

C N N

- Course4 -

1. 컴퓨터 비전의 발전이 AI 플랫폼의 발전
2. 창의적, 다른 곳에 적용 가능

$$64 \times 64 \times 3 (\text{RGB}) \rightarrow 12288$$

$$1000 \times 1000 \times 3 (\text{RGB}) \rightarrow 3\text{million}$$

## Convolutional Neural Network

$$3x1+0 + 1x-1 + 1x1 + 5x0 + 8x-1 \dots$$

3	0	1
1	5	8
2	1	2

$6 \times 6$

CONVOLUTION

$\downarrow$

$*$

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

$3 \times 3$

filter 2 부름

$=$

-5	-4	0	8
-10	-2	2	3
0	-2	-4	7
-3	-2	-3	-6

$4 \times 4$ .

Python : conv-forward  
tensorflow :

## Vertical edge detection

1	0
0	

$6 \times 6$

$*$

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

$3 \times 3$

$=$

0	30	30	0
0	"	"	0
0	"	"	0
0	"	"	0

$4 \times 4$ .

이미지에서 수직 경계선을 찾을 수  
있게 된다.

$$\begin{bmatrix} 0 & 10 \end{bmatrix} * \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} = \begin{bmatrix} 0 & -30 & \cdot \end{bmatrix}$$

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

vertical

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

horizontal

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

sharp filter

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

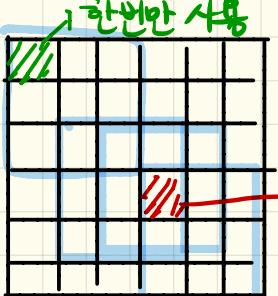
sharp filter

서로 유포인 경계를  
유지해, 90° 회전하면  
가로 유포인 경계

## Padding.

$$6 \times 6 \times 3 \times 3 = 4 \times 4 \quad (n+2p-f+1) \times (n+2p-f+1)$$

$$(n \times n) \quad (f \times f) \quad (n-f+1) \times (n-f+1)$$



- shrinkly output
- throw away into from edge

Padding = 1      0을 넣어줌

Valid:  $n \times n * f \times f \rightarrow n-f+1 \times n-f+1$

Same: pad so that output size is the same as the input size

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

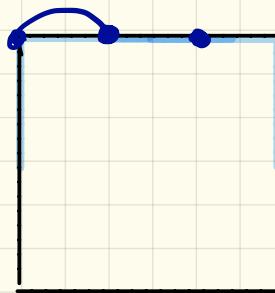
$$P = \frac{f-1}{2}$$

**filter**

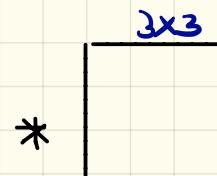
$f = \text{필수}$

- ① 짝수면 padding 비대칭
- ② 중앙 위치가 존재

Strided convolution

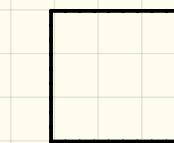


$n \times n * f \times f$



Stride = 2

=

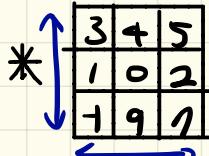
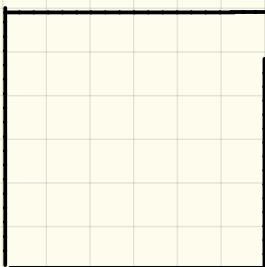


$3 \times 3$

$$\left(\frac{n+2p-f}{s}+1\right) \times \left(\frac{n+2p-f}{s}+1\right)$$

$\lceil z \rceil = \text{floor}(z)$  가장 작은 정수  
계산 결과는 내림 정수.

Technical note on cross-correlation vs convolution.

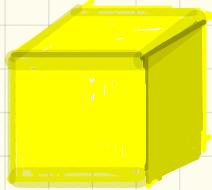
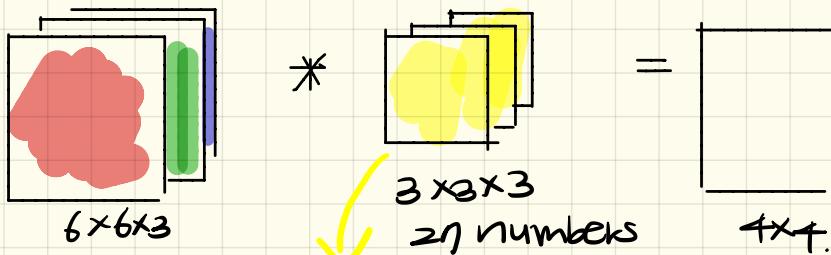


→ 가로의 서로 다른 텁텁기

$$\begin{matrix} 1 & 2 & 5 \\ 9 & 0 & 4 \\ -1 & 1 & 3 \end{matrix}$$

합성곱 기준 이진곱셈  
그리고 시그모이드 계층

# Convolutions on RGB images.



$$R \quad \begin{matrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{matrix}$$

$$G \quad \begin{matrix} 0 \end{matrix}$$

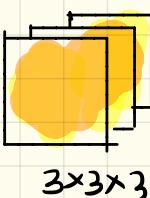
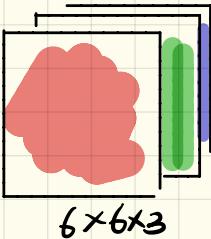
$$B \quad \begin{matrix} 0 \end{matrix} \rightarrow 3 \times 3 \times 3$$

빨간색의  
서브스텝을

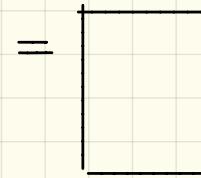
$$R \quad \begin{matrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{matrix}$$

$$G \quad \begin{matrix} 1 & -1 \\ 1 & -1 \\ 1 & -1 \end{matrix}$$

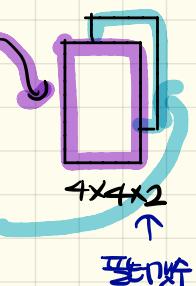
$$B \quad \begin{matrix} 1 & -1 \\ 1 & -1 \\ 1 & -1 \end{matrix} \rightarrow 모든 색  
3 \times 3$$



$3 \times 3 \times 3$



$4 \times 4$



4x4x1

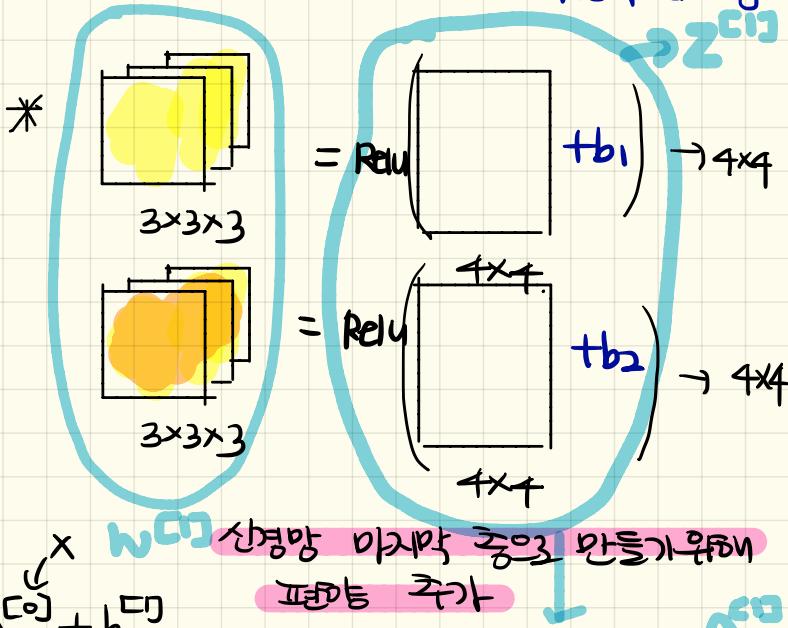
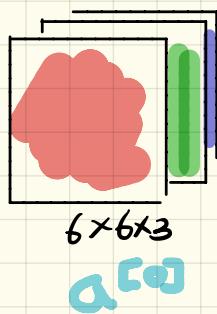
$$\text{Summary} = n \times n \times n_c * f \times f \times n_c$$

$$6 \times 6 \times 3 \quad 3 \times 3 \times 3$$

$$\rightarrow n-f+1 \times n-f+1 \times n_c'$$

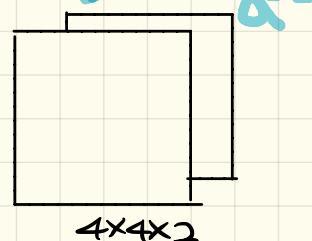
$$4 \times 4 \times 2$$

$\hookrightarrow$  # filters  
no stride  
no padding



$$z^{[c_1]} = w^{[c_1]} \alpha^{[c_0]} + b^{[c_1]}$$

$$\alpha^{[c_0]} = g[z^{[c_1]}]$$



$3 \times 3 \times 3$

27 Para  
+ bias  
 $\rightarrow$  28 Para

...

