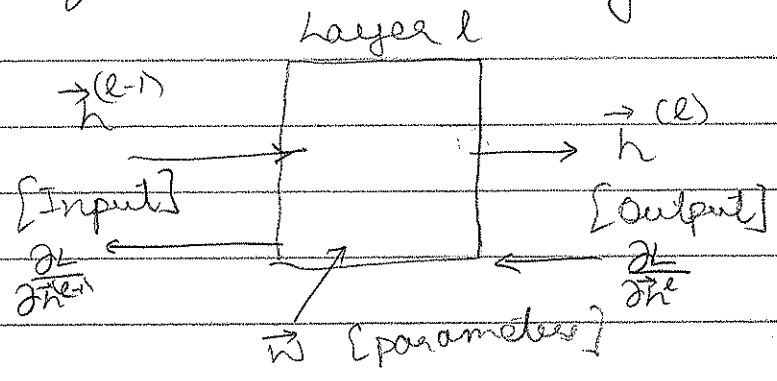


9/03/15

(1)

MLPs + CNNs

① Recall from last time, key abstraction:



$$\vec{h}^{(l)} = g(\vec{h}^{(l-1)}, \vec{w})$$

$$\text{Loss } L = f(\vec{h}^l)$$

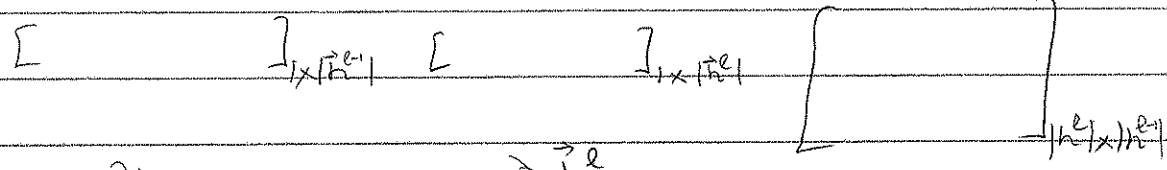
Abstractly: (a) $\frac{\partial L}{\partial \vec{in}} = \frac{\partial L}{\partial \vec{out}} \cdot \frac{\partial \vec{out}}{\partial \vec{in}}$

(b) $\frac{\partial L}{\partial \vec{w}} = \frac{\partial L}{\partial \vec{out}} \cdot \frac{\partial \vec{out}}{\partial \vec{w}}$

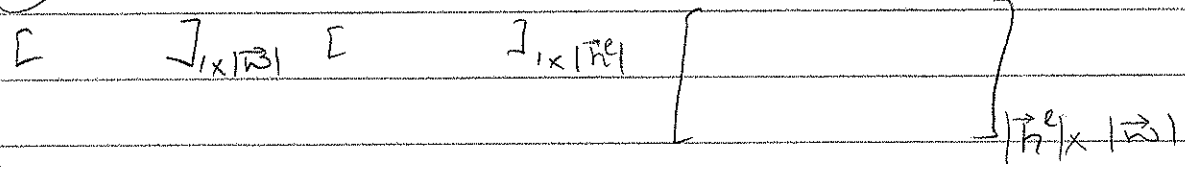
[A bit more]

Concretely:

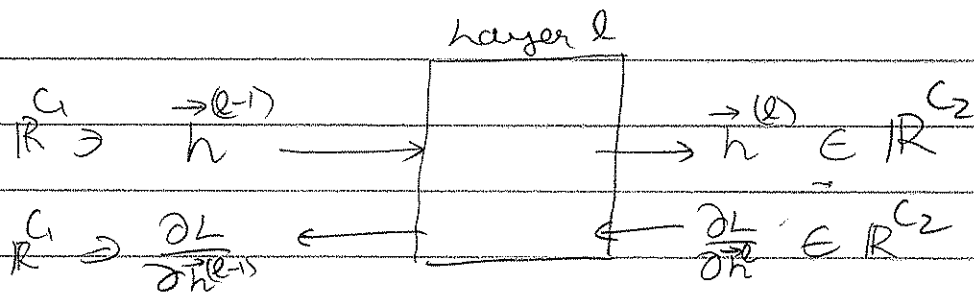
(a) $\frac{\partial L}{\partial \vec{h}^{(l-1)}} = \frac{\partial L}{\partial \vec{h}^l} \times \frac{\partial \vec{h}^l}{\partial \vec{h}^{(l-1)}}$



