

Computer Vision.

e.g. Image classification.

Object detection.

Neural Style Transfer.

Deep learning on large images:

$$64 \times 64 \times 3 = 12288. \text{ it's ok.}$$

↑
3 RGB channel.

$1000 \times 1000 \times 3 = 3 \text{ million. which mean the input has } 3m \text{ } x_i.$

Suppose next layer has 1000 units.

Then the first weight matrix W will be $1000 \times 3m$ matrix. (HUGE).

So we need CONVOLUTION technique.

Edge detection example.

Given a picture.

① Detect vertical edges.

② Detect horizontal edges.

element-wise product.

Vertical edge detection.

$$3 \times 1 + 1 \times 1 + 2 \times 1 + 0 + 0 + 0 + 1 \times (-1) + 8 \times (-1) + 2 \times (-1) = -5$$

$\begin{matrix} 3 & 0 & 1 & 2 & 7 & 4 \end{matrix}$

kernel matrix.

$$\begin{matrix} & & & \\ & & & \\ & & & \\ -5 & -4 & 0 & 8 \end{matrix}$$

$\begin{matrix} 1 & 5 & 8 & 9 & 3 & 1 \end{matrix}$

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

$$\begin{matrix} -10 & -2 & 2 & 3 \end{matrix}$$

$\begin{matrix} 2 & 7 & 2 & 5 & 1 & 3 \end{matrix}$

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

$$\begin{matrix} 0 & -2 & -4 & -7 \end{matrix}$$

$\begin{matrix} 0 & 1 & 3 & 1 & 7 & 8 \end{matrix}$

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

$$\begin{matrix} -3 & -2 & -3 & -16 \end{matrix}$$

$\begin{matrix} 4 & 2 & 1 & 6 & 2 & 8 \end{matrix}$

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

$$\begin{matrix} 4 \times 4 \text{ result.} \end{matrix}$$

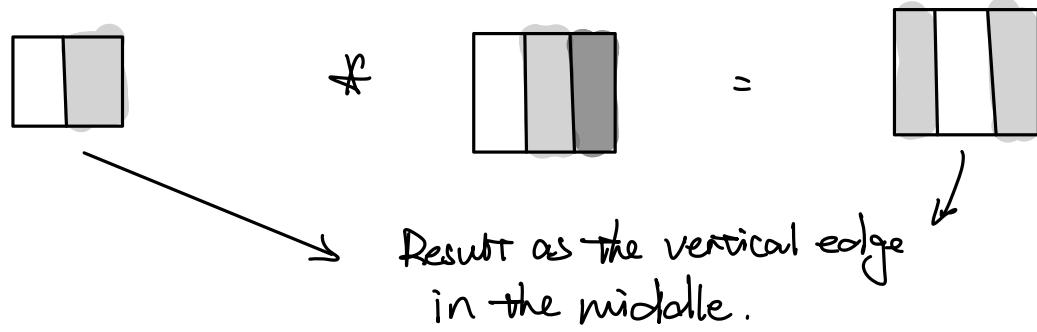
$\begin{matrix} 2 & 4 & 5 & 2 & 3 & 9 \end{matrix}$

Convolution operation

6x6 gray scale image.

Code e.g. Python: conv-forward. Tf: tf.nn.conv2d. keras: Conv2D.

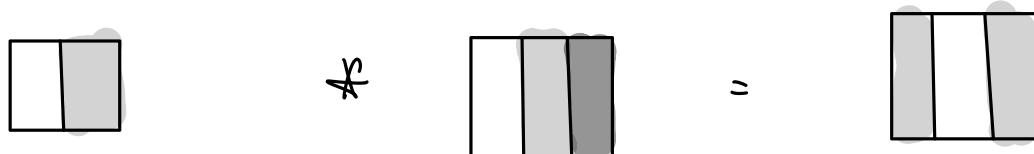
$$\begin{array}{ccccccc} 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \end{array} \quad * \quad \begin{array}{ccc} 1 & 0 & -1 \end{array} = \begin{array}{ccccccc} 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \end{array}$$



More edge detection.

Previous example.

$$\begin{array}{ccccccc} 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \end{array} \quad * \quad \begin{array}{ccc} 1 & 0 & -1 \end{array} = \begin{array}{ccccccc} 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \end{array}$$



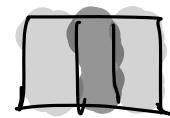
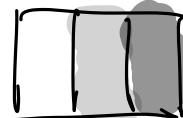
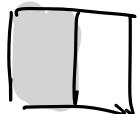
What about we reverse the pos. of 10 and 0.

0 0 0 10 10 10

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

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

$$\begin{matrix}
 0 & -30 & -30 & 0 \\
 0 & -30 & -30 & 0 \\
 0 & -30 & -30 & 0 \\
 0 & -30 & -30 & 0
 \end{matrix}$$



It looks darker.

Means from dark to light.

