

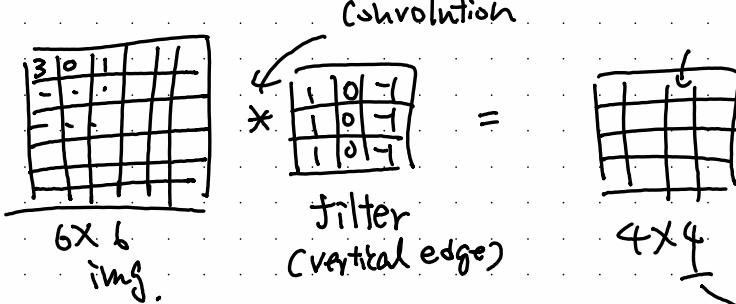
Course 4

Convolutional Neural Networks

CoVNet

Computer Vision

Edge Detection



Python: conv-forward

tensor : tf.nn.conv2d

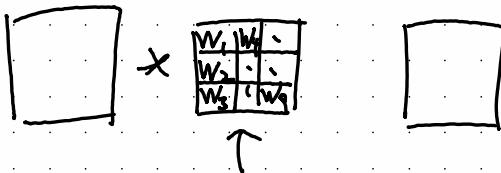
$n-j+1$

Sobel filter

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

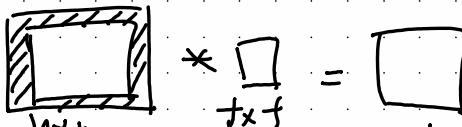
Scharr filter.

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



use NN to learn filter

Padding



valid - $n \times n$ $\times f \times f = n-f+1 \times n-f+1$

(no padding)

 $(n+p) \times (n+p) \times P_{fxf} = n-f+2p+1 \times n-f+2p+1$

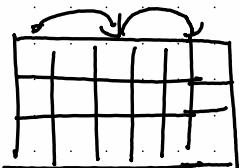
same
(padding)

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

e.g. $f=3 \rightarrow p=1$, $f=5 \rightarrow p=2$

f is almost always odd $\begin{matrix} 3 \times 3 \\ 5 \times 5 \\ \vdots \end{matrix}$

Strided convolution



$$n \times n \times f \times f$$

pad: p , $\Rightarrow L \left\lceil \frac{n+2p}{g} + 1 \right\rceil \times L \left\lceil \frac{n+2p+f}{g} + 1 \right\rceil$

stride: s

Convolution in math

$$\begin{matrix} 3 & 4 & 5 \\ 1 & 0 & 1 \\ 9 & 7 & 8 \end{matrix} \xrightarrow{\text{flip before conv.}} \begin{matrix} 8 & 7 & 9 \\ 1 & 0 & 1 \\ 5 & 4 & 3 \end{matrix}$$

$\xrightarrow{\text{flip}}$

cross-corr = no flip filter
conv = flip filter

ML community just call cross-corr as convolution

Conv in volume

$$\begin{matrix} 6 \times 6 \times 3 \\ \downarrow z \end{matrix} \times \begin{matrix} 3 \times 3 \times 3 \\ \downarrow z \end{matrix} = \begin{matrix} 4 \times 4 \\ \triangle \text{channels} \end{matrix}$$

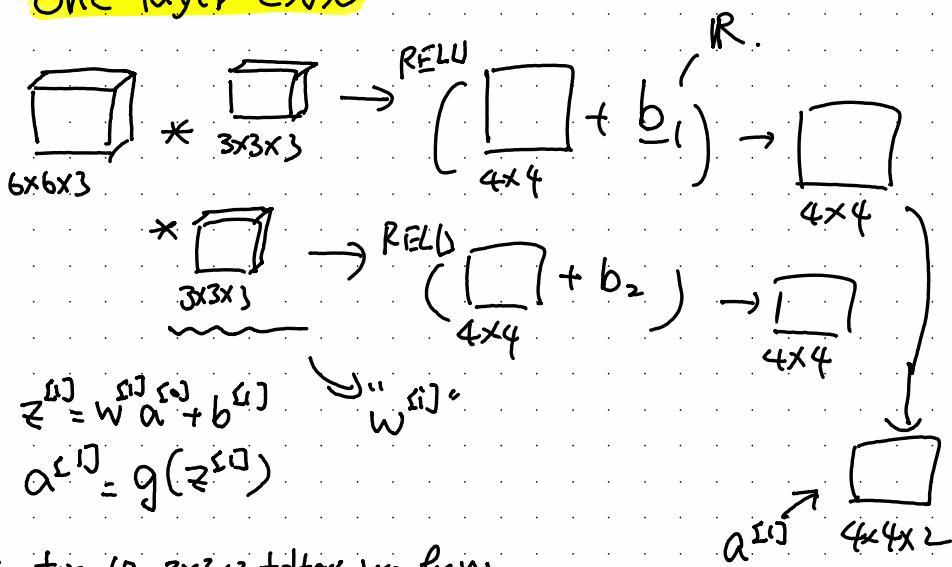
multiple filters

$$\begin{matrix} 6 \times 6 \times 1 \\ \times \end{matrix} \begin{matrix} 3 \times 3 \times 3 \\ \times \end{matrix} = \begin{matrix} 4 \times 4 \\ \triangle \text{horizontal} \end{matrix}$$

$$\begin{matrix} 3 \times 3 \times 3 \\ \times \end{matrix} = \begin{matrix} 4 \times 4 \\ \triangle \text{vertical} \end{matrix} \xrightarrow{\quad} \begin{matrix} 4 \times 4 \times 2 \end{matrix}$$

$$n \times n \times n_c \times f \times f \times n_c := n-f+1 \times n-f+1 \times n_c$$

One-layer CNN



e.g. for 10 $3 \times 3 \times 3$ filters, we have

$$(3 \cdot 3 \cdot 3 + 1) \times 10 = 280 \text{ parameters}$$

* parameters are independent from image resolution

Notation

$f^{[l]}$ - filter

$p^{[l]}$ - padding

$s^{[l]}$ - stride

$n_c^{[l]}$ - number of filters
each filter.

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

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

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

$$n_w \left\lfloor \frac{n_w^{[l-1]} + 2p^{[l-1]} - f^{[l-1]}}{s^{[l-1]}} + 1 \right\rfloor$$

Activation

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

weights

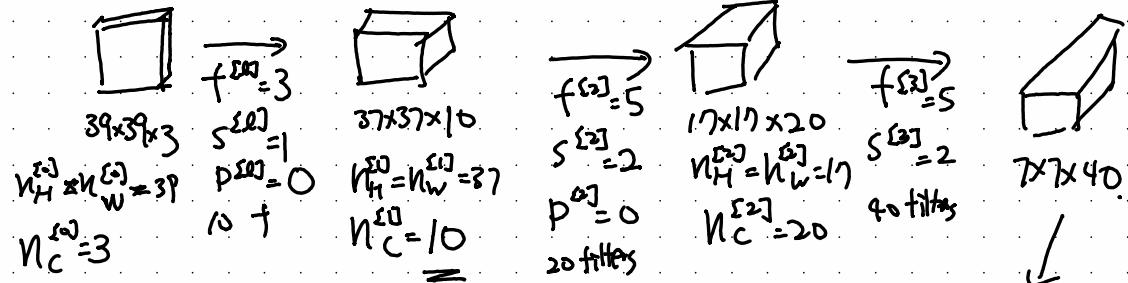
$$f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$$

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

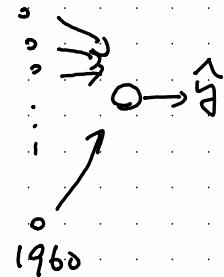
bias:

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

CNN



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

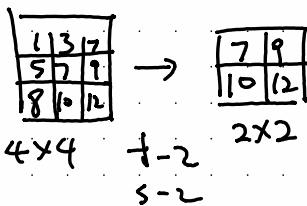


TYPE of LAYERS in CNN

- convolution (Conv)
- pooling (Pool)
- fully connected (FC)

POOLING LAYER

Max pooling - preserve feature?



