

BLG 561 E FALL 2021

Deep Learning

09.11.2021

Görde ÜNAL

3 components of a general ML algorithm : DL is a sub-class

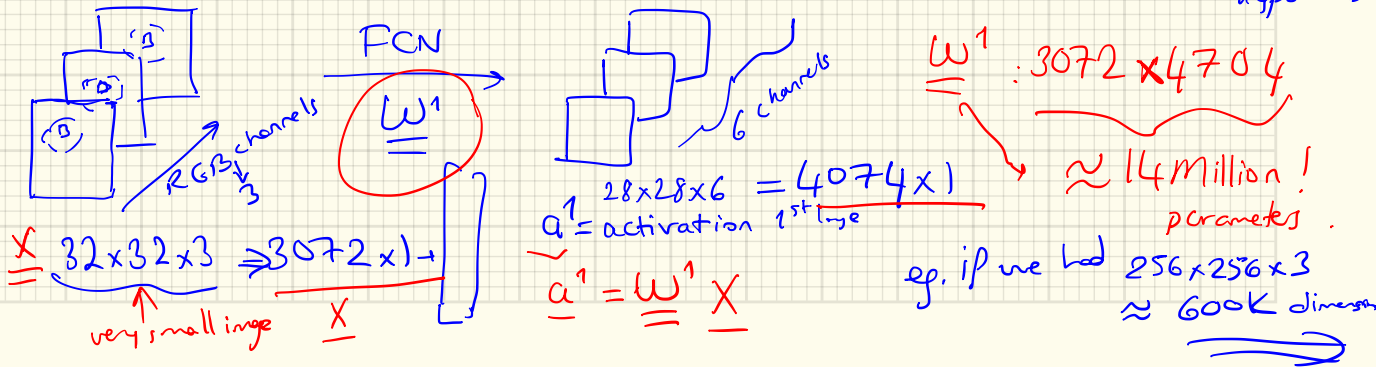
- 1) Hypothesis function: ↗
- 2) Loss function: ↘
- 3) Optimization : optime loss fr.

Today: CONVOLUTIONAL NEURAL NETWORKS (CNNs)

★ Only the hypothesis functions change, i.e. the ①st component ONLY

Then you can pick ② & ③ freely according to your task.

Now → Say, you have large (spatial) data, let's use a FCN / Fully Connected hypothesis

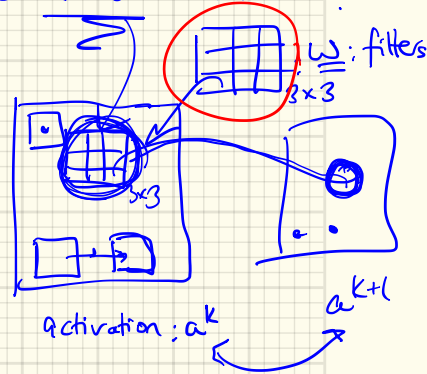


→ FCNs would need a huge # parameters!!

⇒ CNNs can handle large spatial data. @: How?

They restrict the weights in 2 ways: (in contrast to FCNs)

- 1) Activations between layers occur in a "local" manner:  
(sparsity):
- 2) Parameter sharing: All activations share the same weight.



① 1st component:

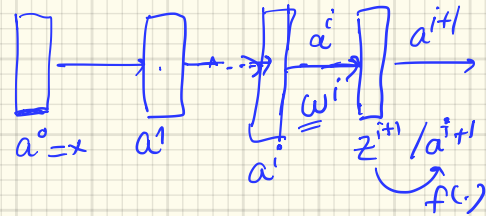
New hypothesis function for CNNs:

$$z^{i+1} = \underline{w}^i * a^i + b^i$$

$$a^{i+1} = f(z^{i+1}) = f(w^i * a^i + b^i)$$

\*: convolution operation

f: ReLU or other nonlinearity



# Compare CNN

convolution

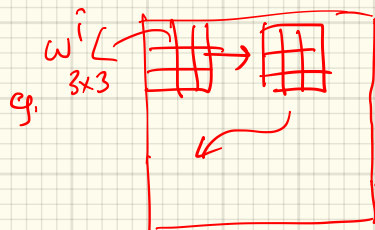
$$z^{i+1} = \underline{\underline{w}}^i * \underline{a}^i + \underline{b}^i$$

$$a^{i+1} = f(\underline{\underline{w}}^i * \underline{a}^i + \underline{b}^i) = f(z^{i+1})$$