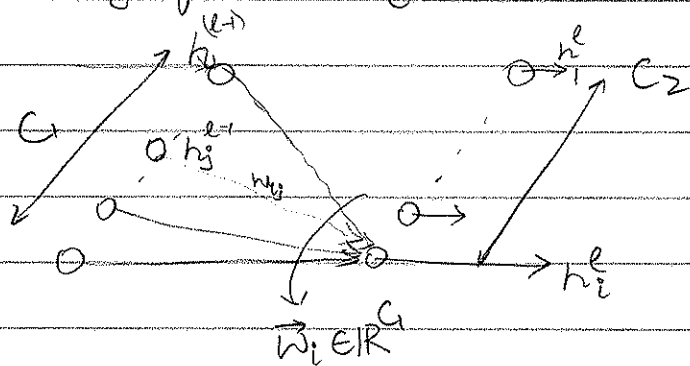
(b) $\frac{\partial L}{\partial \vec{w}} = \frac{\partial L}{\partial \vec{h}^l} \times \frac{\partial \vec{h}^l}{\partial \vec{w}}$



② Multi-Layer Perceptrons / Fully-Connected Layer / Inner-Product Layer



Another way of drawing this



[Why draw in 3D? You'll see. Wait for CNNs!]

Forward-Pass:

$$h_i^l = \sum_{j=1}^{C_1} w_{ij} h_j^{(l-1)}$$

$$= \vec{w}_i^T \vec{h}^{(l-1)}$$

Let's vectorize some more

$$\begin{bmatrix} h_1^l \\ \vdots \\ h_{C_2}^l \end{bmatrix} = W^{(l)} \vec{h}^{(l-1)}$$

$$= \begin{bmatrix} \leftarrow w_1^T \rightarrow \\ \leftarrow w_2^T \rightarrow \\ \vdots \\ \leftarrow w_{C_2}^T \rightarrow \end{bmatrix} \begin{bmatrix} h_1^{(l-1)} \\ \vdots \\ h_{C_1}^{(l-1)} \end{bmatrix}$$

So all of F-PASS
= 1 matrix-mult

Why did we write this as matrix mult?

① Because it's cooler!

② Gradients, my dear Watson, Gradients...

$$\vec{h}^{(l)} = W^{(l)} \vec{h}^{(l-1)}$$

$$\rightarrow \frac{\partial \vec{h}^{(l)}}{\partial \vec{h}^{(l-1)}} = W^{(l)} \quad \left[\text{Now, isn't that beautiful?} \right]$$

What about $\frac{\partial \text{out}}{\partial \vec{w}}$?

Simple
$$h_i^l = \vec{w}_i^T \vec{h}^{(l-1)}$$

$$\rightarrow \frac{\partial h_i^l}{\partial \vec{w}_i} = \vec{h}^{(l-1)T}$$

Now put it all together:

$$\textcircled{a} \frac{\partial L}{\partial \vec{h}^{(l-1)}} = \frac{\partial L}{\partial h^l} \frac{\partial \vec{h}^l}{\partial \vec{h}^{(l-1)}} = \left[\frac{\partial L}{\partial h^l} \times W^{(l)} \right] \text{ Very nice}$$

So B-PASS (a) \equiv 1 matrix mult

Left multiply \equiv gradient
Right multiply \equiv output

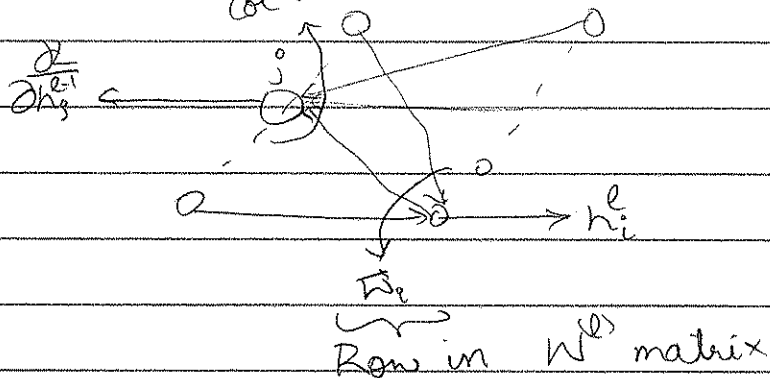
$$\textcircled{b} \frac{\partial L}{\partial \vec{w}_i} = \frac{\partial L}{\partial h_i^l} \frac{\partial h_i^l}{\partial \vec{w}_i} = \underbrace{\frac{\partial L}{\partial h_i^l}}_{\text{scalar}} \times \underbrace{\vec{h}^{(l-1)T}}_{1 \times n}$$