△ no parameters to learn △

Average pooling - not often used

hyperparameters

f - filter size $f=2, s=2$

s - stride $f=3, s=2$

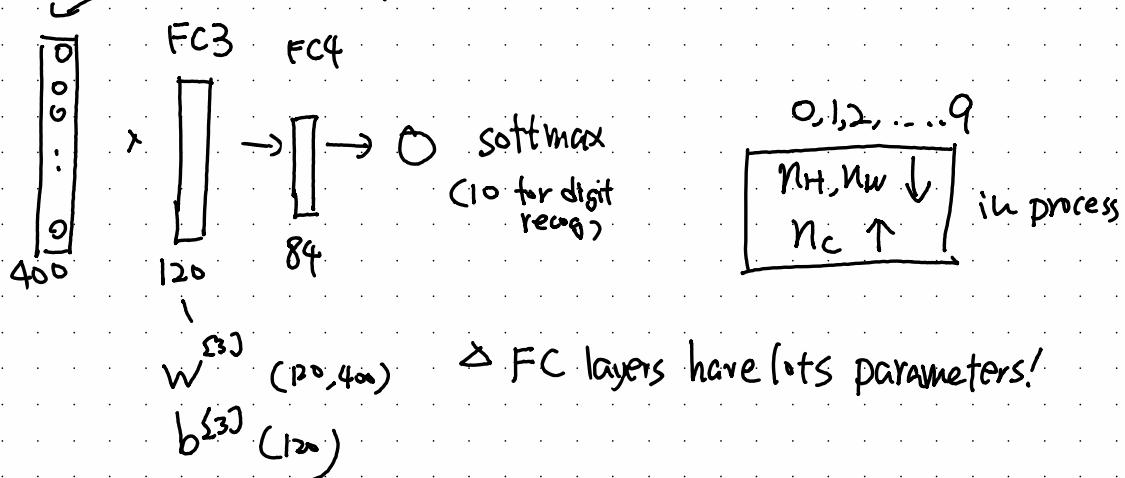
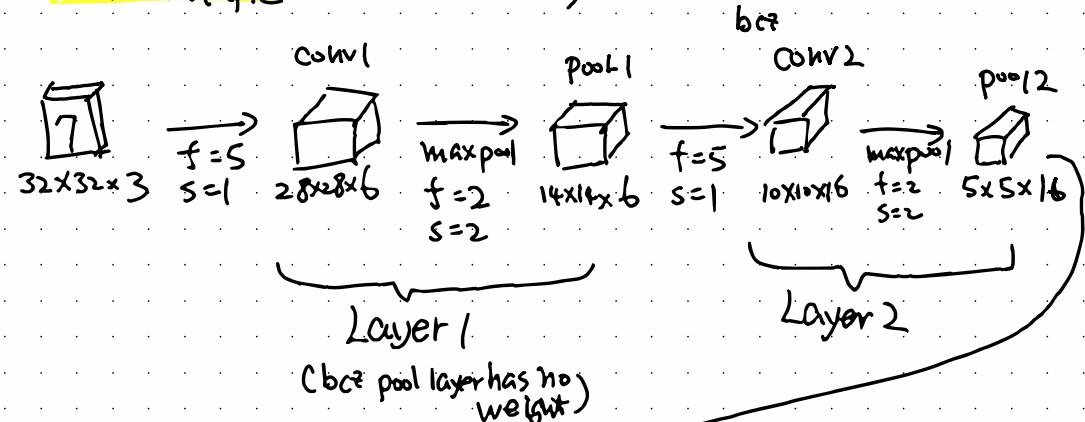
type of pooling layer

$$\left\lfloor \frac{n-f}{s} + 1 \right\rfloor$$

* No parameter to learn

NN example

(LeNet - 5)



* Why Conv?

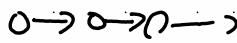
Reduce #parameters compared

- parameter sharing - feature detector useful for all parts of image
- sparsity of connections - each output value depends on small part of image

\triangle translation invariant.
 (pooling layer)

Keras

- sequential API

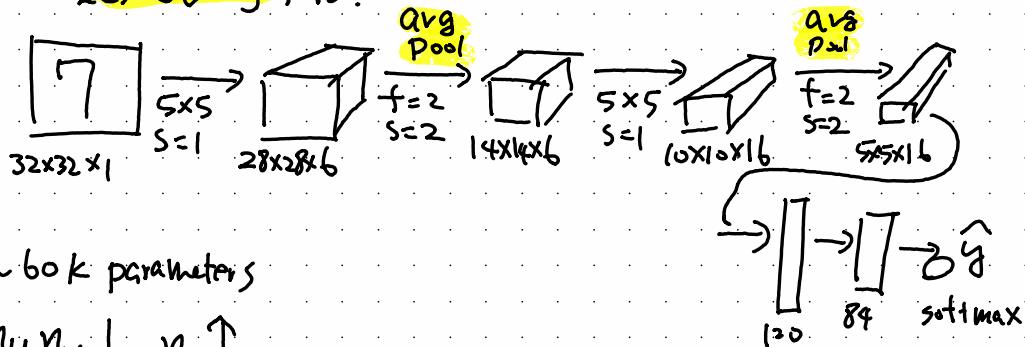


- functional API



CASE STUDY

LeNet-5, 98.



~60k parameters

$n_H n_W \downarrow n_c \uparrow$

conv pool / conv pool + fc pool

original

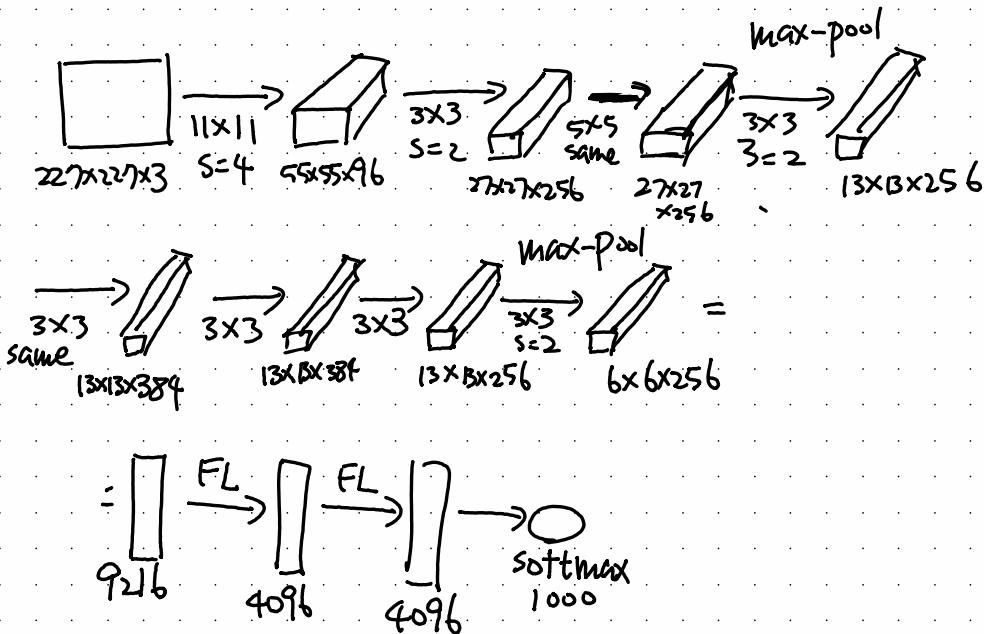
Sigmoid/tanh

different filter

nonlinearity after pooling

section II, III

AlexNet. 2012



$\sim 60M$ parameters

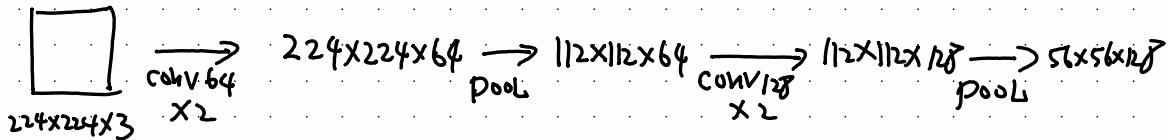
- ReLU

Orig

- multiple GPU
- local response normalization (LRN)
 - normalize along axis
 - less often used.

