

$\text{py} \Rightarrow *$  (element wise multiplication)

Page No.

Date

## CONVOLUTIONAL NEURAL NETWORKS

### • computer vision

- Image classification (0/1)
- Object detection (identify objects & draw boxes around them)
- Neural style transfer (required image = content image + style image)

RGB images (3 channels  $\Rightarrow 64 \times 64 \times 3$ )

Grayscale images (1 channel  $\Rightarrow 64 \times 64 \times 1$ )

### • Edge detection

CONVOLUTION  
 $\downarrow$   
 6x6 grayscale image       $\circledast 3 \times 3$  filter or kernel       $\rightarrow 4 \times 4$  matrix (image)

vertical edge detection:

$$\begin{matrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{matrix}$$

10	10	10	0	0	0	*	1	0	-1		0	30	30	0
10	10	10	0	0	0	*	1	0	-1	=	0	30	30	0
10	10	10	0	0	0	*	1	0	-1		0	30	30	0
10	10	10	0	0	0	*	1	0	-1		0	30	30	0
10	10	10	0	0	0	*	1	0	-1		0	30	30	0
10	10	10	0	0	0	*	1	0	-1		0	30	30	0
Light $\rightarrow$ Dark							3	0	0					

dark-light-dark

0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10

0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10

Dark  $\rightarrow$  Light      3    0    0    Grey-Black

Dark  $\rightarrow$  Light      3    0    0    -Grey

Horizontal edge detection      1    1    1    w

0	0	0	G
-1	-1	-1	8

Alternatives:

### • vertical edge

1    0    -1     $\Rightarrow$  Sobel : more weight given to filter central row

2    0    -2    filter

1    0    -1

3    0    -3     $\Rightarrow$  Scharr

10    0    -10    filter

3    0    -3

\* without hand-coding, the NN can learn weights assigned to detect any angle rotation.

### • padding

$$\text{input image } n \times n \quad * \text{ filter kernel } f \times f = \text{output image } (n-f+1) \times (n-f+1)$$

disadvantages: (w/o padding)

→ shrunked output

→ Throw away info from edge.

padding input image preserves the original size after convolution

$6 \times 6$

$$\downarrow \text{padding} \quad \text{filter } 3 \times 3 \quad * \text{ o/p image } 6 \times 6$$

padding ( $p=1$ ) ;  $n=6$ ,  $f=3$

$$o/p = (n+2p-f+1) \times (n+2p-f+1)$$

$$= (6+2-3+1) \times (6+2-3+1)$$

$$o/p = 6 \times 6$$

\* valid & same convolutions

(a) valid  $\rightarrow$  no padding ( $p=0$ )

$$(m \times n) * (f \times f) = (n-f+1) \times (n-f+1)$$

(b) same  $\rightarrow$  pad so that output size is same as input size

padding =  $p$

$$o/p = (n+2p-f+1) \times (n+2p-f+1)$$

$$n+2p-f+1 = ? \Rightarrow p = \frac{f-1}{2}$$

NOTE  $f$  is usually odd

→ has a central pixel

- Strided convolution  
STRIDE=2

$$7 \times 7 * 3 \times 3 = 3 \times 3$$

input image	kernel	output image
-------------	--------	--------------

$n \times n$  = input image

$p$  = padding

$s$  = stride

$f \times f$  = kernel

$$7+0-3 = 3 \times 3$$

$$2 \times 2$$

IMP  $\therefore$  output image =  $\left[ \frac{n+2p-f}{s} + 1 \right] \times \left[ \frac{n+2p-f}{s} + 1 \right]$

$\lfloor z \rfloor = \text{floor}(z) \approx \text{rounding } z \text{ down to nearest integer}$