Vertical and Horizontal Edge detection.

$$\begin{matrix}
 1 & 0 & -1 \\
 1 & 1 & 1
 \end{matrix} \text{ lighter top.}$$

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

$$\begin{matrix}
 1 & 0 & -1 \\
 -1 & -1 & -1
 \end{matrix} \text{ darker bottom.}$$

Vertical . Horizontal .

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

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

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

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

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

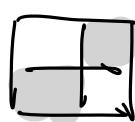
$$\begin{matrix}
 0 & 0 & 0 & 10 & 10 & 10
 \end{matrix}$$

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

positive edge neg. edge

blending edge.

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



Learning to detect edges.

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

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

$$2 \quad 0 \quad -2$$

$$1 \quad 0 \quad -1$$

$$10 \quad 0 \quad -10$$

$$3 \quad 0 \quad -3$$

Sobel filter.

Scharr filter.

More weight on middle.

Actual, you can learn all these 9 numbers. Treating them as params

Padding.

We see: $(6 \times 6) * (3 \times 3) = (4 \times 4)$.

Math: $(n \times n) * (f \times f) = [(n-f+1) \times (n-f+1)]$

When looking the corner, these pixels only being used once.

Edge pixels received less attention. (Throw informations).

Solution: pad the image with additional border.

6×6 image into 8×8 image.

And $(8 \times 8) * (3 \times 3) = (6 \times 6)$. Preserved the original image size.

Valid and Same Convolution.

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

no padding.

"Same": Pad so that output size is the same as the input size.

P: padding size. Say the kernel is $f \times f$.

So if $P = \frac{f-1}{2}$. Then the image will preserve to be the same size.

e.g. $f=3$. $P = \frac{3-1}{2} = 1$. $f=5$. $P = \frac{5-1}{2} = 2$.

By convention, the width of filter is usually odd.

(So that you can have a central pixel, to talk about position)

Strided Convolution.

We take element-wise product as usual. But we take a bigger step.

$$(n \times n) * (f \times f)$$

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

padding p . Stride s .

$$\text{e.g. } (7 \times 7) * (3 \times 3). : \frac{7+0-3}{2} + 1 = 3. \text{ final } 3 \times 3 \text{ result.}$$

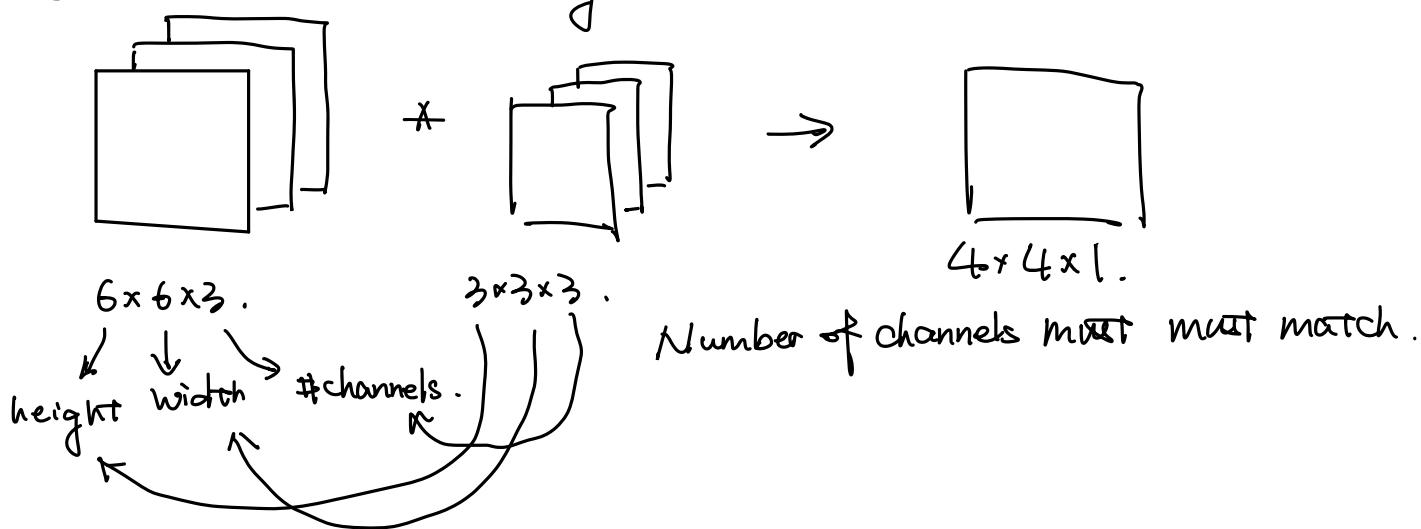
padding 0. stride 2.

$$\text{We can add floor operation } \left\lfloor \frac{n+2p-f}{s} \right\rfloor$$

By convention, if the image was not fully land (include padding) inside the filter, we just don't do operation.

Convolutions over volumes.

Convolution on RGB images.