# FCN

matrix-vector multiplication

$$\begin{cases} z^{i+1} = \underline{\underline{w}}^i \cdot \underline{z}^i + \underline{b}^i \\ a^{i+1} = f(z^{i+1}) \end{cases}$$



- ① locality (sparsity)
- ② shared connections

Note: other components can be added  
maxpool  
BN etc.


1D Conv

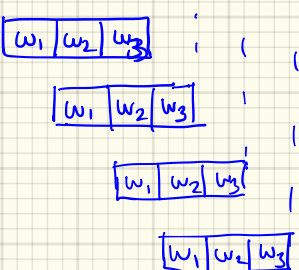


For illustration:

→ 1-D CNN.  $a^{i+1} = f(w^i * a^i + b^i)$

Say  $w^i = (w_1, w_2, w_3)$ : 1D filter

Input:  $a^i \Rightarrow$  

$w^i$  filter:  $3 \times 1$  

4 possible "valid" shifts

Write convolution to a matrix operator

$$\underline{w^i} = \begin{bmatrix} w_1 & w_2 & w_3 & 0 & 0 & 0 \\ 0 & w_1 & w_2 & w_3 & 0 & 0 \\ 0 & 0 & w_1 & w_2 & w_3 & 0 \\ 0 & 0 & 0 & w_1 & w_2 & w_3 \end{bmatrix}_{4 \times 6}$$

CNN:  $a^{i+1} = f(w^i \cdot a^i + b^i)$

FCN

What happens at backprop?

Recall for FC: backprop:

$$\text{upstream gradient: } \frac{\partial \mathcal{L}}{\partial a^i} = \underline{g^i} = \underbrace{(w^{i+1})^T}_{6 \times 4} \cdot \underbrace{g^{i+1}}_{4 \times 1} \cdot \underbrace{f'(z^{i+1})}_{1 \times 1}$$

→ let's transpose

in CNN backprop :

$$\underline{W_i^T} = \begin{bmatrix} w_1 & 0 & 0 & 0 \\ w_2 & w_1 & 0 & 0 \\ w_3 & w_2 & w_1 & 0 \\ 0 & w_3 & w_2 & w_1 \\ 0 & 0 & w_3 & w_2 \\ 0 & 0 & 0 & w_3 \end{bmatrix}_{6 \times 4}$$

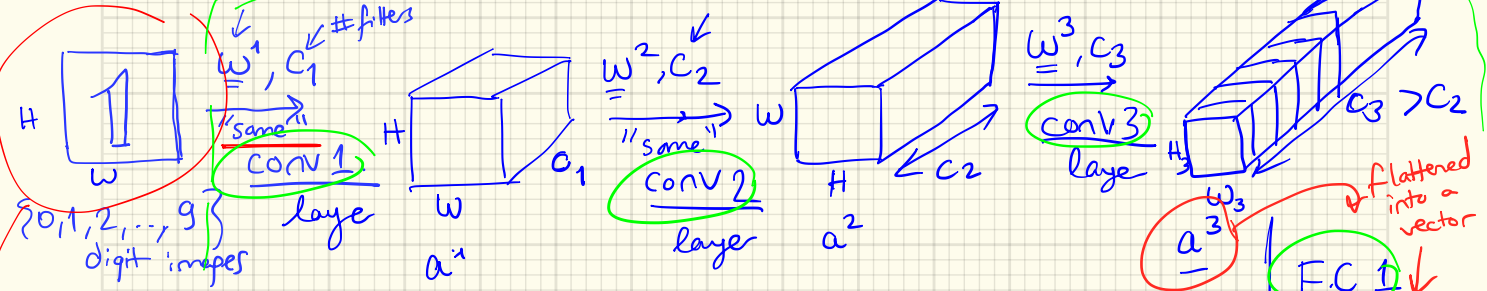
→ filter is flipped

∴ In CNNs, backprop just flips convolutions.