VGG-16, 2015

CONV = 3x3, S=1, same, MAX-Pool = 2x2, S=2



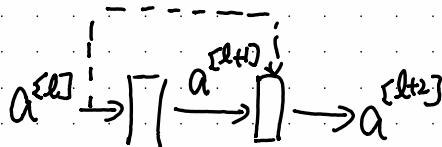
~138M parameters LARGE!

~Uniform

$\Delta w \downarrow n \downarrow c \uparrow$

Residual Network (ResNets) 2015

Residual Block.

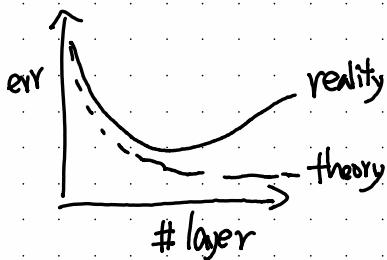


SHORT CUT / SKIP CONNECTION

* $a^{[l]}$ $\xrightarrow{\text{linear}} \text{ReLU} \xrightarrow{\text{linear}} \text{ReLU} \xrightarrow{\oplus} \text{ReLU} \rightarrow a^{[l+2]}$

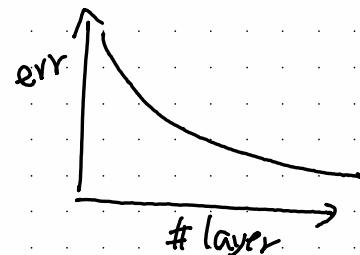
$$z^{[l+1]} = w^{[l+1]} a^{[l]} + b^{[l+1]}, a^{[l+1]} = g(z^{[l+1]}), z^{[l+2]} = w^{[l+2]} a^{[l+1]} + b^{[l+2]}, a^{[l+2]} = g(z^{[l+2]})$$

plain



$$\text{ResNet} \Rightarrow a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$$

Res



△ help training deeper network.

Why ResNet Work?



$$a^{l+2} = g(z^{l+2} + a^{l+1})$$

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

IF $\begin{matrix} \downarrow \\ 0 \end{matrix}$ $\begin{matrix} \downarrow \\ 0 \end{matrix} \Rightarrow a^{l+2} = g(a^{l+1})$

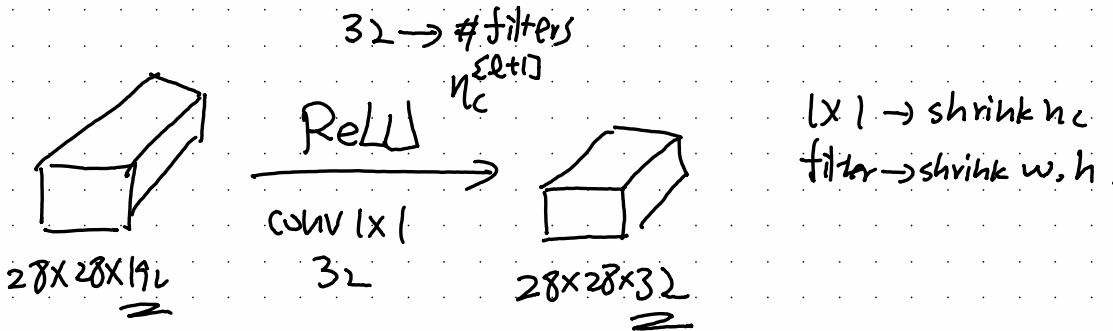
→ identity function is easy to learn

→ ① res block doesn't hurt.

→ ② prevent deeper layers making result worse.

1×1 convolution / network in network

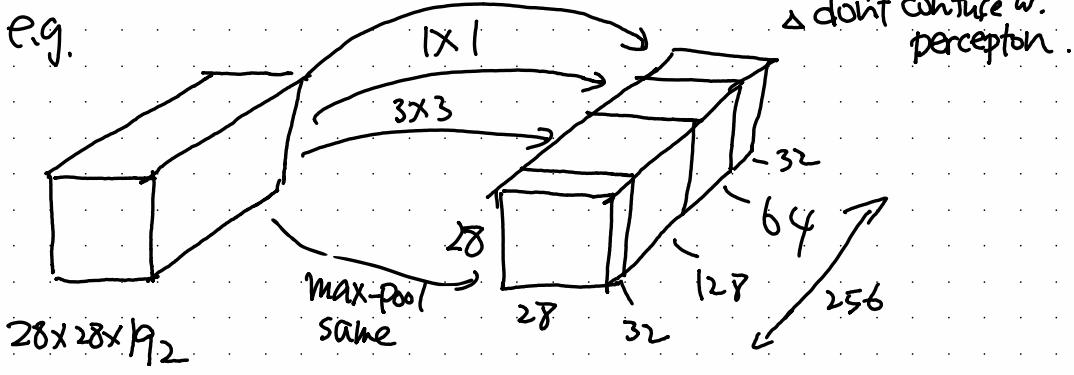
$$\begin{array}{c} \text{3D tensor} \\ 6 \times 6 \times 32 \end{array} * \begin{array}{c} \text{filter} \\ 1 \times 1 \times 32 \end{array} \xrightarrow{\text{ReLU}} = \begin{array}{c} \text{3D tensor} \\ 6 \times 6 \times \# \text{filters} \end{array}$$



shrink n_c , if we like.