Consider $3 \times 3 \times 3 \Rightarrow 1$ numbers in the filter. Put this "Cube" in the 3 channel image. Element-wise product. Summation. \Rightarrow Result.

e.g. You might want to detect red edges.

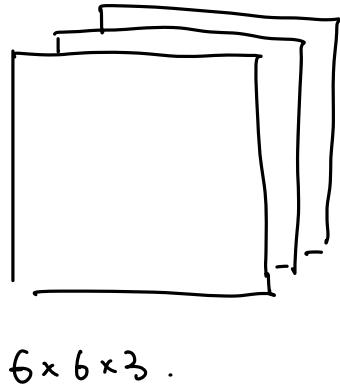
$$R : \begin{matrix} 1 & 0 & -1 \\ 0 & -1 & 0 \\ 1 & 0 & -1 \end{matrix}$$

$$Gr: \begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix}$$

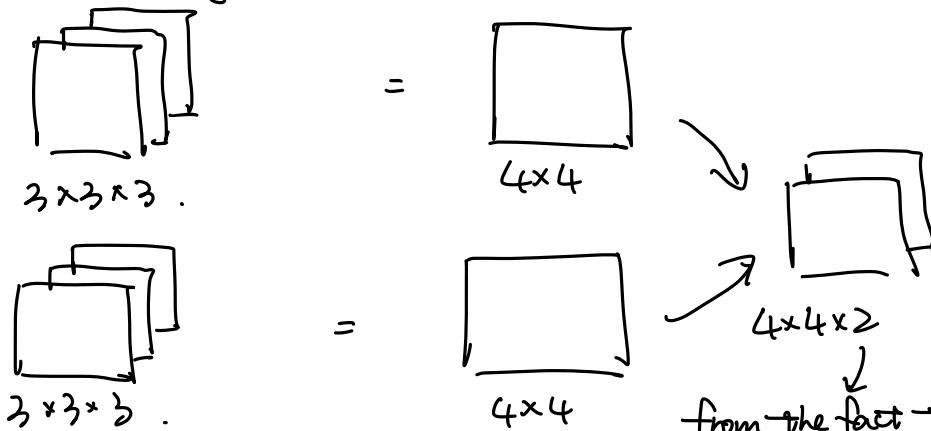
$$B : \begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix}$$

Edge detector only in red.

Multiple filters.



Vertical edge detector



from the fact that
we use 2 filters

Horizontal edge detector.

Number of channels

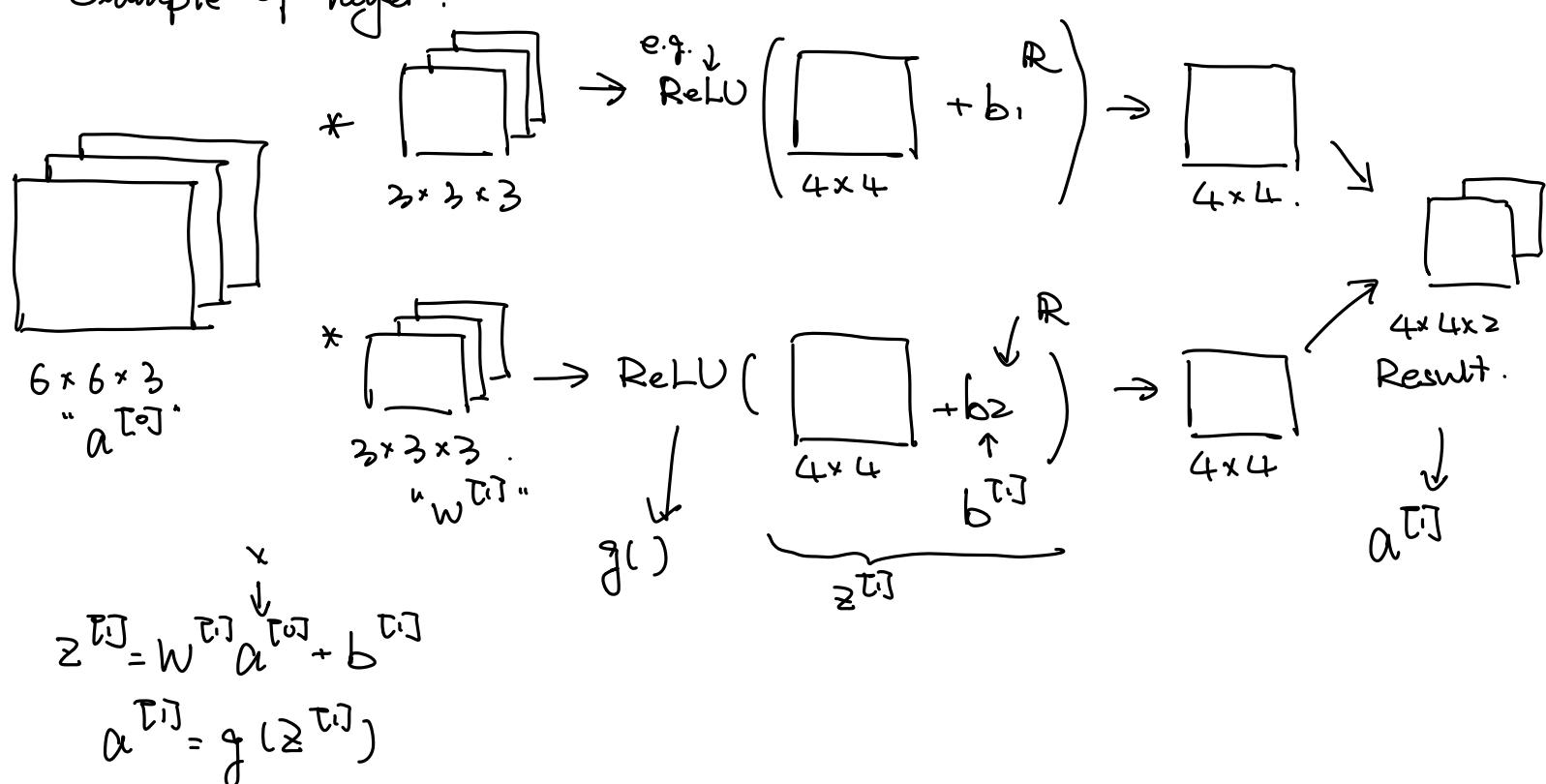
$$\text{Summary: } n \times n \times n_c * f \times f \times n_c \rightarrow n-f+1 \times n-f+1 \times n_c$$