10

$\therefore$  280 parameters

If layer  $l$  is a convolutional layer

$f^{[l]}$  = filter size

$P^{[l]}$  = Padding

$S^{[l]}$  = stride

$n_c^{[l]}$  = number of filters

Input:  $n_H^{[l-1]} \times n_w^{[l-1]} \times n_c^{[l-1]}$

Output:  $n_H^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

$$n^{[l]} = \left[ \frac{n_H^{[l-1]} + 2P^{[l]} - f^{[l]}}{S^{[l]}} + 1 \right]$$

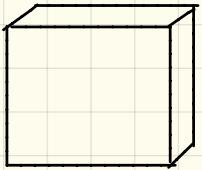
Each filter is:  $f^{[l]} \times f^{[l]} \times n_c^{[l]}$

Activation:  $a^{[l]} \rightarrow n_H^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

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

weights:  $f^{[l]} \times f^{[l]} \times f^{[l]} \times n_c^{[l]}$

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



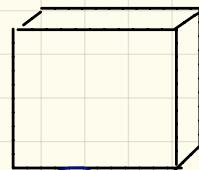
$39 \times 39 \times 3$

$$n_H^{[0]} = n_w^{[0]} = 39$$

$$n_C^{[0]} = 3$$

$$\begin{aligned} f^{[1]} &= 3 \\ s^{[1]} &= 1 \\ p^{[1]} &= 0 \end{aligned}$$

10 filters

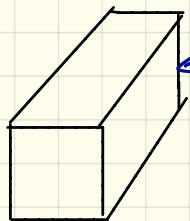


$37 \times 37 \times 10$

$$n_H^{[1]} = n_w^{[1]} = 37$$

$$n_C^{[1]} = 10$$

$$\frac{n + 2p - f}{s} + 1$$



$11 \times 11 \times 20$

$$\begin{aligned} f^{[2]} &= 5 \\ s^{[2]} &= 2 \\ p^{[2]} &= 0 \end{aligned}$$

20 filters.

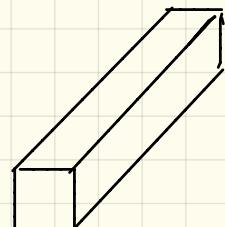
stride를 줄여서  
넓이를 줄였다.

$$\begin{aligned} n_H^{[2]} &= n_w^{[2]} = 11 \\ n_C^{[2]} &= 20 \end{aligned}$$

$\alpha^{[2]}$

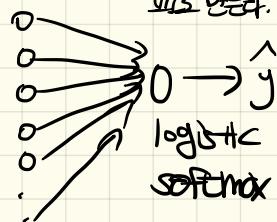
$$\begin{aligned} f^{[3]} &= 5 \\ s^{[3]} &= 2 \end{aligned}$$

to filters



$1 \times 1 \times 40$

$= 1960$   
풀자서  
196개의  
값을 만든다.



Types of layer in convolutional network:

- Convolution (CONV)
- Pooling (POOL)
- Fully connected (FC)

## Max Pooling

1	3	2	1
2	9	1	1
1	3	2	3
5	6	1	2

$4 \times 4$



9	2
6	3

$$\begin{aligned} &2 \times 2 \\ &f=2 \\ &s=2 \end{aligned}$$




9	9	5
9	9	5
8	6	9

$$\begin{aligned} &f=3 \\ &s=1 \end{aligned}$$

$$\left( \frac{n-2p+f}{s} + 1 \right)$$

## Average Pooling

1	3	2	1
2	9	1	1
1	3	2	3
5	6	1	2

$4 \times 4$



3.75	1.25
4	2

$$\begin{aligned} &2 \times 2 \\ &f=2 \\ &s=2 \end{aligned}$$

신경망에서는 최대풀링을 평균풀링 보다 많이 사용한다.

# Summary of Pooling

Hyper Parameters :

F: filter size

$$f=2, S=2$$

S: stride

$$F=2, S=2$$

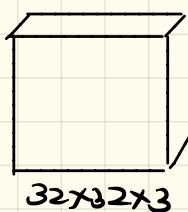
Max or Average pooling

~~P: Padding~~

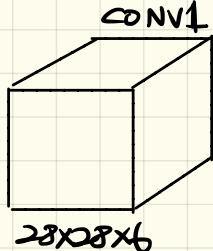
padding은 사용하지 않는다.  
일반적인 값은 0.

$$n_H \times n_W \times n_C \rightarrow \left[ \frac{n_H - f}{s} + 1 \right] \times \left[ \frac{n_W - f}{s} + 1 \right] \times n_C$$

## Neural Network example (LeNet-5)

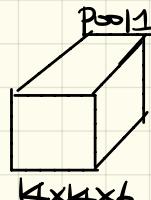


$$\begin{matrix} f=5 \\ S=1 \end{matrix}$$

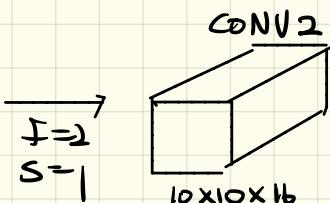


CONV1

$$\begin{matrix} \text{maxPool} \\ f=2 \\ S=2 \end{matrix}$$



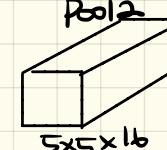
Pool1



$$\begin{matrix} f=2 \\ S=1 \end{matrix}$$

CONV2

$$\begin{matrix} \text{maxPool} \\ f=2 \\ S=2 \end{matrix}$$



Pool2

$$\begin{matrix} b[3]=120 \\ W[3] (120, 40) \\ 120 \end{matrix} \rightarrow \begin{matrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{matrix}$$

Layer 2

# Neural Network example

	Activation shape	Activation size	# parameters
Input	(32, 32, 3)	$32 \times 32 \times 3$	0
CONV1(f=5, s=1)	(28, 28, 8)	6272	208
Pool L1	(14, 14, 8)	1568	0
CONV2(f=5, s=1)	(10, 10, 16)	1600	416
Pool L2	(5, 5, 16)	400	0
FC3	(120, 1)	120	12000 } 대체로의 연산은
FC4	(84, 1)	84	10081 } FC 층 에 있으므로
Softmax	(10, 1)	10	84

## why convolution ?

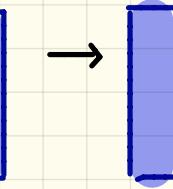
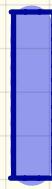
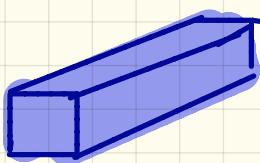
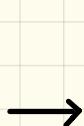
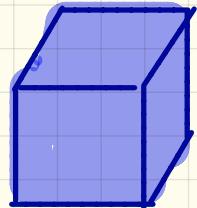
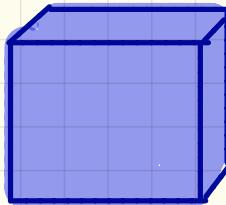
① FC로 하면 변수가 너무 많아진다.

Parameter sharing : A feature detector that's useful in one part of the image is probably useful in another part of the image.

② sparsity of connections : In each layer, each output value depends on a small number of inputs.

# Putting it together

Training set  $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$



운전여행  
두어는 소프트웨어  
출력값인  $y$ 의  
예측값

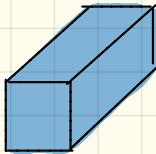
$$\text{Cost } J = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(g^{(i)}, y^{(i)})$$