So $W^{(e)} = W^{(e)} - \eta \frac{\partial L}{\partial W^{(e)}}$

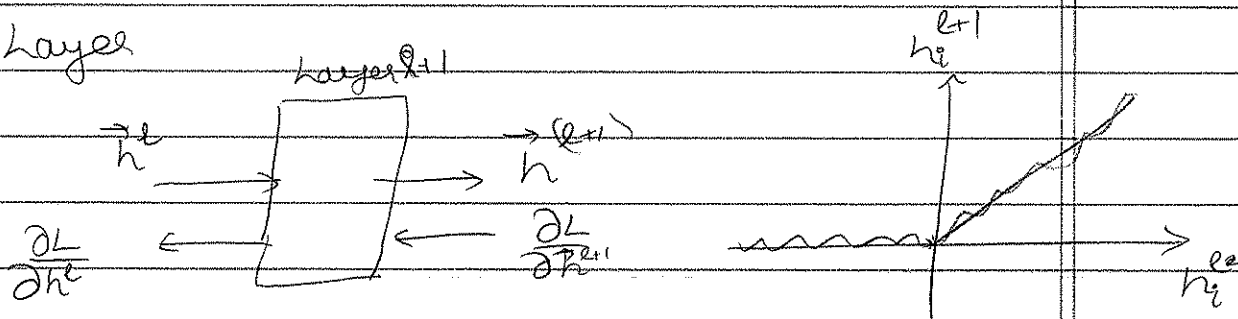
where

$$\frac{\partial L}{\partial W^{(e)}} = \begin{bmatrix} \leftarrow \frac{\partial L}{\partial h_1^{(e-1)}} \cdot h^{(e-1)T} \rightarrow \\ \vdots \\ \leftarrow \frac{\partial L}{\partial h_2^{(e-1)}} \cdot h^{(e-1)T} \rightarrow \end{bmatrix}$$

One more intuition
col in $W^{(e)}$ matrix



③ ReLU Layer



F-PASS $h_i^{(l+1)} = \max\{0, h_i^e\}$

B-PASS $\frac{\partial h_i^{(l+1)}}{\partial h_i^e} = \begin{cases} +1 & \text{if } h_i^e \geq 0 \\ 0 & \text{else} \end{cases}$
[No params]

$\Rightarrow \frac{\partial L}{\partial h_i^e} = \frac{\partial L}{\partial h_i^{(l+1)}} \times \mathbb{1}\{h_i^e \geq 0\}$

④ Recall: Convolutions! [Yes, that nightmare from signal processing class]

→ Let's start with pure math then move to math → CS → programming

Pure Math

→ Let $y(t)$, $x(t)$, $w(t)$ be continuous-time signals

1-D convolution $y(t) = (x * w)(t)$

$$\Rightarrow y(t) = \int_{a=-\infty}^{\infty} x(a) w(t-a) da$$

Intuition: to compute $y(t)$

→ Flip: $w(a) \rightarrow w(-a)$
[why? ∵ math]
we won't care in ML DL

→ Shift: $w(-a) \rightarrow w(-(a-t))$
by t

→ Compute area of product $\int_{a=-\infty}^{\infty} x(a) w(-(a-t)) da$

Repeat & Repeat $\forall t$

BTW: $y(t) = \int x(a) w(t-a) da = \int x(t-a) w(a) da$
[commutative]

Similarly 2D convolution

$$y(t_1, t_2) = \int_{a=-\infty}^{\infty} \int_{b=-\infty}^{\infty} x(t_1, t_2) w(a-t_1, b-t_2) da db$$

Now, lets
move to
CS

There is no infinite precision! Let's go discrete

$$y[r, c] = \sum_{a=-\infty}^{\infty} \sum_{b=-\infty}^{\infty} x[r, c] w[a-r, b-c]$$

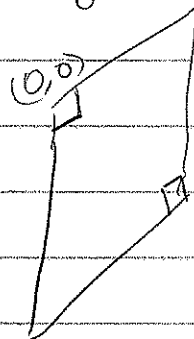
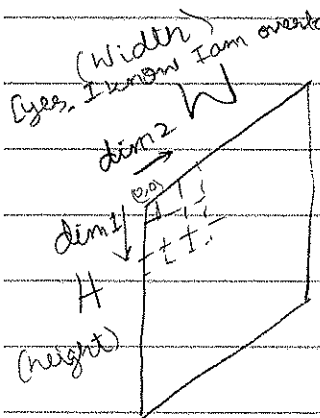
No inf memory either, Has to be finite matrices

$$y[r, c] = \sum_{a=-\frac{H+1}{2}}^{\frac{H+1}{2}} \sum_{b=-\frac{W+1}{2}}^{\frac{W+1}{2}} x[r, c] w[a-r, b-c]$$

lets move
to programming

Arrays don't index with negatives!

And we don't really care about commutivity!