$6 \times 6 \times 3 * 3 \times 3 \times 3$

Number of filters.

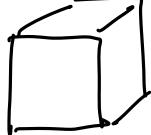
One layer of a convolutional network.

example of layer.



Number of params. in one layer.

If you have 10 filters that are $3 \times 3 \times 3$ in one layer of a NN.



$3 \times 3 \times 3$ 27 params.
+ bias: 28 param

in total 280 params.

No matter how large the input is, you can always use these filters to detect features.

Summary of Notation.

If layer l is a conv. layer.

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

Input: $n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$.

$p^{[l]}$ = padding size.

Height . Width . #Channel.

$s^{[l]}$ = Striding size.

Output: $n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$.

$n_C^{[l]}$ = number of filters.

Each filter is: $f^{[l]} \times f^{[l]} \times n_C^{[l]}$ Notice: $n_H^{[l]} = \left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$

Activation: $a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

$A^{[l]} \rightarrow m \times n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

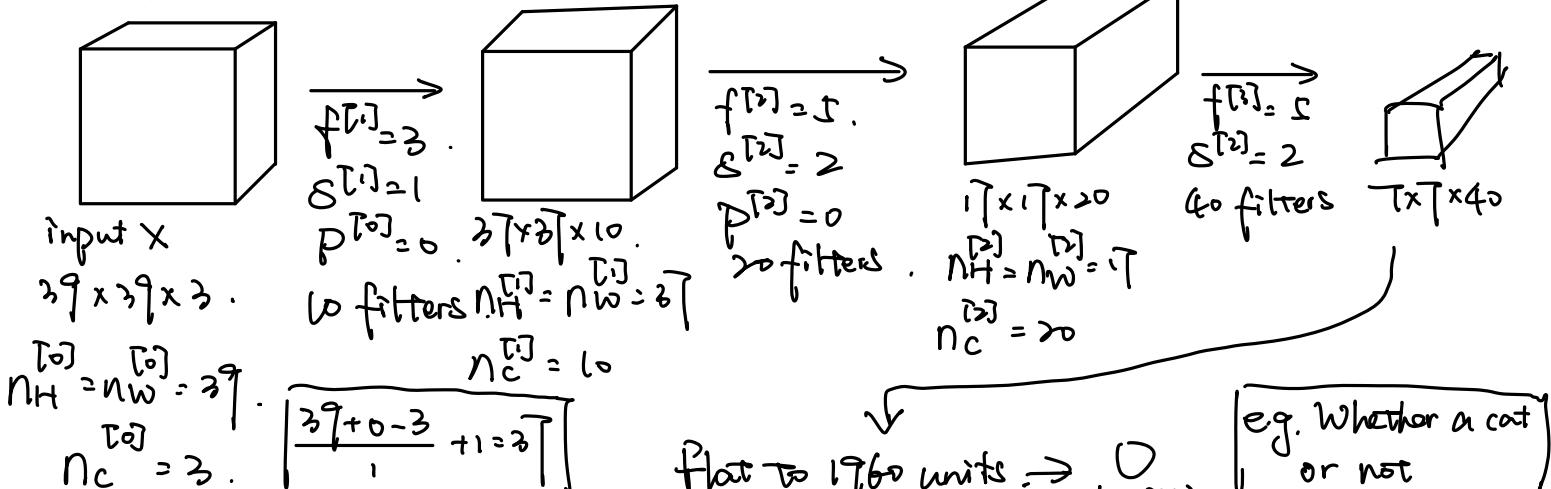
where m is the batch size.