use gradient descent to optimize parameters  
to reduce  $J$

# LeNet - 5



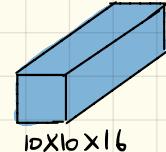
$\xrightarrow{5 \times 5}$   
 $S=1$



avg pool  
 $f=2$   
 $S=2$

$14 \times 14 \times 6$

$\xrightarrow{5 \times 5}$   
 $S=1$



$10 \times 10 \times 16$

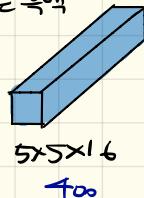
$32 \times 32 \times 1$

LeNet은 흑백

avg pool

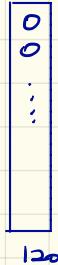
$\xrightarrow{f=2}$

$S=2$

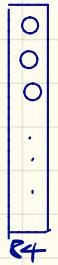


$\xrightarrow{FC}$

$5 \times 5 \times 16$   
400



120

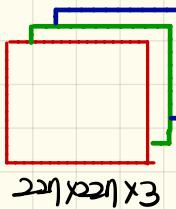


84

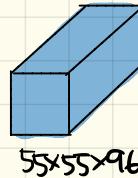
$\xrightarrow{\text{Softmax}}$

60k parameters  
 $N_{in}, N_{out} \downarrow N_c \uparrow$

## AlexNet



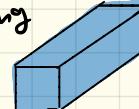
$\xrightarrow{11 \times 11}$   
 $S=4$



$55 \times 55 \times 96$

Max Pooling

$\xrightarrow{3 \times 3}$   
 $S=2$

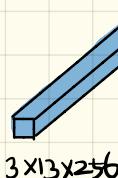


$27 \times 27 \times 96$

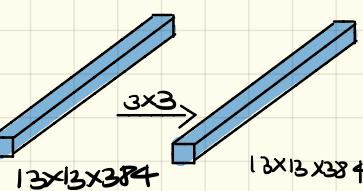
$\xrightarrow{5 \times 5}$   
some

$227 \times 227 \times 3$

MAXpooling  
 $\xrightarrow{3 \times 3}$   
 $S=2$



$\xrightarrow{3 \times 3}$   
Same



$13 \times 13 \times 256$

$\xrightarrow{3 \times 3}$   
 $S=2$

$\xrightarrow{3 \times 3}$   
 $S=2$

=

$9216$

$4096$

$4096$

$4096$

Soft  
max  
1000

- Similar to LeNet, but much bigger

- ReLU  $\frac{2}{3}$  AF

- Multiple GPUs

- Local Response normalization (LRN)

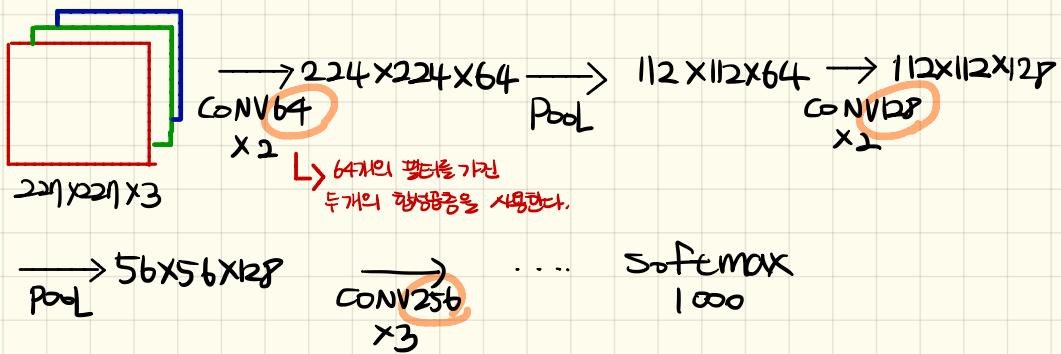
Less important

# VGG-16

→ 16개의 가중치를 가진 층이 있다.

$\text{CONV} = 3 \times 3 \text{ filter}, s=1, \text{Same}$

$\text{Maxpool} = 2 \times 2, s=2$



- 꼬리나 균일하고 간결하다.

## Residual block

$$a^{[l]} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \end{bmatrix} \rightarrow a^{[l+2]}$$

"main path"

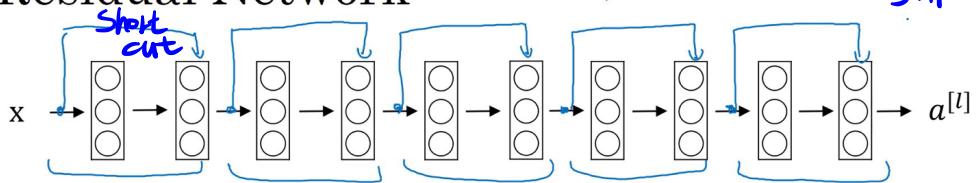
skip connection / "short cut"

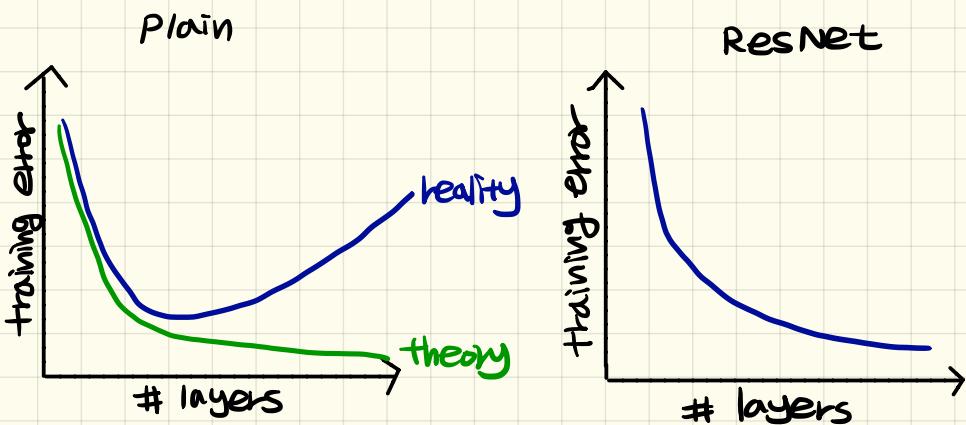
$$a^{[l]} \xrightarrow{\text{Linear}} \text{ReLU} \xrightarrow{a^{[l+1]}} \text{Linear} \xrightarrow{a^{[l+2]}}$$

$$\text{ReLU} \xrightarrow{a^{[l+2]}}$$

$$z^{[l+1]} = w^{[l+1]} a^{[l]} + b^{[l+1]} \quad a^{[l+1]} = g[z^{[l+1]}] \quad z^{[l+2]} = w^{[l+2]} a^{[l+1]} + b^{[l+2]} \quad a^{[l+2]} = g[z^{[l+2]}]$$

## Residual Network





## why do residual networks work?

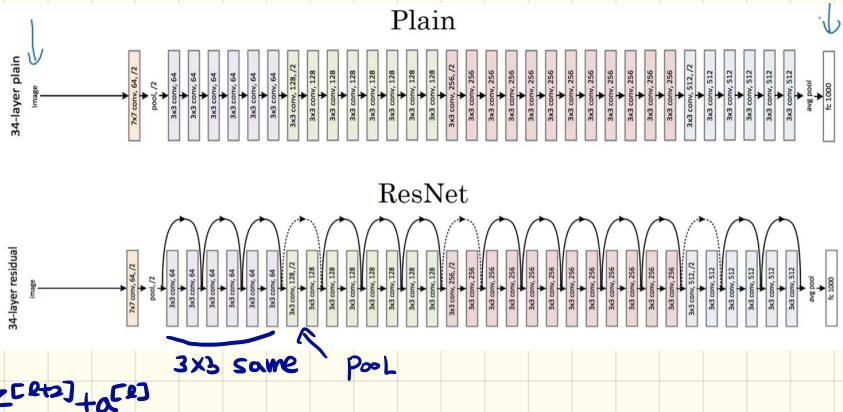
$X \rightarrow \text{Big NN} \xrightarrow{\alpha^{[l]}} \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} \rightarrow \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} \rightarrow \alpha^{[l+2]}$  "Same"

ReLU  $\alpha \geq 0$

If  $\alpha^{[l+2]} = g(z^{[l+2]} + \alpha^{[l]})$   $z^{[l+2]} \text{ et } \alpha^{[l]} \text{ 이 같은 차원을 가진다고 가정}$

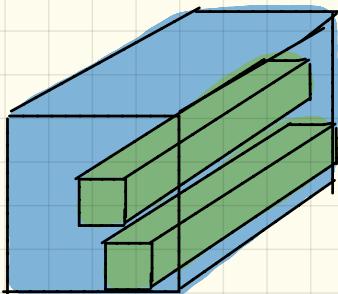
$$\begin{aligned} &= g(w^{[l+2]} \alpha^{[l+1]} + b^{[l+2]} + \underbrace{(w \alpha^{[l]})}_{\text{If } w^{[l+2]}=0 \text{ or } b^{[l+2]}=0 \text{ IR}^{1056 \times 1028} = g(\alpha^{[l]})} \\ &= \alpha^{[l]} \end{aligned}$$

어느곳에 ResNet을 추가해도 성능 저하가 없다.  
경사화 방법으로 더 나아질 뿐이다.