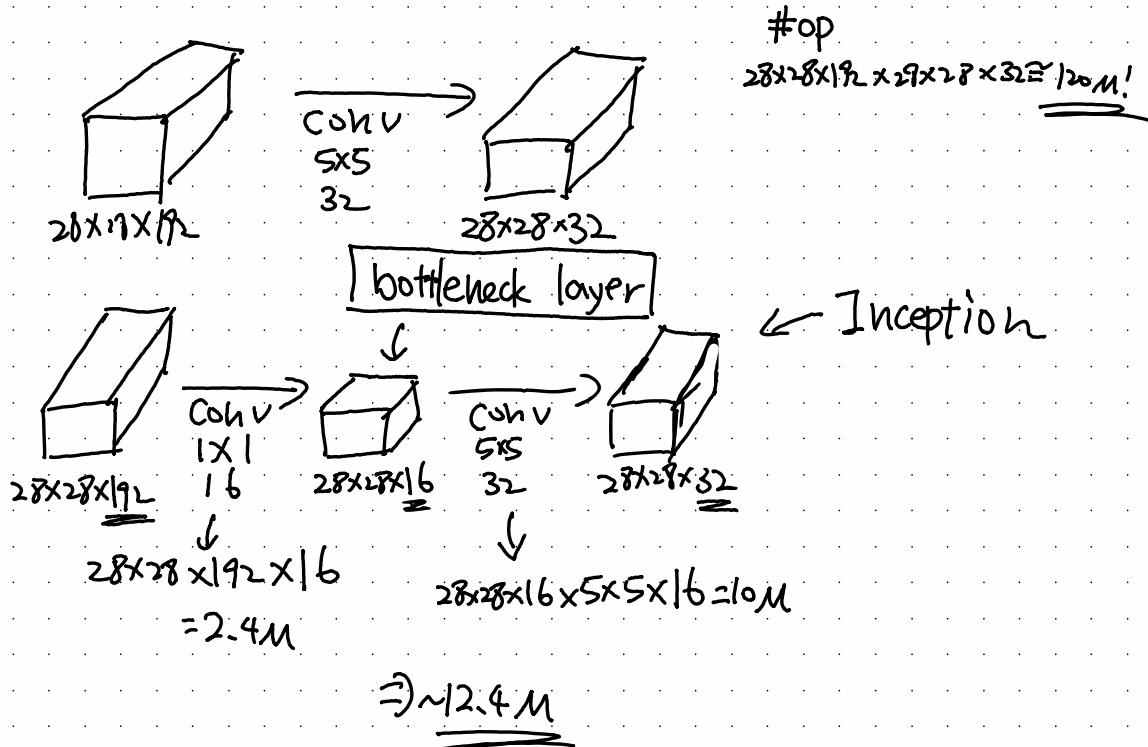
Inception Network, 2014

e.g.

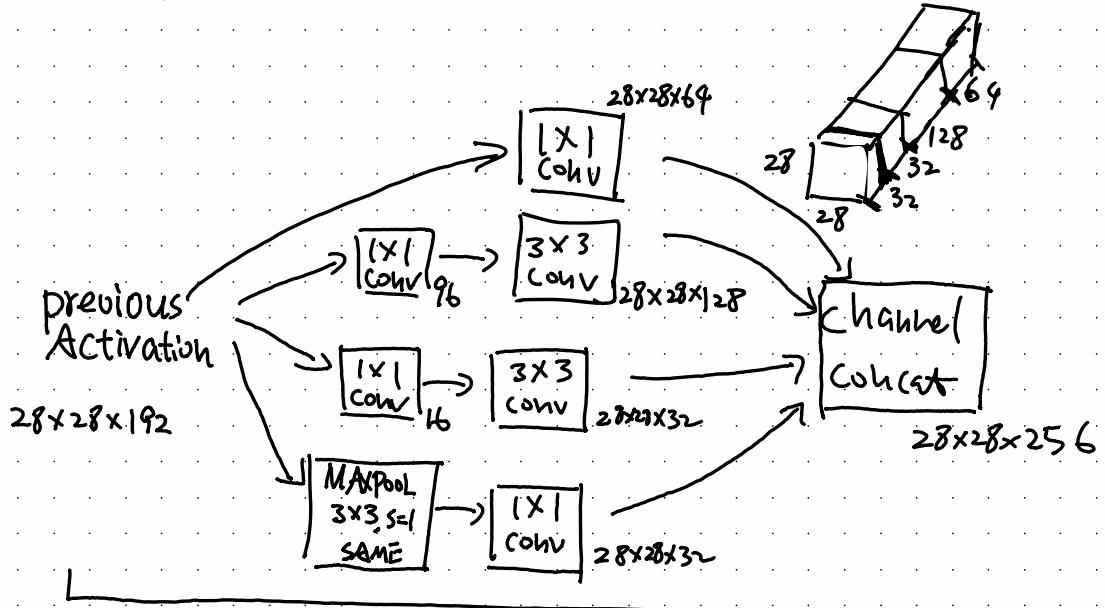


△ don't confuse w.
perception

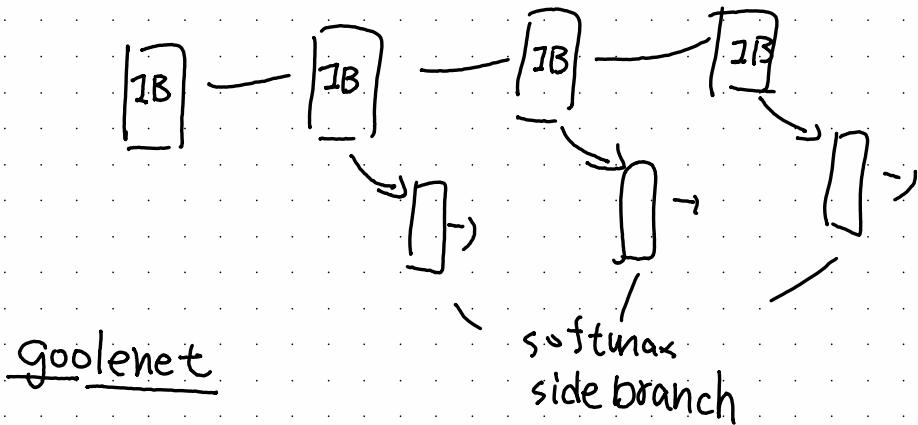
△ try all different size at once!



Inception Module



Inception Block



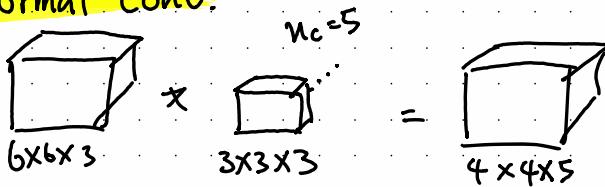
"we need to go deeper"

Inception

MobileNet, 2017

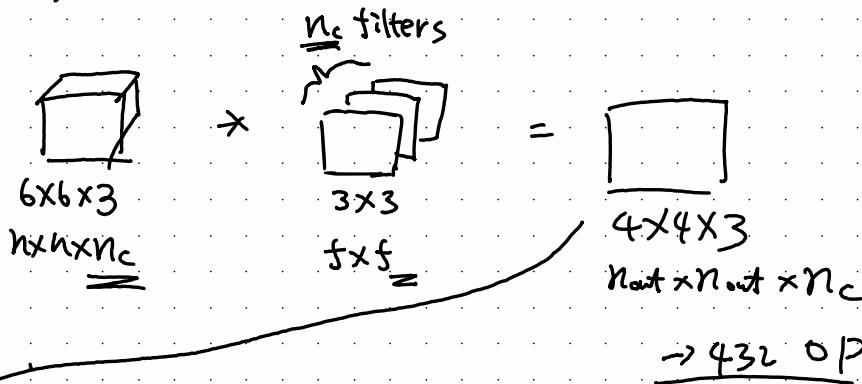
- low computational cost.
- Normal v.s depthwise - separable convolution
⇒ depth-wise + point-wise.

Normal Conv.