Weights: $\underbrace{f^{[l]} \times f^{[l]} \times n_C^{[l-1]}}_{\text{One filter size.}} \times \underbrace{n_C^{[l]}}_{\text{Number of filters in layer } l}$.

Bias: $n_C^{[l]} \rightarrow \underbrace{(1, 1, 1, n_C^{[l]})}_{\text{later for convenience}}$

A Simple Conv. NN example.

Example. ConvNet.



Types of layer in a conv. net.

- Convolution. (conv.) .
- Pooling. (POOL) .
- Fully connected. (FC) .

Pooling layers

Max pooling.

$$\begin{array}{cccc}
 1 & 3 & 2 & 1 \\
 2 & 9 & 1 & 1 \\
 1 & 3 & 2 & 3 \\
 5 & 6 & 1 & 2
 \end{array} \xrightarrow[f=2]{s=2} \begin{array}{cc}
 9 & 2 \\
 6 & 3
 \end{array}$$

2x2.

4x4

Notice: No params to learn. A fixed layer.

e.g. 1 3 2 1 3

$$\begin{array}{ccccc}
 2 & 9 & 1 & 1 & 5 \\
 1 & 3 & 2 & 3 & 2 \\
 8 & 3 & 5 & 1 & 0 \\
 5 & 6 & 1 & 2 & 9
 \end{array} \xrightarrow[f=3]{s=1} \begin{array}{ccc}
 9 & 9 & 5 \\
 9 & 9 & 5 \\
 8 & 6 & 9
 \end{array}$$

5x5.

Average Pooling .

$$\begin{array}{cccc}
 1 & 3 & 2 & 1 \\
 2 & 9 & 1 & 1 \\
 1 & 3 & 2 & 3 \\
 5 & 6 & 1 & 2
 \end{array} \xrightarrow[f=2]{s=2} \begin{array}{cc}
 3.75 & 1.25 \\
 4 & 2
 \end{array}$$

Max pooling use much more often than Average pooling.

Summary of Pooling :

Hyperparams : Common value : $f=2, s=2$; $f=3, s=2$.

f : filter size.

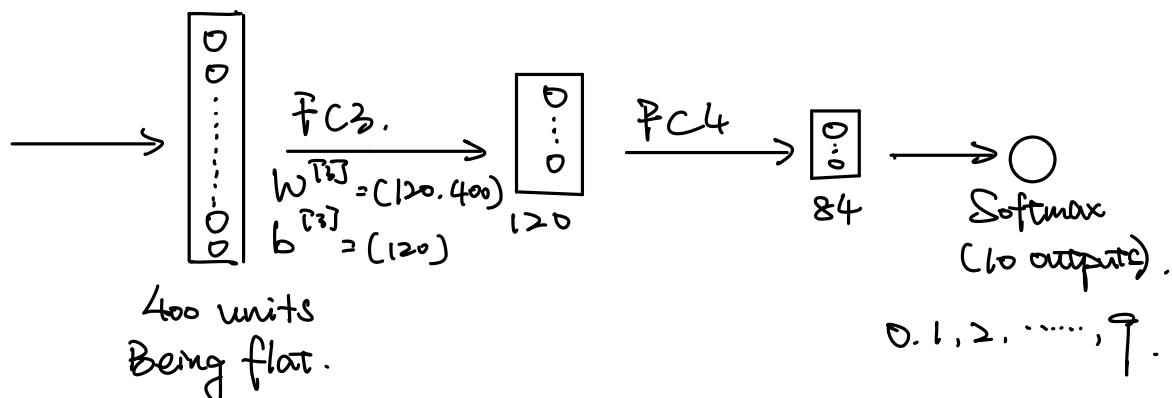
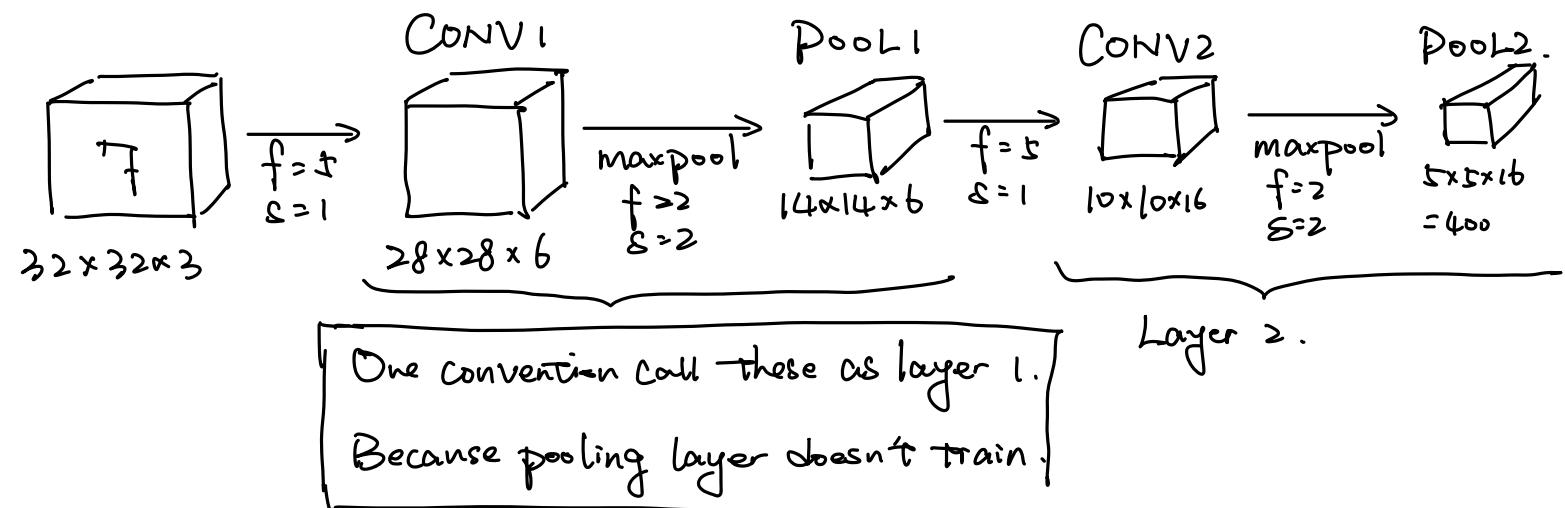
$$n_H \times n_W \times n_C \rightarrow \left\lfloor \frac{n_H - f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n_W - f}{s} + 1 \right\rfloor \times n_C.$$