Why do a  $1 \times 1$  Convolution do?

$$\begin{array}{|c|c|} \hline & 1 & 2 & 3 \\ \hline \end{array} \quad * \quad \boxed{2} \quad = \quad \begin{array}{|c|c|} \hline & 2 & 4 & 6 \\ \hline \end{array}$$

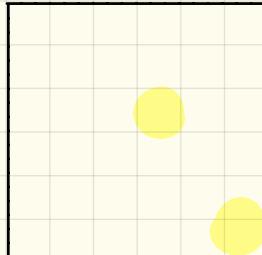


$$6 \times 6 \times 32$$

$1 \times 32 \times 32$

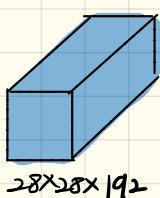
32 → #filters

$$n_e^{[e+1]}$$

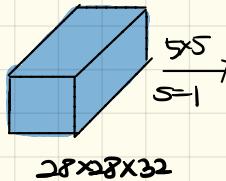


$6 \times 6 \times \# \text{ filters}$

# Using 1x1 convolutions

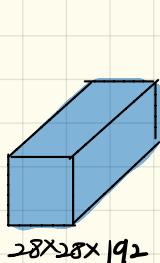


ReLU  
→  
CONV1X1  
32  
 $1 \times 1 \times 192$   
32 filters

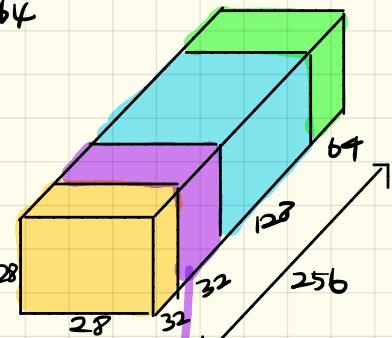


- 네트워크가 아인형성을 더하고 채널 수를 조절할 수 있게 된다.

## Motivation for inception network.

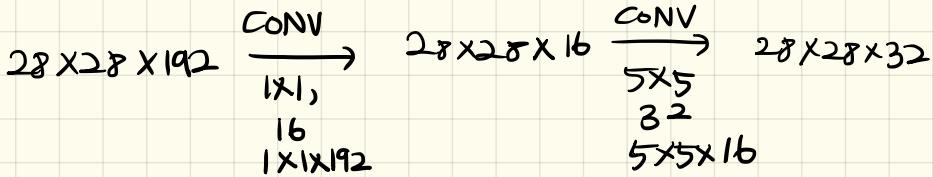


→  $1 \times 1$   $28 \times 28 \times 64$   
→  $3 \times 3$   $28 \times 28 \times 128$   
Same  
→  $5 \times 5$   $28 \times 28 \times 32$   
Same  
→ MaxPool  $28 \times 28 \times 32$   
Same  $s=1$



32-filters filters are  $5 \times 5 \times 192$

$$28 \times 28 \times 32 \times 5 \times 5 \times 192 = 120m$$



비용

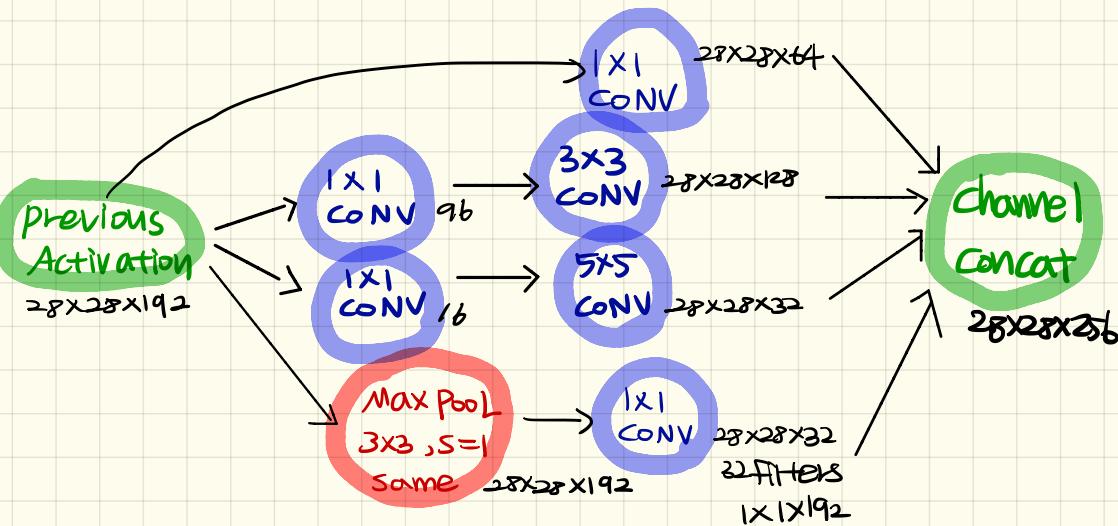
$$28 \times 28 \times 192 \times 192 = 2.4m$$

$$28 \times 28 \times 32 \times 5 \times 5 \times 16 = 10m$$

12.4m

필터  $1 \times 1, 3 \times 3, 5 \times 5$ , 팔링층인가 고민하고 살피면  
인셉션 모델은 그걸 다 엮은 것이다.

## Inception module



# Transfer Learning

- 기존 신경망을 다운로드 받을 때 가중치도 받아야 한다.
- 마지막 소프트맥스층만 수정.
- trainable parameter = 0, freeze=1로 두면 훈련층은 그대로 유지된다.
- 데이터가 아주 많으면 네트워크를 다시 훈련.

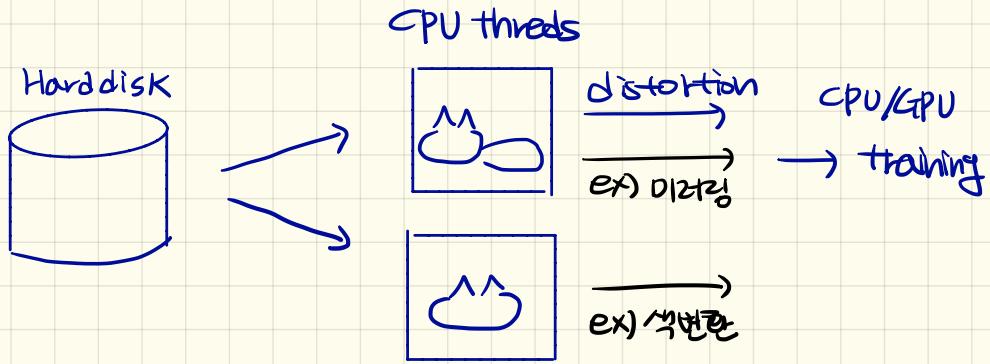
## Common argumentation method

- mirroring
- Random Cropping
  - Rotation
  - shearing
  - Local wraping

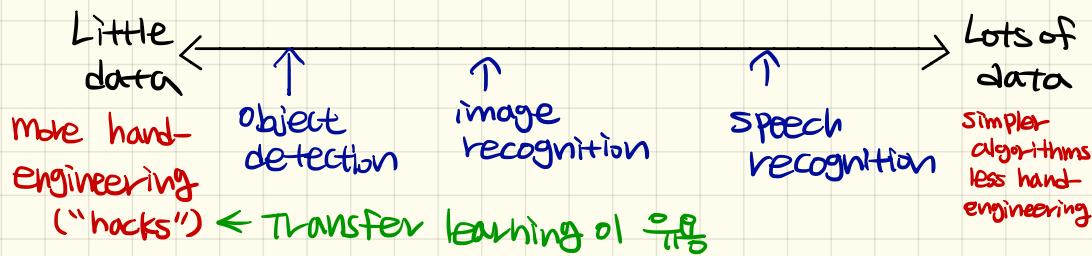
## Color shifting

- RGB 조절
- PCA color argumentation.

# Implementing distortions during training



# Data vs hand-engineering



## Two Sources of Knowledge

- Labeled data ( $x, y$ )
  - Hand engineered features/network architecture/other components.

레이어별 데이터간 충분하지 않다면 즉각 실행이 의존해야 함

# 번지마킹을 잘 하는 팁

## Ensembling

- 어떤 선경망을 사용할지 결정하고, 몇개의 선경망을 독립적으로 훈련시킨 뒤 평균을 내는 것.  
성능이 좋아 디리에스 쓰이지만 서비스용으로 쓰이지 않음.  
~~많은 양의 메모리가 필요.~~

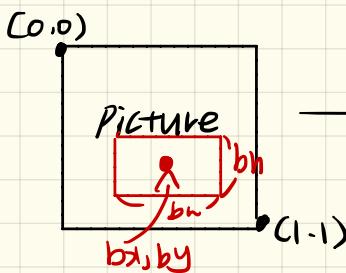
## Multi-crop at test time

- Run classifier on multiple versions of test images and average results.

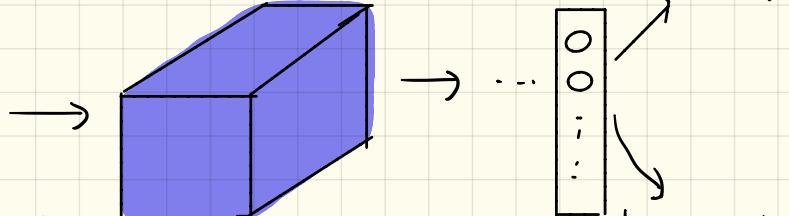
이것도 실제 제품 보다는 번지마킹에 사용된다.