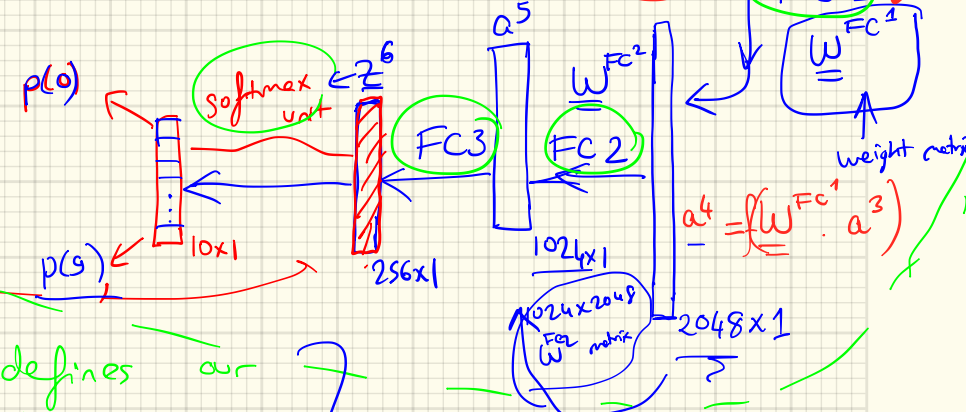
2D Conv: flip vertically & horizontally.

# Encoder Module

Ex: CNN <sup>mainly</sup> for images. Build a CNN for Digit Classification (10 class)



{0,1,2,...,9} digit images

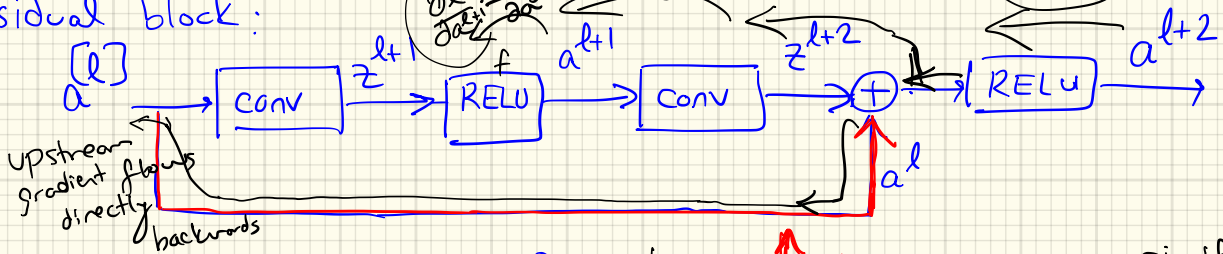


- ① This whole model defines our hypothesis function  $\rightarrow$
- ② loss fn: cross-entropy, hinge loss, multi-class
- ③ pick an optimizer, e.g. Adam.

✓ You have a deep learning CNN classifier.

Ex: Res-Nets <sup>add</sup> Residual Blocks:

1 residual block:



$$a^{l+1} = f(w^l * a^l + b^l)$$

$$z^{l+2} = w^{l+1} * a^{l+1} + b^{l+1}$$

$$a^{l+2} = f(z^{l+2} + a^l)$$

$\rightarrow$  newly added in a ResNet block

$$a^{l+2} = f(w^{l+1} * f(w^l * a^l + b^l) + b^{l+1} + a^l)$$

$\uparrow$  skip connection : significance.

- think of vanishing gradient problem
- $\therefore$  ResNet : its skip connections helped reduce vanishing grad. problem, particularly in deep networks

(ResNet 2015) :  $\approx$  150 layers !! very deep network