Input

Kernel/
Filter / weight

Output [size depends on
'valid' or 'same'
conv]

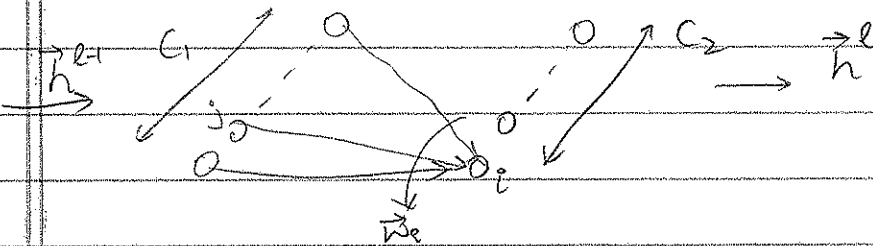
What we actually implement:

(4)

$$y[a, c] = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} x[a+a, c+b] w[a, b]$$

(5) CNNs / Conv Nets / Convolutional Neural Nets

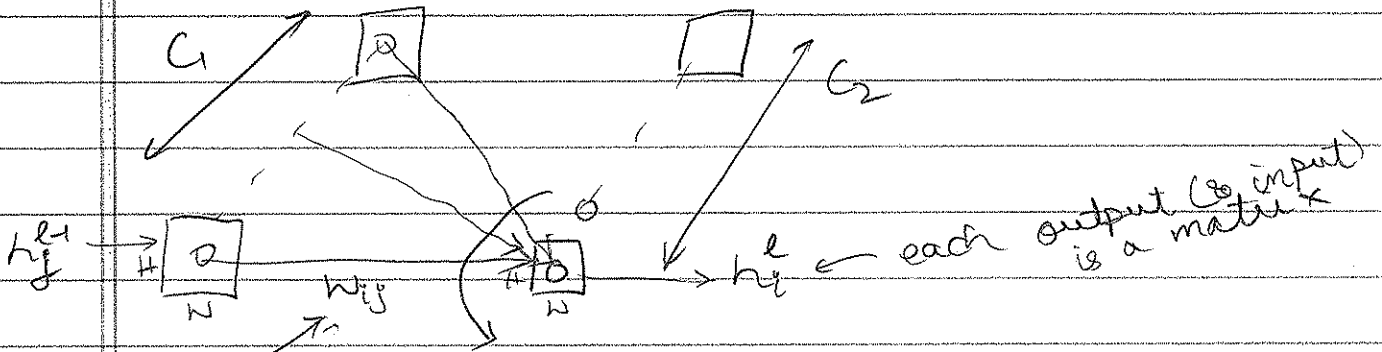
→ Recall MLP



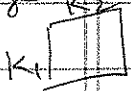
$$h_i^l = \sum_{j=1}^{C_1} w_{ij}^{(l-1)} h_j^{(l-1)}$$

↑
scalar product

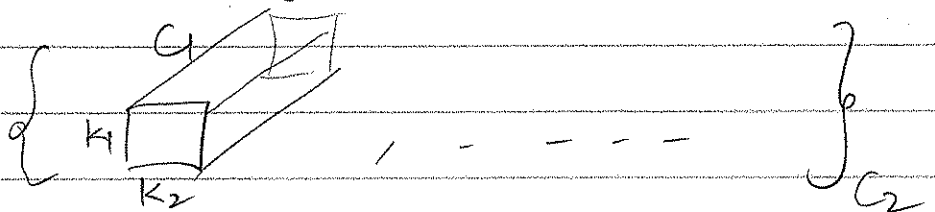
→ CNN [Make each neuron a 2D matrix]



Now a matrix / filter on every edge



All weights at output used i



CNN - F-Pass

Compose:

$$\text{MLP: } h_i^l = \sum_{j=1}^{C_1} w_{ij}^{(l)} \cdot h_j^{l-1}$$

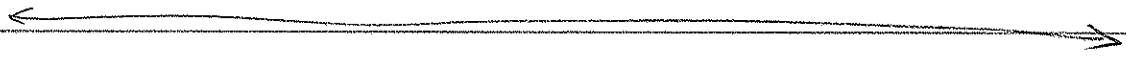
↙ scalar product
↖ scalar

$$\text{CNN: } h_i^l = \sum_{j=1}^{C_1} w_{ij}^{(l)} * h_j^{l-1}$$

↖ matrix
↖ 2D conv

$$= \left[\sum_{j=1}^{C_1} \right] \left[\sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} \right] h_j^{(l-1)} [a+a, c+b] w_{ij}^{(l)} [a,b]$$

Input Channels
2D conv



CNN - B-Pass

→ This is going to be a notation nightmare

→ So let's simplify thing a bit

→ Say $C_1 = 1$ $C_2 = 1$ [we can always generalize] its just an extra sum

→ And let's not over-load h