네트워크 하나만으로 유저에게 있어서 메모리를 많이 차지  
하지 않지만 여전히 실행 시간을 낮춘다.

## classification with localization



1. Pedestrian
2. car
3. motorcycle
4. background



$$\begin{aligned}bx &= 0.5 \\by &= 0.7 \\bh &= 0.3 \\bw &= 0.4\end{aligned}$$

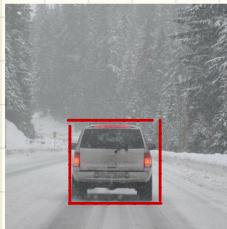
bbox, bwh, bws, by

Need to output  $b_x, b_y, b_h, b_w \rightarrow$  class label (1-4)

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

is there any object?

$x =$



$$\begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

1. Pedestrian
2. Car
3. motorcycle
4. background

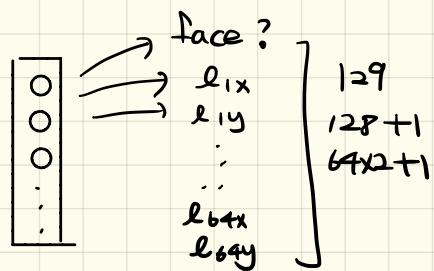
$$L(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_8 - y_8)^2 & \text{if } y=1 \\ (\hat{y}_1 - y_1)^2 & \text{if } y_1=0 \end{cases}$$

## Land mark detection.

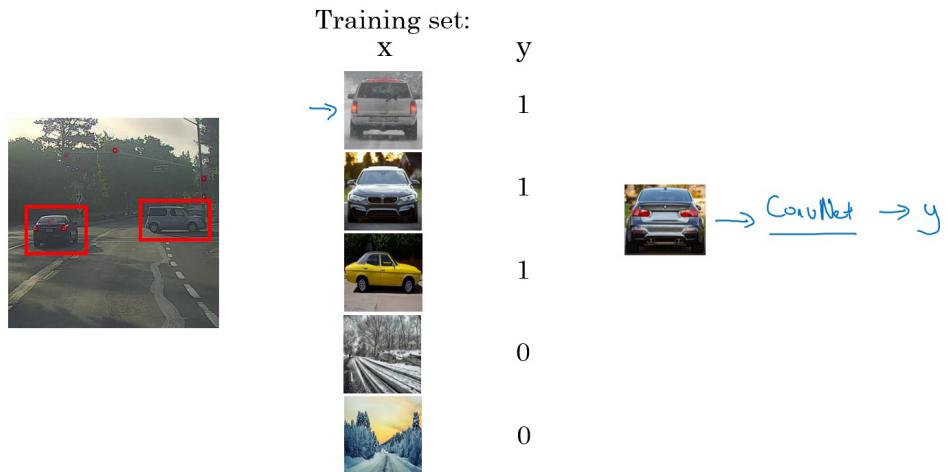
만약 얼굴이 64개의 랜드마크가 있다면?



$\rightarrow$  CONVNet  $\rightarrow$

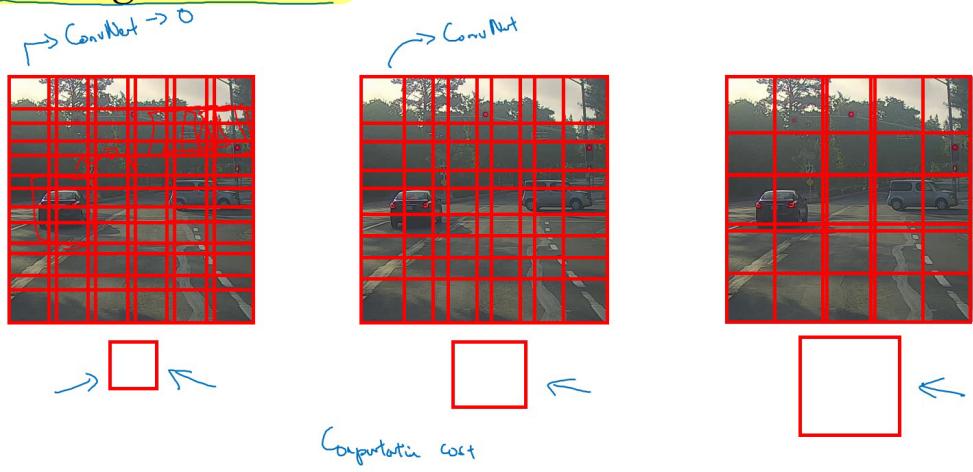


# Car detection example



Andrew Ng

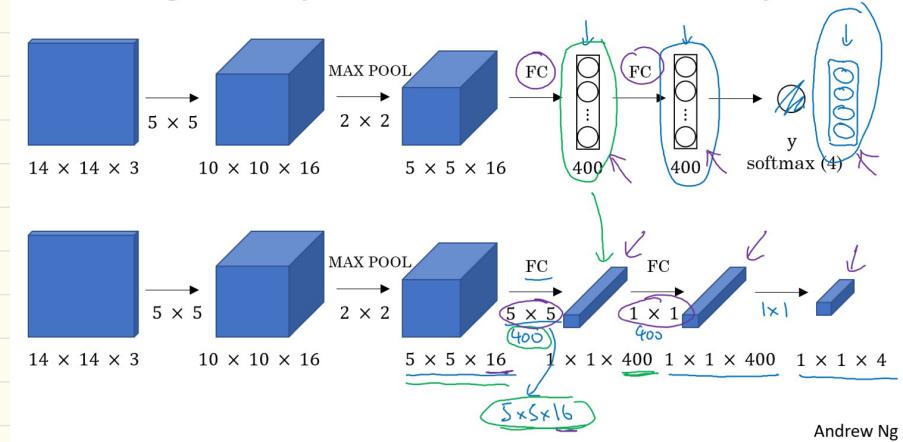
## Sliding windows detection



Andrew Ng

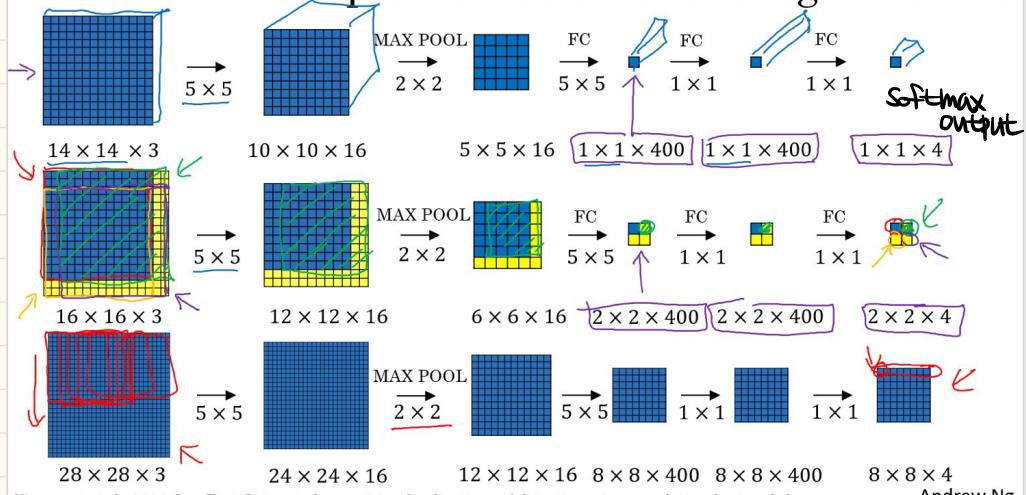
계산비용이 많음, 아주 큰 스트라이드는 정확도↓ 아주 작은 스트라이드는 아주 큰 계산비용.

# Turning FC layer into convolutional layers



Andrew Ng

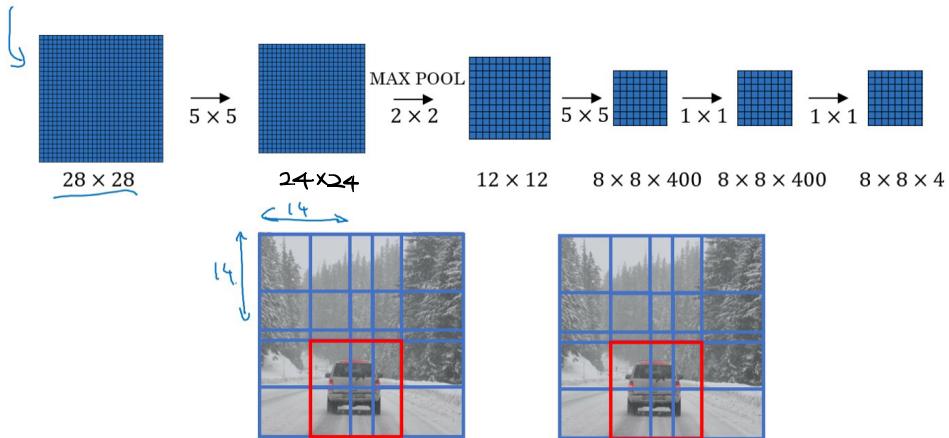
# Convolution implementation of sliding windows



[Sermanet et al., 2014, OverFeat: Integrated recognition, localization and detection using convolutional networks]

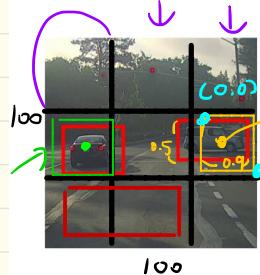
Andrew Ng

# Convolution implementation of sliding windows



Andrew Ng

가장 큰 단점은 바운딩 박스의 위치가 그리 정확하지  
않다는 것.