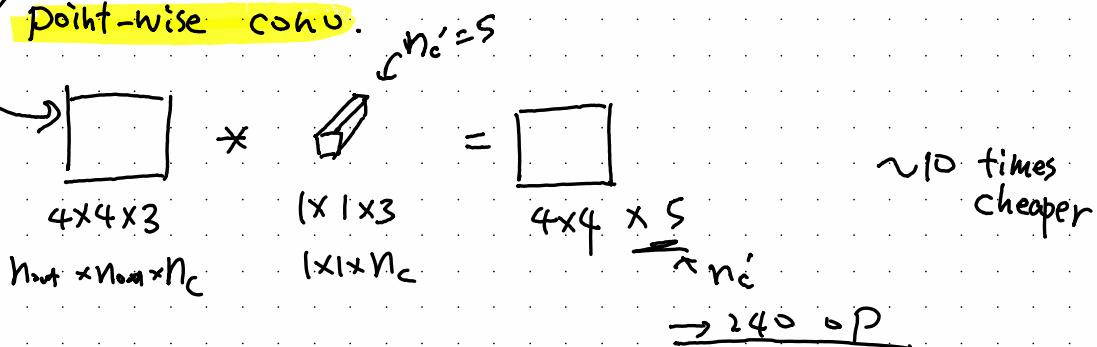


$$\text{cost} = \# \text{filter para} \times \# \text{filter pos} \times \# \text{filter} \\ = 2160 \text{ op.}$$

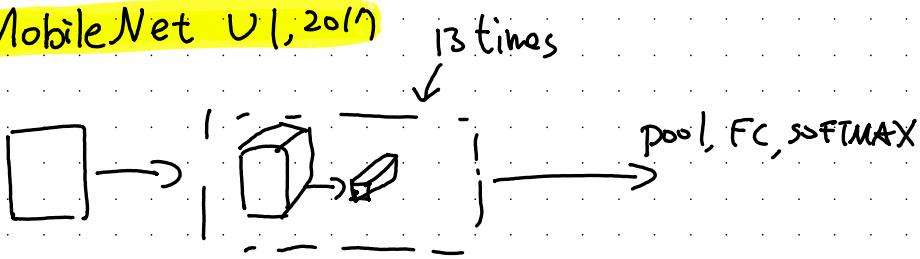
Depth-wise Conv.



Point-wise conv.



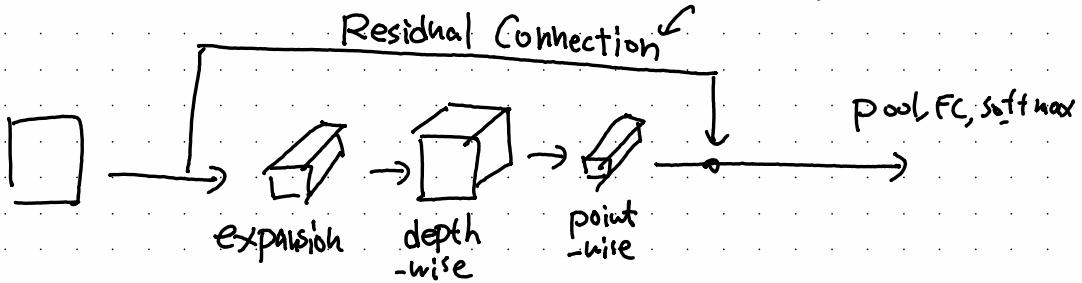
MobileNet V1, 2017



MobileNet V2, 2019

bottleneck blocks

17 times



eg. $n \times h \times 3 \rightarrow n \times h \times 18 \rightarrow n \times h \times 18 \rightarrow n \times h \times 3$

expansion

(1×3)
18 filters

depth
wise



point-wise

$1 \times 1 \times 18$
3 filters

(projection)

- expansion → increase representation in bottleneck block

- projection → reduce bandwidth between layers

EfficientNet. 2019

- automatic scale up/down model

scale r resolution

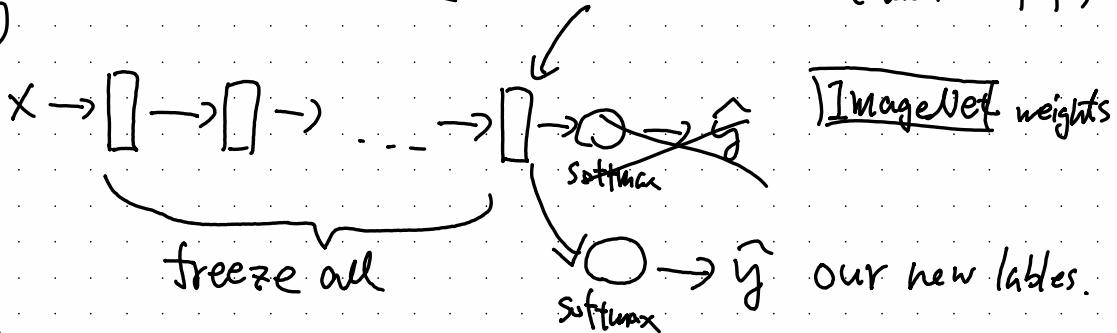
d depth

w width the model selects r,d,w that suits device.

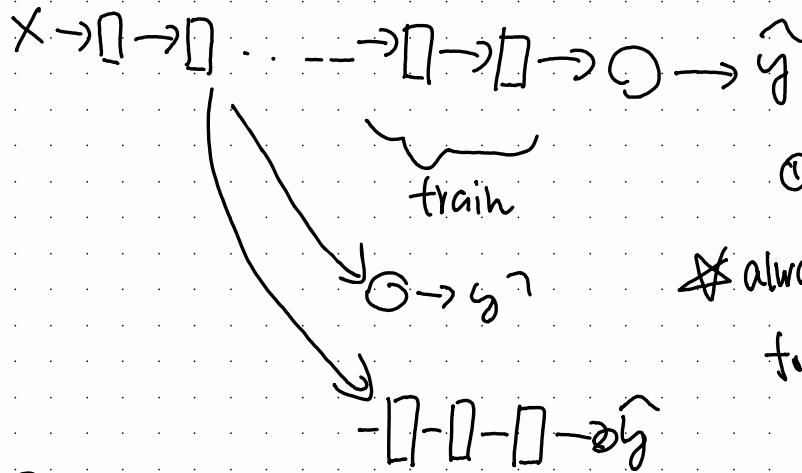
Transfer Learning

precompute this, to accelerate
(Avoid forward prop)

①



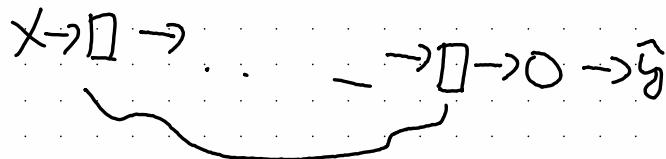
②



①, ② if small dataset.

* always start from existing weights.
for CV problems

③

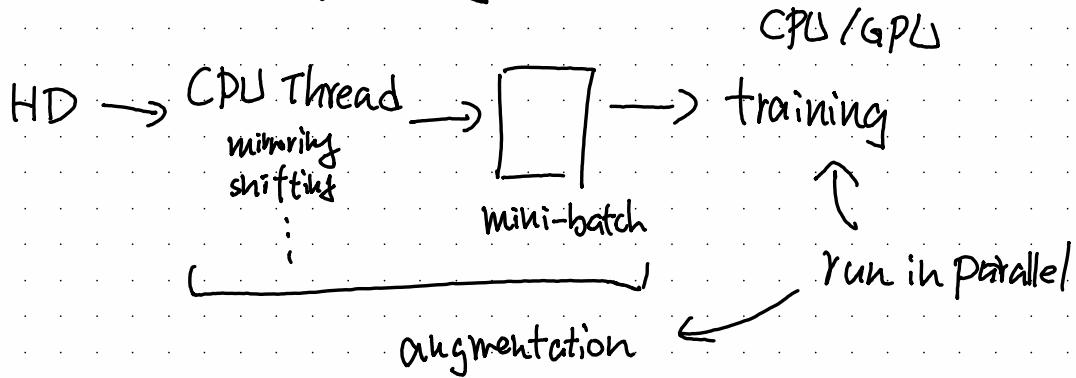


use previous weights
as initialization
(if we have a large data set)

Data Augmentation

- mirroring
- random cropping
 - might crop some less informative area
- color shifting
 - PCA color augmentation
 - simple RGB shifting, e.g. different sunlight.