• Associativity  $\Rightarrow A * (B * C) = (A * B) * C$

• Convolutions over volumes (3D)  
eg: convolutions over RGB images

ht → width, (filter)  $\rightarrow$  height, (channel)  $\rightarrow$  depth

$$6 \times 6 \times 3 * 3 \times 3 \times 3 = 4 \times 4$$

(image) no. of channels (equal)  $\rightarrow$  depth

3D

3D

2D

→ Filter to detect vertical edges in red channel: 10100000000000000000000000000000

R G B

1	0	-1	0	0	0	0	0
1	0	-1	0	0	0	0	0
1	0	-1	0	0	0	0	0

→ Detect edges in any colour:

filter = 1 0 -1

1 0 -1

1 0 -1

NOTE  $(n \times n \times n_c) * (f \times f \times n_c) = (n-f+1 \times n-f+1 \times n_c)$

$$6 \times 6 \times 3 * 3 \times 3 \times 3 = 4 \times 4 \times 2$$

(ch.) features  $\leftarrow$  (no. of filters)

... case for padding = 0 & stride = 1

• Build one layer of CNN

$$6 \times 6 \times 3 * 3 \times 3 \times 3 \rightarrow \text{RELU}(4 \times 4 + b_1) \rightarrow 4 \times 4$$

$$\underbrace{\text{input}}_{a^{[0]}} * \underbrace{\text{filter}}_{w^{[0]}} \rightarrow \underbrace{\text{RELU}(4 \times 4 + b_2)}_{g(z^{[0]})} \rightarrow 4 \times 4$$

$$\underbrace{\text{input}}_{a^{[0]}} * \underbrace{\text{filter}}_{w^{[1]}} \rightarrow \underbrace{\text{RELU}(4 \times 4 + b_3)}_{g(z^{[1]})} \rightarrow 4 \times 4$$

REVIEW

forward propagation:

$$z^{[i]} = w^{[i]} a^{[0]} + b^{[i]}$$

x

$$a^{[i]} = g(z^{[i]})$$

Q. 10 layers filters that are  $3 \times 3 \times 3$   
 10 one layer of a neural network.  
 $\text{no. of parameters} = 10 \times 28 = 280$

↑ filters      ↑ parameters

$3 \times 3 \times 3 = 27$   
 $\downarrow +1 \text{ (bias)}$   
28

\* layer  $l$  is convolutional layer:  
 $f^{[l]}$  = filter size —  $f \times f$  in layer  $l$   
 $p^{[l]}$  = padding ;  $n_c^{[l]}$  = no. of filters  
 $s^{[l]}$  = stride

Input:  $n_h^{[l-1]} \times n_w^{[l-1]} \times n_c^{[l-1]}$  ← from prev. layer

Output:  $n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]} \dots \left[ \frac{n_h^{[l-1]} \cdot n_w^{[l-1]} \cdot n_c^{[l-1]}}{s^{[l]} + 2p - f^{[l]} + 1} \right]$

Filter:  $f^{[l]} \times f^{[l]} \times n_c^{[l-1]}$   
 (weights)

Activations:  $a^{[l]} \rightarrow n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

$A^{[l]} \rightarrow m \times n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

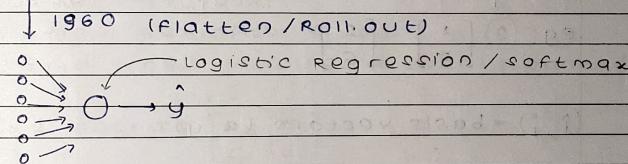
Bias:  $n_c^{[l]} \cdot (1, 1, 1, \dots, n_c^{[l]})$