Target output

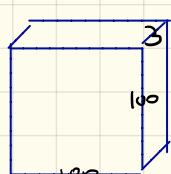
$3 \times 3 \times 8$

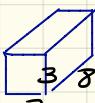
Labels for training  
For each grid cell:

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

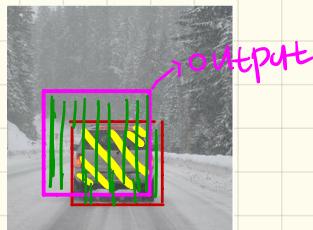
0	?	?	1	b <sub>x</sub>	b <sub>y</sub>	b <sub>h</sub>	b <sub>w</sub>
?	?	?	b <sub>x</sub>	b <sub>y</sub>	b <sub>h</sub>	b <sub>w</sub>	0.4
?	?	?	b <sub>y</sub>	b <sub>h</sub>	b <sub>w</sub>	0.3	0.3
?	?	?	b <sub>h</sub>	b <sub>w</sub>	0.9	0.9	0.9
?	?	?	b <sub>w</sub>	0.5	0.5	0.5	0.5
0	0	0	0	0	0	0	0

between 0,1 ) could be >1



$\rightarrow$  CONV  $\rightarrow$  MaxPool  $\rightarrow$  

# Evaluating object localization



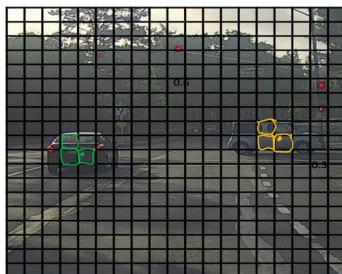
Intersection over Union (IoU)

Size of  $\frac{\text{Size of } \text{Yellow}}{\text{Size of } \text{Green}}$

Size of  $\frac{\text{Size of } \text{Green}}{\text{Size of } \text{Yellow}}$

IoU↑  $\rightarrow$  정답과 같은 박스

## Non-max suppression example



19x19

Each output prediction is:

$$\begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \end{bmatrix}$$

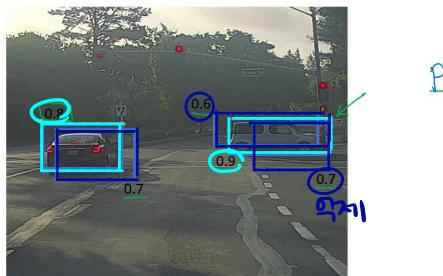
Discard all boxes  
with  $P_c \leq 0.6$ .

While there are any  
remaining boxes:

- 가장 높은  $P_c$ 를 가지고 있는  
상자들을 반복적으로 제거해 예측
- 가장 높은 IoU 제외하고 넘기기

Andrew Ng

## Non-max suppression example

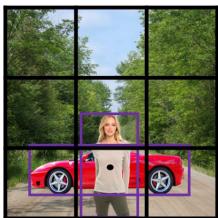


Non-Max Suppression은 여러번 템포를 깨끗하게 정리하는 것

Andrew Ng

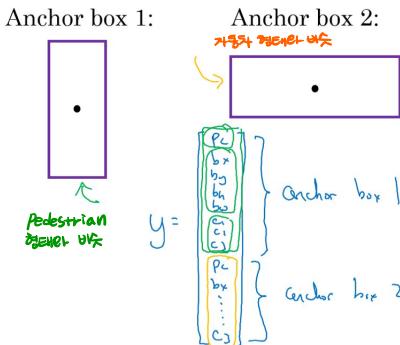
## Anchor Boxes

## Overlapping objects:



$$y = \{ p_c, b_x, b_y, b_h, b_w, c_1, c_2, c_3 \}$$

[Redmon et al., 2015, You Only Look Once: Unified real-time object detection]



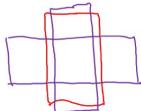
Andrew Ng

## Anchor box algorithm

Previously:

Each object in training image is assigned to grid cell that contains that object's midpoint.

Output y:  
3 x 3 x 8



With two anchor boxes:

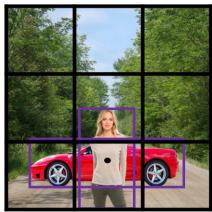
Each object in training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IoU.

Output y: box)

$$3 \times 3 \times \underline{16}$$

$$3 \times 3 \times 2 \times 8$$

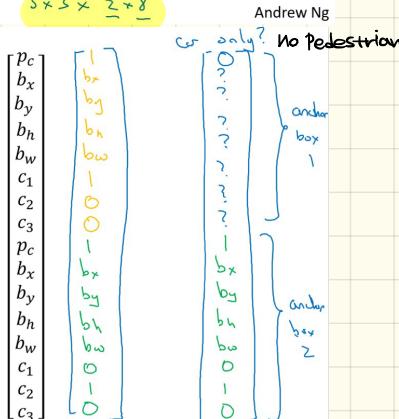
## Anchor box example



Anchor box 1: Anchor box 2:



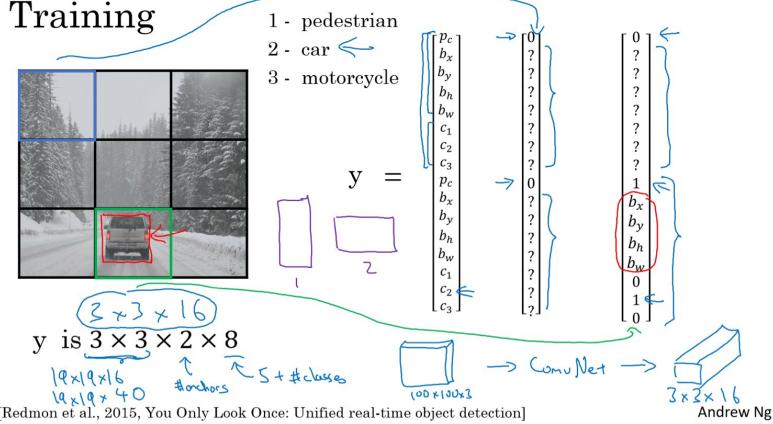
$$y =$$



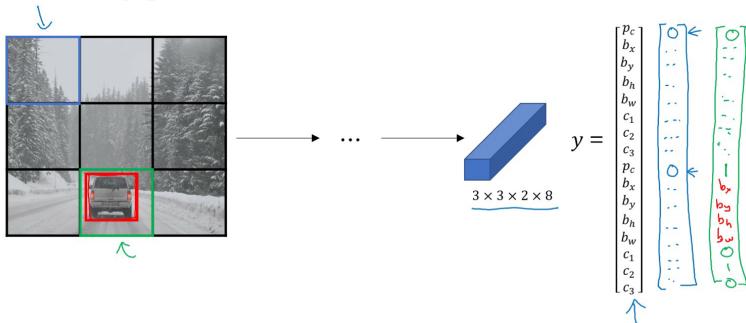
Andrew Ng

# YOLO Algorithm

## Training



## Making predictions



## Outputting the non-max suppressed outputs



- For each grid cell, get 2 predicted bounding boxes.
- Get rid of low probability predictions.
- For each class (pedestrian, car, motorcycle) use non-max suppression to generate final predictions.

Andrew Ng

# Face verification vs face recognition

→ Verification 실제로 잘 적용하기 위해서는 99% 이상이야함

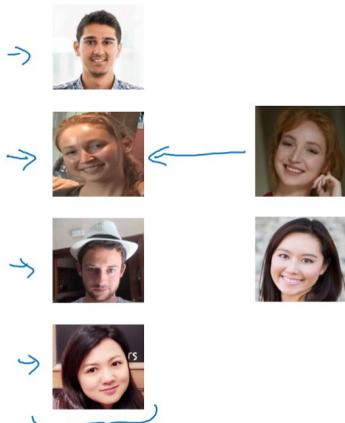
- Input image, name/ID
- Output whether the input image is that of the claimed person.