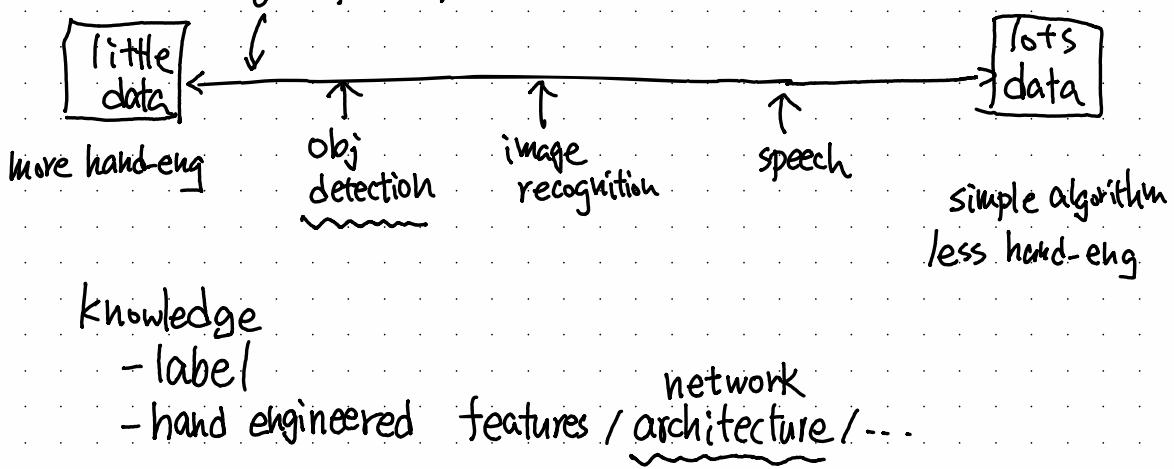
Distortion during training



State for CV.

data vs hand-engineering
(hack)

your specific problem (TL)



Tips for doing well on benchmarks

- Ensembling (3~15 networks)
 - train several networks and average outputs (\hat{y})
Class useful in practical)
- Multi-crop at test time
 - run classifier on multiple vers of test images and avg results
e.g. (10-crops)

Use Open source code

- use network arch.
- use implementation
- use pre-trained model

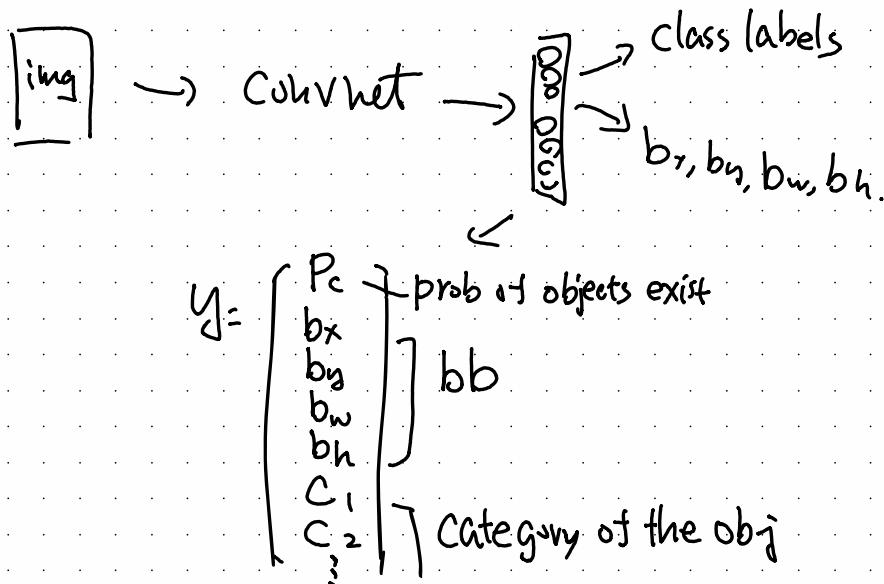
Object detection

classification - car, cat etc

localization - bounding box

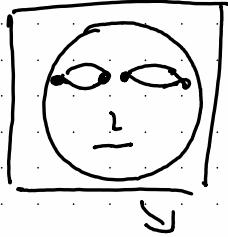
detection - multiple objs of different categories

classification with localization



$$\mathcal{L}(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + \dots + (\hat{y}_j - y_j)^2 & \text{if } y_i = 1 \\ (\hat{y}_i - y_i)^2 & \text{if } y_i = 0 \end{cases}$$

Landmark detection



eyes: l_{1x}, l_{1y}
 l_{2x}, l_{2y}
 l_{4x}, l_{4y}

additional
landmarks
e.g. mouse.
nose
:

CuhVnet



each out a landmark on face

- same for pose detection

skeleton landmarks in a photo.

Obj detection

sliding window detection

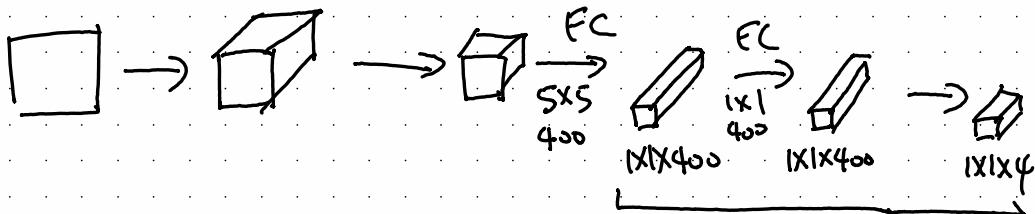
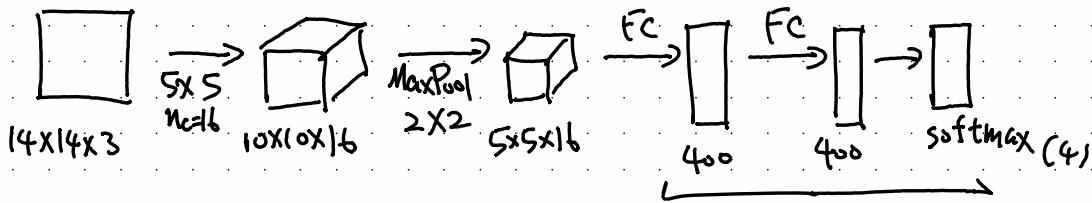


send each window through ConvNet.
from small to large window size

- high computation cost.

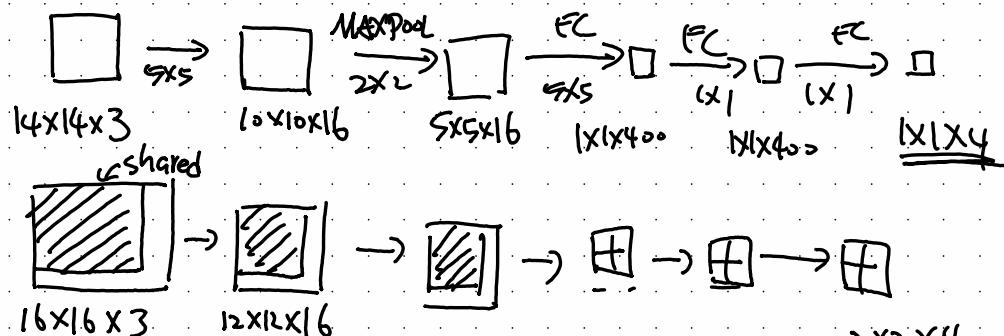
Convolutional Implementation

turn FC to convolutional layer



OverFeat. 2014.