### Example of convnet

$$\begin{array}{lll} 39 \times 39 \times 3 & \xrightarrow{\quad} & 37 \times 37 \times 10 \\ n_h^{[0]} = n_w^{[0]} = 39 & f^{[0]} = 3 & n_h^{[1]} = n_w^{[1]} = 37 \\ s^{[0]} = 1 & n_c^{[0]} = 3 & n_c^{[1]} = 10 \\ n_c^{[0]} = 3 & p^{[0]} = 0 & \end{array}$$

$$\begin{array}{lll} \text{filters} = 10 & f^{[2]} = 5 & \\ s^{[2]} = 2 & p^{[2]} = 0 & \end{array}$$

$$\begin{array}{lll} 7 \times 7 \times 40 & \xleftarrow{\quad} & 17 \times 17 \times 20 \xleftarrow{\quad} 20 \text{ filters} \\ n_h^{[2]} = 7 & f^{[2]} = 5 & n_h^{[3]} = n_w^{[3]} = 17 \\ n_w^{[3]} = 7 & s^{[3]} = 2 & n_c^{[2]} = 20 \\ n_c^{[3]} = 40 & 40 \text{ filters} & \end{array}$$



- Layers:
- 1) convolution
  - 2) Pooling
  - 3) Fully connected

- Pooling layer (down sampling) [FIXED FUNCTION]
  - 1) Reduce size of convolution
  - 2) Speed up computation
  - 3) Make system robust
- \* MAX POOLING → usually doesn't use padding

1	3	2	1		→	9	2
2	9	1	1	<small>max pool</small>		6	9
1	3	9	3				
5	6	1	2	(2x2)			

I/P.

Hyperparameters:  $f_c = 2 \leftarrow$  filter size  
 $s = 2 \leftarrow$  stride

- ⇒ feature detected in any quadrant
- high number
- not detected
- low number

\* Has hyperparameters but no parameters to learn (fixed computation, gradient descent doesn't change anything)  
 (No parameters that backprop. will adapt through output size =  $\frac{n+2p-f+1}{s} + 1$  max pooling)

3-D:  $5 \times 5 \times n_c \xrightarrow{\text{M.P.}} 3 \times 3 \times n_c$   
 max pooling done independently on each of the  $n_c$  channels

#### \* AVERAGE POOLING

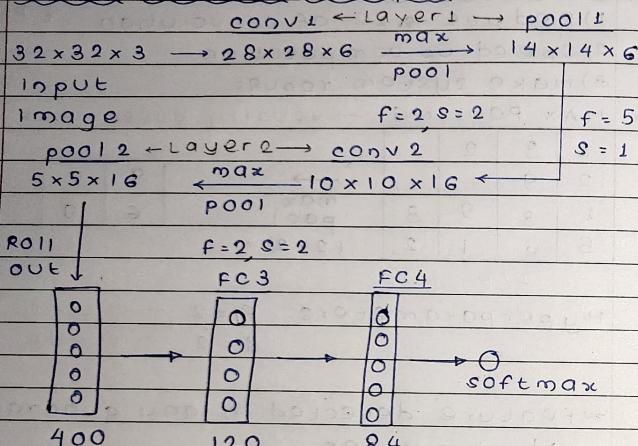
I/p	average pooling	→	3.75	1.25	
			4	2	

parameters =  $(9+1) \times 5 = 50$

one filter  $(3 \times 3) = 9$  no. of filters =  Page No.  Date

Bias = 1 per filter

### Neural network example



As layers progress:  $n_H, n_W \downarrow$  &  $n_C \uparrow$

→ conv - pool - conv - pool - FC - FC - FC - softmax

	activation shape	activation size	# parameters
Input	$(32, 32, 3)$	3072	0
conv1 ( $f=5, s=1$ )	$(28, 28, 6)$	4704	
pool1	$(14, 14, 6)$	1176	0
conv2 ( $f=5, s=1$ )	$(10, 10, 16)$	1600	
pool2	$(5, 5, 16)$	400	0
FC3	$(120, 1)$	120	
FC4	$(84, 1)$	84	
softmax	$(10, 1)$	10	

### why convolutions?

Advantages of using convolutions over just fully connected layers:

- 1) Parameter sharing (feature detector used)
- 2) Sparsity of connections (each layer  $\rightarrow$  each output depends only on a small number of inputs)

\* training set:  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$

$$\text{cost } J = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

→ use gradient descent to optimize parameters to reduce  $J$ .