→ Recognition

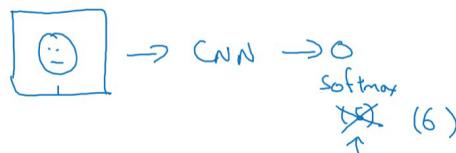
$$k=100$$

- Has a database of  $k$  persons.
- Get an input image.
- Output ID if the image is any of the  $k$  persons (or

## One-shot learning



Learning from one example to recognize the person again



training set이 1이랑만 작으면  
초연이 걸 안될것

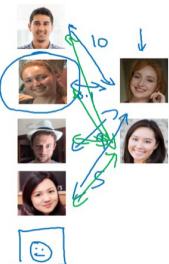
Andrew Ng

# Learning a “similarity” function

→  $d(\text{img1}, \text{img2})$  = degree of difference between images

If  $d(\text{img1}, \text{img2}) \leq \tau$       "same"  
 $> \tau$       "different"

} Verification.

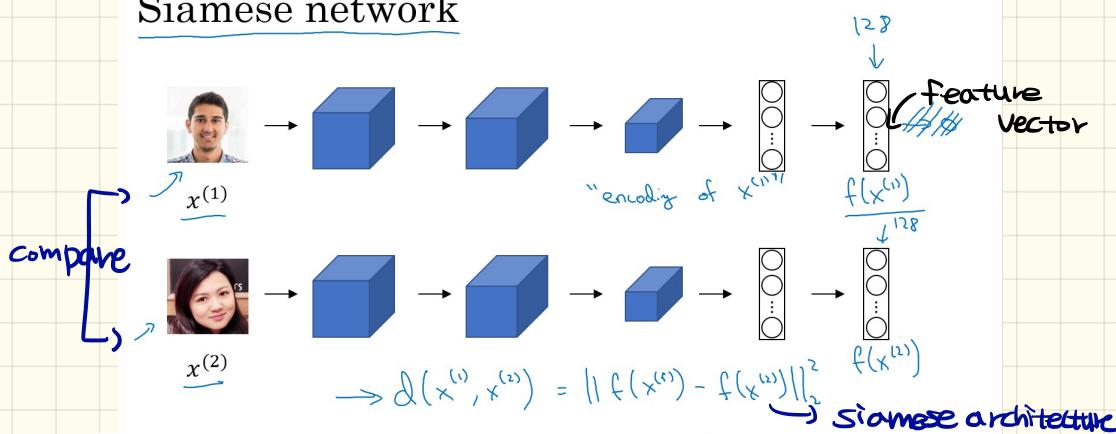


$$d(\text{img1}, \text{img2})$$

두 사진이 같으면 같은수 출력, 두 사진이 다르면  
큰 수 출력하게 하여라

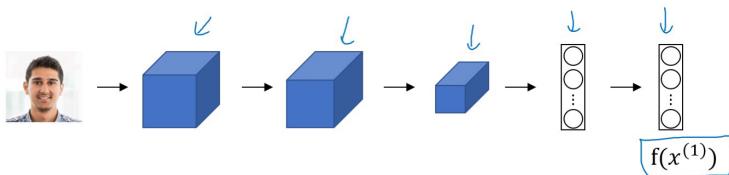
Andrew Ng

## Siamese network



## Goal of learning

Andrew Ng



Parameters of NN define an encoding  $f(x^{(i)})$

Learn parameters so that:

If  $x^{(i)}, x^{(j)}$  are the same person,  $\|f(x^{(i)}) - f(x^{(j)})\|^2$  is small.

If  $x^{(i)}, x^{(j)}$  are different persons,  $\|f(x^{(i)}) - f(x^{(j)})\|^2$  is large.

Andrew Ng

# Learning Objective



same person



different person

Anchor

$$\frac{d(A, P)}{d(A, P) + \alpha} = 0.5$$

Want:  $\frac{\|f(A) - f(P)\|^2}{d(A, P)} + \alpha \leq 0.2$

Positive

Anchor

$$\frac{d(A, N)}{d(A, N) + \alpha} = 0.5$$

$\frac{\|f(A) - f(N)\|^2}{d(A, N)} + \alpha \leq 0.7$

Negative

$$\frac{\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha}{\text{margin}} \leq 0 \quad \text{if } f(\text{img}) = \vec{0}$$

[Schroff et al., 2015, FaceNet: A unified embedding for face recognition and clustering]

Andrew Ng

## Loss function

Given 3 images  $A, P, N$ :

$$L(A, P, N) = \max(\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha, 0) = 0$$

$$J = \sum_{i=1}^m L(A^{(i)}, P^{(i)}, N^{(i)}) > 0$$

$A, P, N$  같은 사건이 일어나는 =

training set: lot pictures of 1k persons.

(사람 당 사건 카운트는 충분히 높기 때문이다)

## Choosing the triplets $A, P, N$

During training, if  $A, P, N$  are chosen randomly,  $d(A, P) + \alpha \leq d(A, N)$  is easily satisfied.

$\sum_i \|f(A) - f(P)\|^2 + \alpha \leq \|f(A) - f(N)\|^2$ , 한집합에서 선택되면 원래의 사건보다 확률 더 큰 학습률을 가질 가능성이 매우 높다.  
따라서 신경망으로부터 배울 수 있는 것의 많이 있다.

choose triplets that're "hard" to train on.

$$d(A, P) + \alpha \leq d(A, N)$$

$$\frac{d(A, P)}{\downarrow} \approx \frac{d(A, N)}{\uparrow}$$
 해서 둘 사이의  $\alpha$ 를 만들기 위해서

Training set using triplet loss

Anchor      Positive      Negative



J

$$d(x^{(i)}, x^{(j)})$$

Andrew Ng

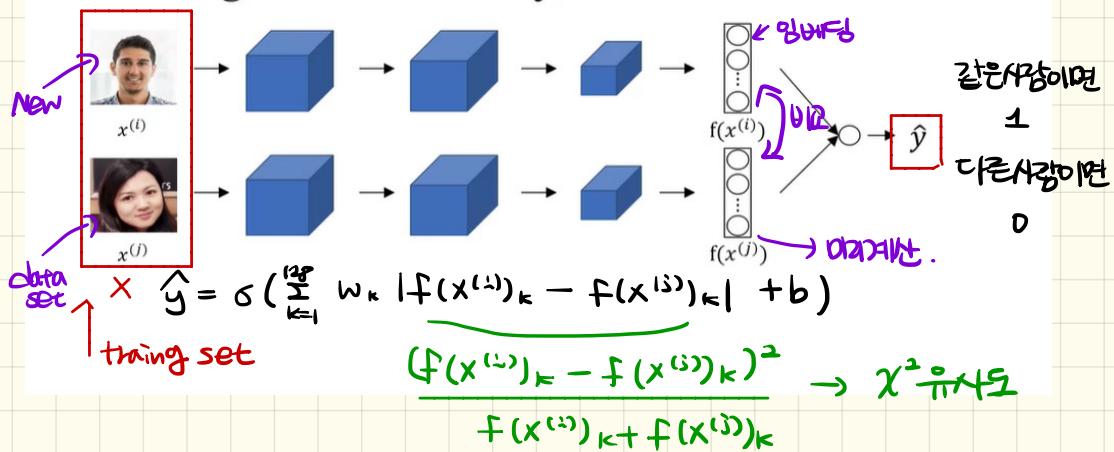
(J) 비용함수를 작게 만들기 위해 기울기 방향을 사용하라.

그리면 신경망의 모든 파라미터에 영향을 미치는 효과가 있다.

$$d(x^{(i)}, x^{(j)}) \downarrow \text{"same"} \quad d(x^{(i)}, x^{(j)}) \uparrow \text{"different"}$$

# Learning the similarity function

simianse network

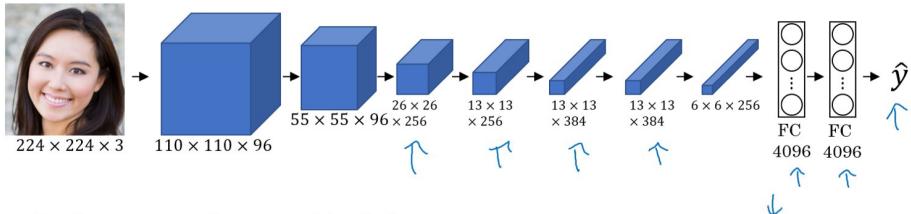


새로운 직원이 들어오면 위쪽 구성 요소를 사용하여 인코딩을 계산하고 사용하면 된다.