Conv implementation

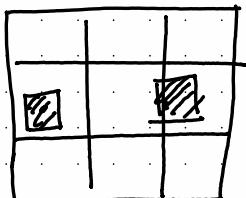


- reuse overlapped computation result sequentially
- instead of send each sliding window through ConvNet, simultaneously do prediction for all window position
- the window position is not precise though.

Output accurate bounding boxes

YOLO - You only look once, 2015 (a difficult paper to read)

(100)



100

target output

grid 3x3x8
y P,b,c

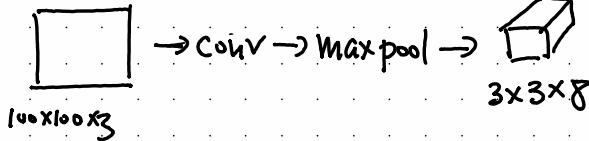
for each grid cell

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

- relative coord

to the grid cell

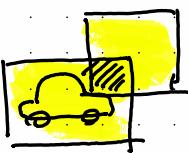
- b_x, b_y - between 0 and 1
 b_w, b_h - could be > 1



⇒ increase grid resolution if need higher precision.

~ still one Cudnn implementation

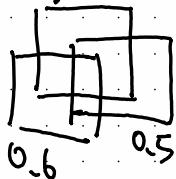
Intersection over Union (IoU)



correct if $\text{IoU} > 0.5$

Non-max Suppression

0.7



find max \rightarrow suppress others with high IoU

each grid cell $\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \end{bmatrix}$

discard all boxes with $p_c \leq 0.6$

for remaining boxes

- pick largest p_c

discard remaining with $\text{IoU} \geq 0.5$ with other boxes.

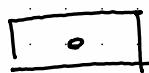
Anchor boxes



predefined shapes



box 2



Anchor box 1

$y =$

$\begin{bmatrix} p_c \\ b_x \\ b_y \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$

$\begin{bmatrix} p_c \\ b_x \\ b_y \\ ; \end{bmatrix}$

Anchor box 2

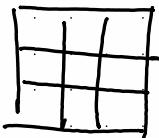
each obj assigned to cell
contain mid-point.

- Anchor box for the cell with high IoU

$$\Rightarrow y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Anchor box 1
Anchor box 2

YOLO algorithm



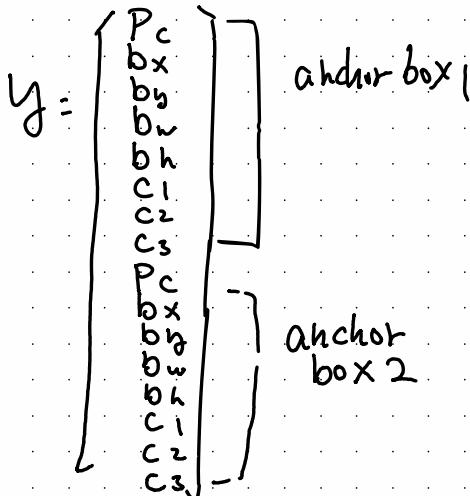
C

- 1 - pedestrian
- 2 - car
- 3 - motorcycle

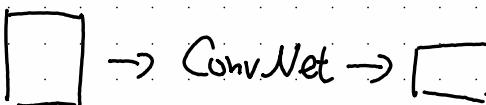
$$y \quad 3 \times 3 \times 2 \times 8$$

`

anchor boxes 5 + #classes

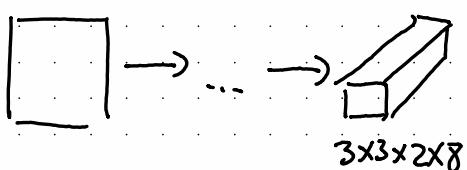


Training



$$3 \times 3 \times 16$$

Prediction



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

$$y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 0.2 \\ 0.3 \\ 0.2 \\ 0.1 \\ 0 \\ 0 \\ \vdots \end{bmatrix}$$

non-max suppression

for each cell,

- get 2 predicted bb (for each anchor boxes)
- get rid of low p_c prediction
- for each class, run non-max suppression

Region Proposals. R-CNN, 2013

- pick interesting region to run CNN

segmentation → run CNN on interesting regions

R-CNN: propose region, CNN, output

Fast R-CNN : propose region, Convolution implementation, output
C20/5)

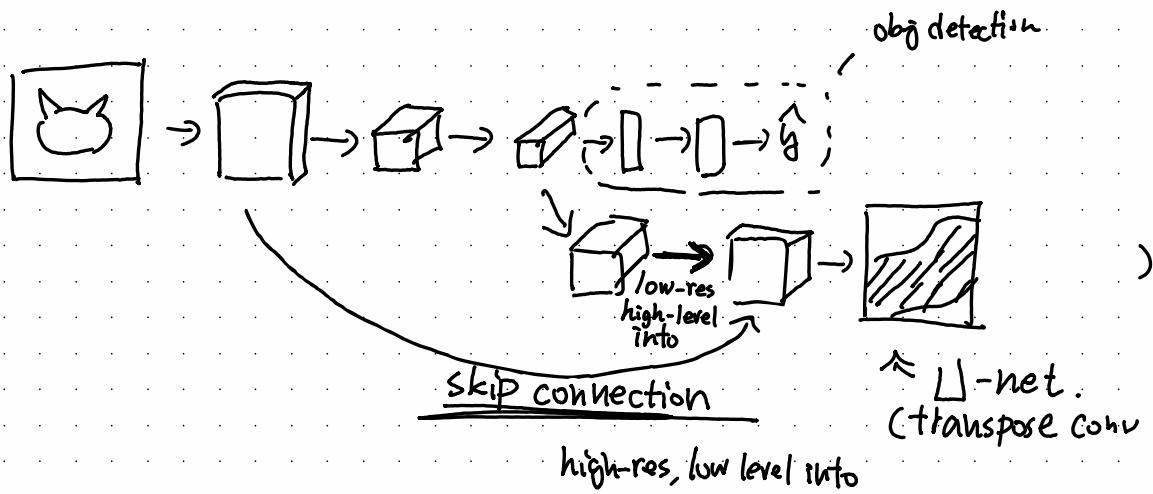
Faster R-CNN : conv net to propose regions (2016)

still slower than YOLO...

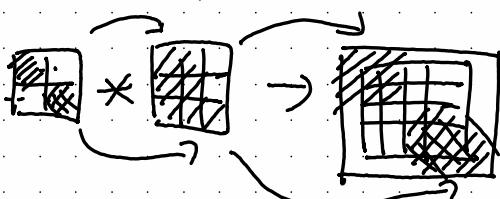
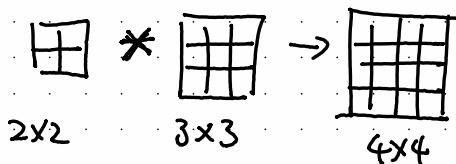
Semantic Segmentation U-net, 2017

obj detection - bb

semantic seg - per-pixel labeling

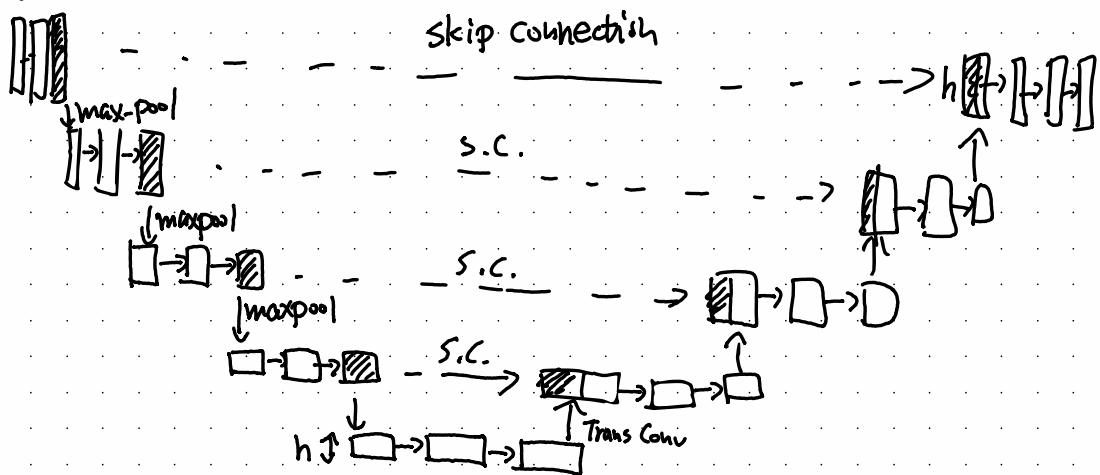


Transpose Conv, deconvolution, upscale conv



if overlap, add values together.