s : stride size.

Max or average pooling.

CNN example . (LeNet-5).

Digit recognition .



Common pattern. n_H, n_W decrease, n_C increase.

CONV - POOL - CONV - POOL - FC - FC - Softmax.

Activation Shape . Activation Size . # Params .

Input

$(32, 32, 3)$

3×7^2

0

CONV1

$(28, 28, 8)$

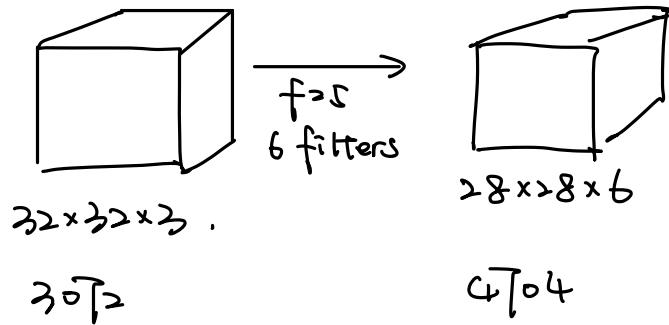
62×7^2

> 208

Pool1	(14, 14, 8)	1568	0
Conv2	(10, 10, 16)	1600	416.
Pool2	(5, 5, 16)	400	0
FC3	(120, 1)	120	4800.
FC4.	(84, 1)	84	10081.
Softmax .	(10, 1) .	10	841.

Why Convolutions.

Say:



If you create a fully connected network between the two.

But by using conv layer, it has $6 \times (5 \times 5 + 1) = 156$ parameters.

You need 14m weights.

Reason ①: Param Sharing : A feature detector (such as vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.

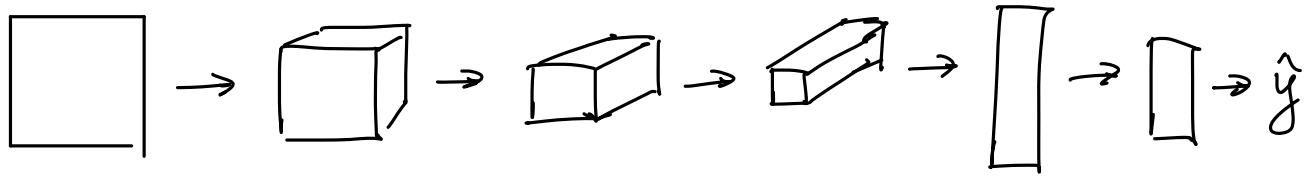
E.G: You don't need to learn multiple edge detectors, one for upper corner and another for bottom corner.

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

e.g Using 3×3 filter, each output only have relation with 9 inputs

Putting it together.

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



$$\text{Cost } J = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (\hat{y}^{(i)} - y^{(i)})^2$$

Use gradient descent to optimize parameters to reduce J .

Classic Networks.

→ LeNet-5.

