

Course 4 : Convolutional Neural Network

Content :

Week 1 : Foundations of CNNs

Week 2 : Deep convolutional models : case studies

Week 3 : Object detection

Week 4 : Special applications : Face recognition ;

Neural style transfer

Course 2

Deep Learning Specialization



Vertical detector

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

Horizontal

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

Sobel filter

→ Taking care of middle row

$$\begin{matrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{matrix}$$

Scharr filter

$$\begin{matrix} 3 & 0 & -3 \\ 10 & 0 & -10 \\ 3 & 0 & -3 \end{matrix}$$

Padding

Imagine size: $n \times n$; convolution size: $f \times f$

Padding size: p

Output size after convolution:

without Padding: $(n-f+1) \times (n-f+1)$

Padding: $(n+2p-f+1) \times (n+2p-f+1)$

Convention:

$$p = (f-1)/2$$

valid convolution / no padding

Same convolution: output size = input size

f usually odd

Strided Convolution

Stride : s

matrix $n \times n$, $f \times f$ filter
padding P , stride s

output size after convolution :

$$L \lfloor (n+2P-f)/s + 1 \rfloor * L \lfloor (n+2P-f)/s + 1 \rfloor$$

1 layer of a Convolutional Network

filter $f^{[l]}$, Padding $P^{[l]}$, stride $s^{[l]}$

filters $n_c^{[l]}$

filter shape $f^{[l]} * f^{[l]} * n_c^{[l-1]}$

input shape $n_h^{[l-1]} * n_w^{[l-1]} * n_c^{[l-1]}$

Output shape $n_h^{[l]} * n_w^{[l]} * n_c^{[l]}$

height $n_h^{[l]} = \lfloor \frac{n_h^{[l-1]} + 2P^{[l]} - f^{[l]}}{s^{[l]}} + 1 \rfloor$

n_w

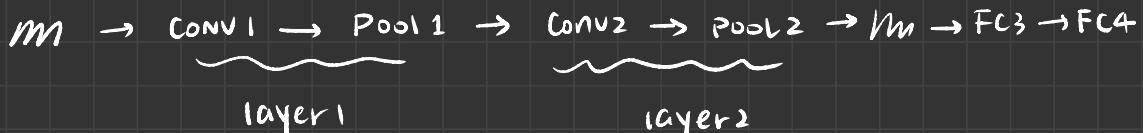
activation $a^{[l]}$ $n_h^{[l]} * n_w^{[l]} * n_c^{[l]}$

$A^{[l]} = m * n_h^{[l]} * n_w^{[l]} + n_c^{[l]}$

weights $f^{[l]} * f^{[l]} * n_c^{[l-1]} * n_c^{[l]}$

bias $(1, 1, 1, n_c^{[l]})$

CNN



	act. shape	act. size	# Par.
Input	(32, 32, 3)	3072	0
CONV1 ($f=5, s=1$)	(28, 28, 8)	6272	$(5*5*3+1)*8$
Pool1	(14, 14, 8)	1568	0
Conv2 ($f=5, s=1$)	(10, 10, 16)	1600	$(5*5*8+1)*16$
Pool2	(5, 5, 16)	400	0
FC3	(120, 1)	120	$400*120 + 120$
FC4	(84, 1)	84	$120*84 + 84$
softmax	(10, 1)	10	$84*10 + 10$

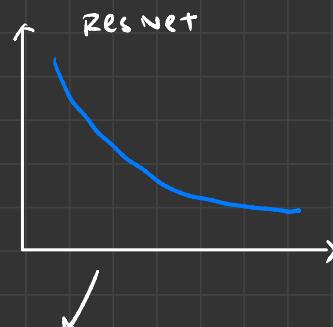
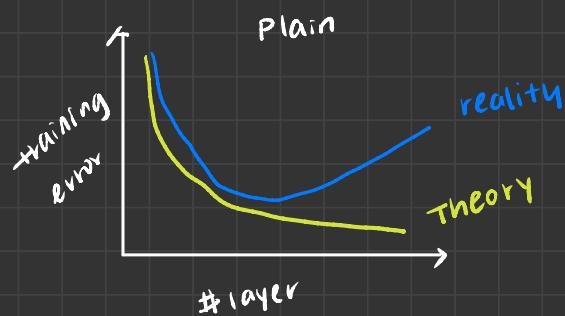
Classic network

} LeNet - 5
AlexNet
VGG

Lecture
for details

ResNet

Inception



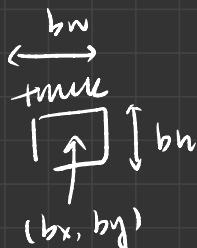
Performance T as network goes deeper

Detection algorithm

Image classification

1 object vs. multi object

Classification with localization



1. pedestrian , 2 car . 3 moto . 4 background

$$y = (P_c, bx, by, bw, bh, c_1, c_2, c_3)^T$$

Bounding Box Predictions (YOLO)

Intersection over Union (IoU)

= size of inter. area / size of union area

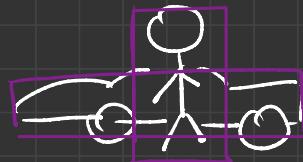
if IoU > 0.5, correct.

Non max Suppression:

Each output prediction is $(P_c \ b_x \ b_y \ b_n \ b_w)^T$

Discard all boxes with $P_c \leq 0.6$

Anchor Boxes



Anchor box 1



Anchor box 2



$$y = (\underbrace{P_c \ b_x \ b_y \ b_n \ b_w \ c_1 \ c_2 \ c_3}_{11} \ \underbrace{P_c \ b_x \ b_y \ b_n \ b_w \ c_1 \ c_2 \ c_3}_{10})^T$$

$$(1 \ 0 \ 0 \ 1 \ 0 \ 0 \ 1 \ b_x \ b_y \ b_n \ b_w \ 0 \ 1 \ 0)^T$$

$$\text{Car only? } (0 \ ? \ ? \ ? \ ? \ ? \ ? \ 1 \ b_x \ b_y \ b_n \ b_w \ 0 \ 1 \ 0)^T$$

R - CNN

Propose regions. Classify proposed regions one at a time.

time

Run "segmentation algorithm"

slow

Fast R-CNN

Propose regions. Use convolution implementation of sliding windows to classify all the proposed regions.

clustering step slow

Faster R-CNN

Use convolutional network to propose regions

still a bit slower than YOLO.

Triplet loss



dis. bit Anchor imagine if a positive or negative image

different
↑ person

same person