U-Net, 2015 (on biomedical image seg.)



Face Recognition

- verification

- input image, name
- output whether image/name match
(1 to 1)

- Recognition

- database of k people
- input image
- output ID of the image among k.
(1 to k)

one-shot learning

learn from one example to recognise the person again

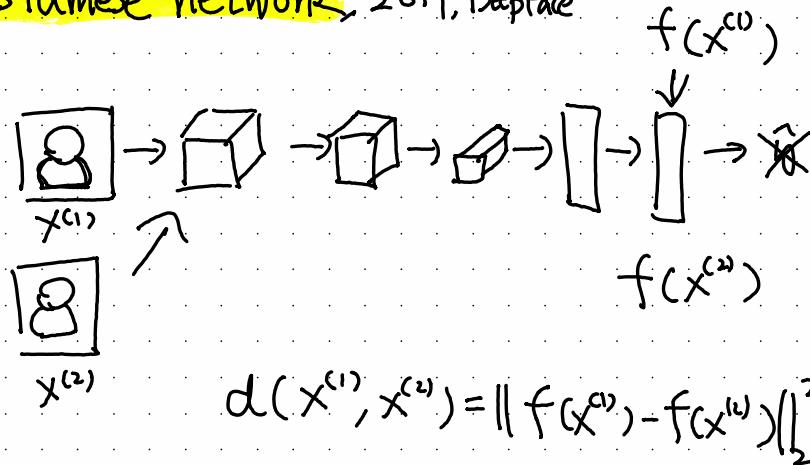
learn a "similarity" function

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

if $d(\text{img1}, \text{img2}) > \tau$, same person

$< \tau$, no match

Siamese network, 2014, DeepFace

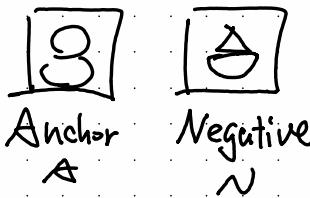
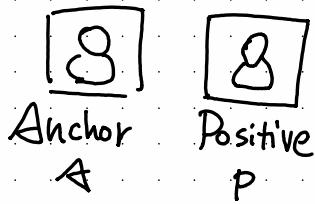


parameters of NN define $f(x^{(1)})$

learn parameters that

if $x^{(1)}, x^{(2)}$ are same person $d(x^{(1)}, x^{(2)})$ small
different large

Triplet Loss FaceNet, 2015



$$\text{WANT: } \|f(A) - f(P)\|^2 + \alpha \leq \|f(A) - f(N)\|^2$$

$$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha \leq 0$$

$\underbrace{}$ margin

to prevent trivial solution

$$\text{e.g. } f(\text{img}) = 0$$

$$\text{or } f(\text{img}^1) = f(\text{img}^2),$$

Loss Function

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

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

training set: 10K pictures of 1K persons

Choosing A, P, N

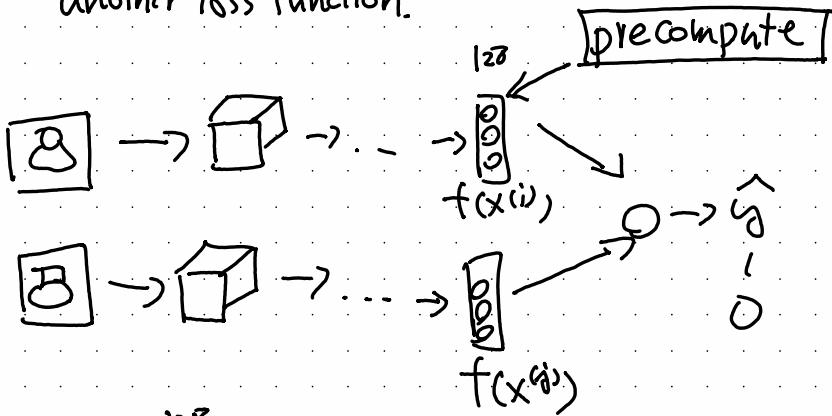
if chosen randomly, $d(A, P) + \alpha \leq d(A, N)$ is too easy.

choose A, P, N that are "hard" (NN learns little)
(check paper)

Current K = 100M images
commercial

DeepFace, 2014

another loss function.



$$\hat{y} = \sigma \left(\sum_{k=1}^{128} w_k |f(x^{(i)})_k - f(x^{(j)})_k| + b \right)$$

- kth feature or x^2

Neural Style Transfer, 2015

Content + style \rightarrow generated
(C) (S) (G)

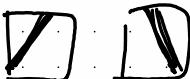
Visualize ConvNet 2013

hidden unit covers larger patch
 \rightarrow

[8] \rightarrow conv \rightarrow conv $\rightarrow \dots \rightarrow$ FC \rightarrow FC $\rightarrow \hat{y}$

pick a unit in layer 1, find image patches that maximize
unit activation
repeat for all units
patch size \uparrow , to deeper layer.

lower lr -



higher level



Neural Style Transfer: 2015

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

1. init G randomly (like noise)
2. gradient descent to minimize $J(G)$

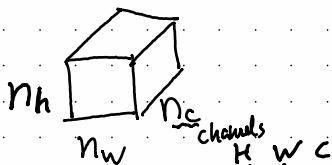
$J_{\text{content}}(C, G)$

- use hidden layer l to compute content cost (usually middle of network)
- use pre-trained ConvNet.
- Let a^{enc} and $a^{\text{enc}(G)}$ be activation of layer l

$$J_{\text{cont}}(C, G) = \frac{1}{2} \| a^{\text{enc}} - a^{\text{enc}(G)} \|^2$$

$J_{\text{style}}(S, G)$

at layer l , define "style" as correlation between activation across channels



intuition: each channel corresponds

Q_{ijk} = activation at (i, j, k) . $G^{\{l\}} : n_c^{\text{enc}} \times n_c^{\text{enc}}$ to hidden units that capture different levels of features

$$G_{kk'}^{\text{enc}(S)} = \sum_i \sum_j Q_{ijk}^{\text{enc}} Q_{ijk'}^{\text{enc}} \leftarrow \text{style} \quad \text{"Gram Matrix"}$$

$$G_{kk'}^{\text{enc}(G)} = \dots \leftarrow \text{generated}$$

$$J_{\text{style}}(S, G) = \| G^{\text{enc}(S)} - G^{\text{enc}(G)} \|_F^2 = \frac{1}{2n_H n_w n_c} \sum_k \sum_{k'} (G_{kk'}^{\text{enc}(S)} - G_{kk'}^{\text{enc}(G)})^2$$