→ AlexNet

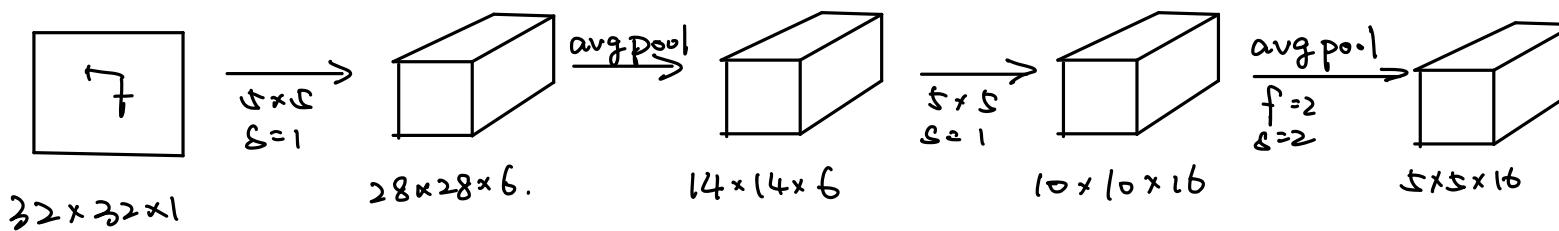
→ VGG.

ResNet. (152). deep

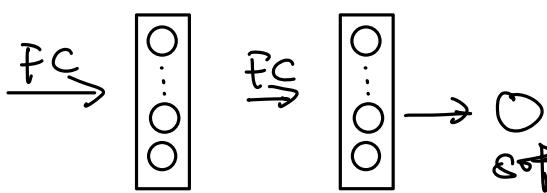
Inception.

Classic Networks.

LeNet-5.