## Face verification supervised learning

$x$	$y$	
	1	"Same"
	0	"Different"
	0	
	1	

# ✓ Visualizing what a deep network is learning

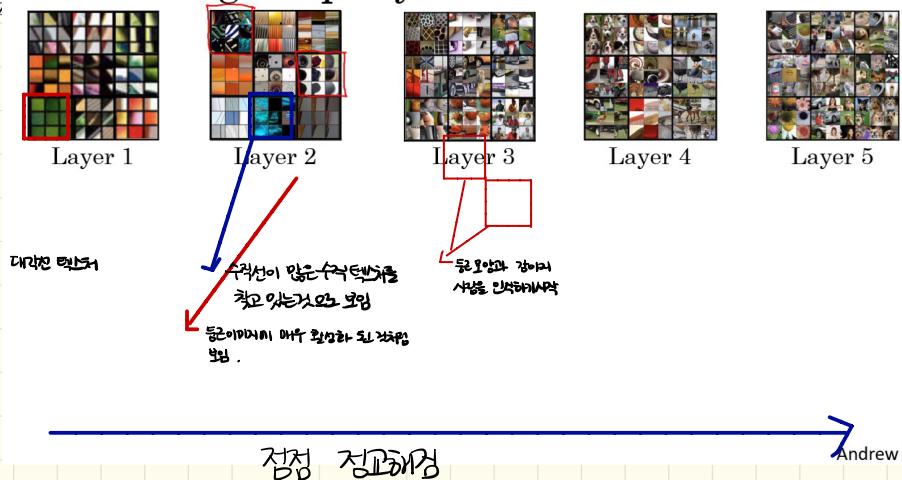


Pick a unit in layer 1. Find the nine image patches that maximize the unit's activation.



Repeat for other units.

## Visualizing deep layers



## Cost function.

$$J(G) = \alpha J_{content}(C, G) + \beta J_{style}(S, G)$$

Find the Generated image  $G$ .

1. Initiate  $G$  randomly  $G: 100 \times 100 \times 3$

2. USE gradient descent to minimize  $J(G)$

$$G := G - \frac{\partial}{\partial G} J(G)$$

# Content cost function

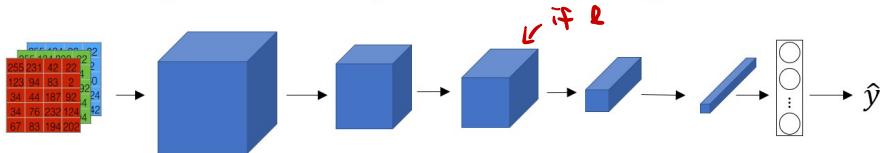
$$J(G) = \alpha J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$$
$$\hookrightarrow \frac{1}{2} \|a^{[l](C)} - a^{[l](G)}\|^2$$

- Say you use hidden layer  $l$  to compute content cost.

너무 많지도 적지도 않게 하든지이어 설정.

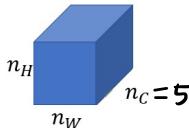
- USE pre-trained ConvNet (E.g. VGG network)
- Let  $a^{[l](C)}$  and  $a^{[l](G)}$  be the activation of layer  $l$  on the images.
- If  $a^{[l](C)}$  and  $a^{[l](G)}$  are similar, both images have similar content.

## Meaning of the “style” of an image



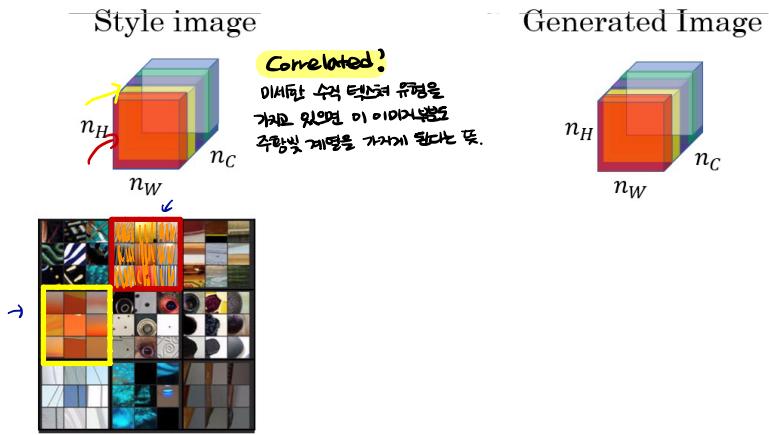
Say you are using layer  $l$ 's activation to measure “style.”

Define style as correlation between activations across channels.



How correlated are the activations across different channels?

# Intuition about style of an image



[Gatys et al., 2015. A neural algorithm of artistic style]

Andrew Ng

## Style matrix

Let  $a_{i,j,k}^{[l]}$  = activation at  $(i, j, k)$ .  $G^{[l]}$  is  $n_c^{[l]} \times n_c^{[l]}$

$$G_{kk'}^{[l](S)} = \sum_{i=1}^{n_h^{[l]}} \sum_{j=1}^{n_w^{[l]}} \alpha_{ikk}^{[l](S)} \alpha_{jjk'}^{[l](S)}$$

이미지의 서로 다른 위치에서 높이와 너비에 대한 합.  
채널  $k$ 와  $k'$ 의 activation을 합한다.

둘다 커지면  $G_{kk'}$ 도 증가  
두개나 상반되는 경우면  $G_{kk'}$ 는 감소

$$G_{kk'}^{[l](G)} = \sum_i \sum_j \alpha_{ikk}^{[l](G)} \alpha_{jjk'}^{[l](G)} \rightarrow \text{Generated (grand) matrix}$$

$$J_{\text{style}}^{[l]}(S, G) = \frac{\|G^{[l](S)} - G^{[l](G)}\|_F^2}{\sum_{k,k'} (G_{kk'}^{[l](S)} - G_{kk'}^{[l](G)})^2}$$

[Gatys et al., 2015. A neural algorithm of artistic style]

Andrew Ng

# Style cost function

$$J_{style}^{[l]}(S, G) = \frac{1}{\left(2n_H^{[l]} n_W^{[l]} n_C^{[l]}\right)^2} \sum_k \sum_{k'} (G_{kk'}^{[l](S)} - G_{kk'}^{[l](G)})^2$$

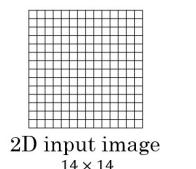
$$J_{style}(S, G) = \sum_l \lambda^{[l]} J_{style}^{[l]}(S, G)$$

$$J(G_T) = \alpha J_{content}(C, G) + \beta J_{style}(S, G)$$

[Gatys et al., 2015. A neural algorithm of artistic style]

Andrew Ng

## Convolutions in 2D and 1D

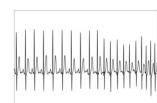


\*



$$14 \times 14 \times 3 * 5 \times 5 \times 3 \rightarrow 10 \times 10 \times 16$$

2D input image  
14 × 14



\*



$$14 \times 1 * 5 \times 1 \rightarrow 10 \times 16$$

14 (1D)

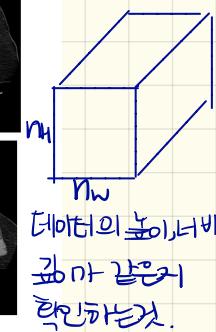
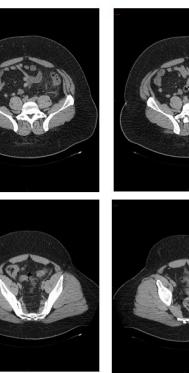
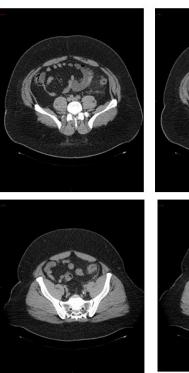
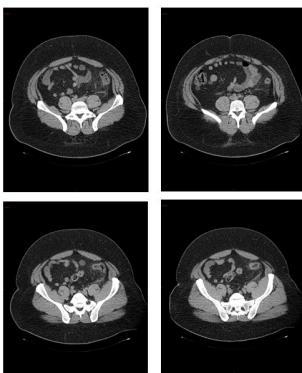
$$\begin{bmatrix} 1 & 2 & 1 & 3 & 1 & 8 & 1 & 2 & 1 & 1 & 1 \end{bmatrix}$$

5

$$10 \times 16 * 5 \times 1 \rightarrow 6 \times 32$$

Andrew Ng

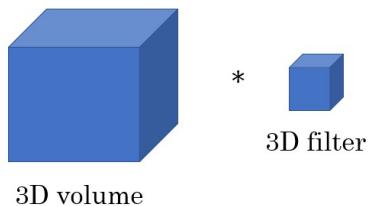
## 3D data



데이터의 높이 너비,  
깊이 갈고리  
혹은 핸드폰.

Andrew Ng

# 3D convolution



$14 \times 14 \times 14 \times 1$

$\ast 5 \times 5 \times 5 \times 1$  16 filters

$\rightarrow 10 \times 10 \times 10 \times 16$

$\ast 5 \times 5 \times 5 \times 16$  32 filters

$\rightarrow 6 \times 6 \times 6 \times 32$

Andrew Ng