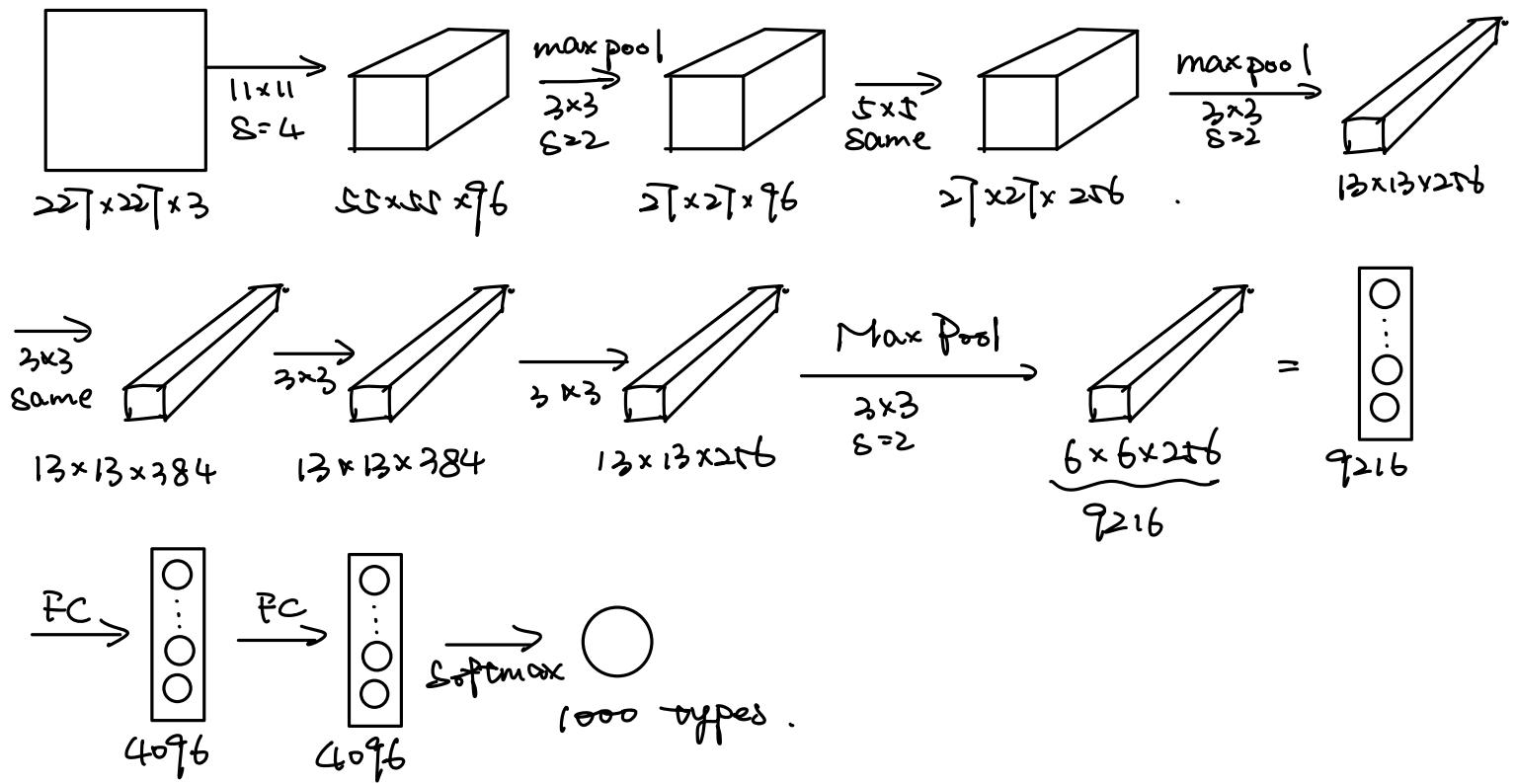


Gray scale



This network is actually small.
It only takes 60k params.
softmax. Trend: $n_h, n_w \downarrow, n_c \uparrow$

AlexNet.



→ Similar to LeNet. But much bigger. ~ 60m params.

→ Using ReLU activation function.

→ Multiple GPUs.

→ Local response normalization (LRN). (Doesn't help much.)

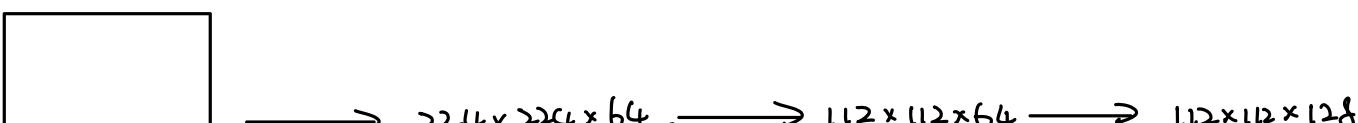
e.g. look one of the blocks. Say $13 \times 13 \times 256$ block. Look at one pixel. and normalize all 256 of them at that point.

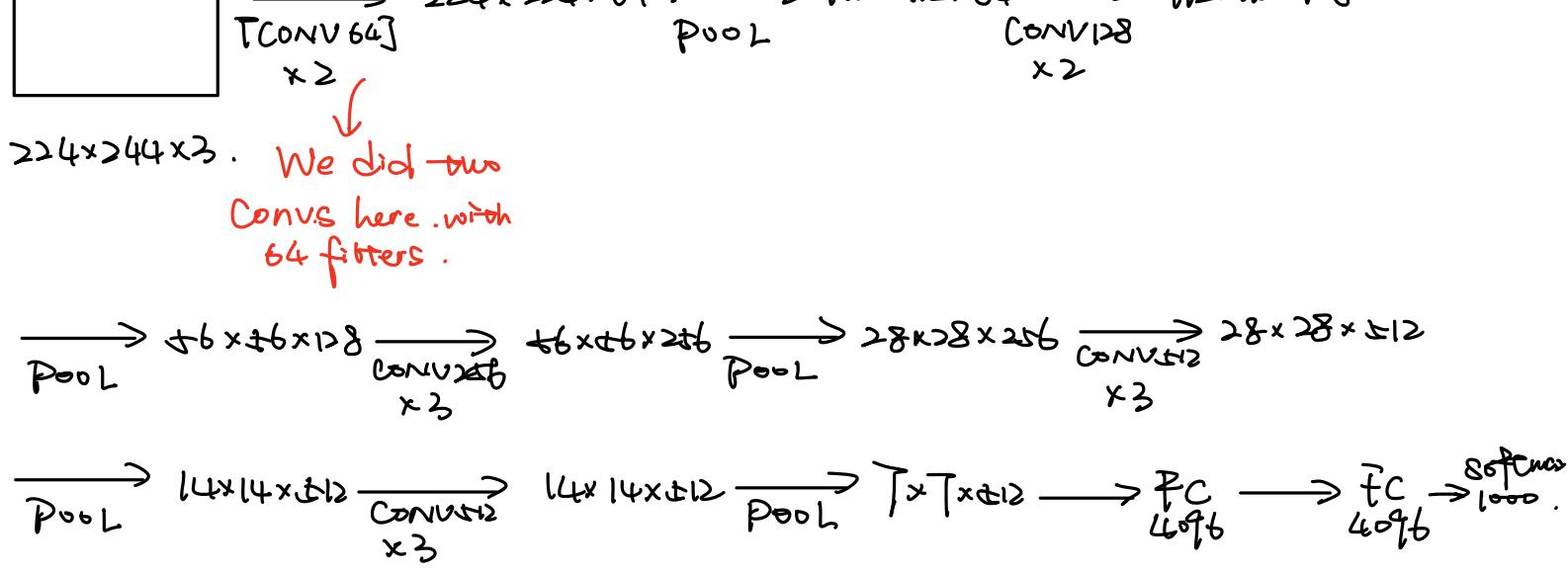
Intuition: For each position, you don't want too large/small activation.

VGG-16.

$\text{CONV} = 3 \times 3$ filters, $s=1$. Same. $\text{MAX-POOL} = 2 \times 2$, $s=2$.

It simplifies the CNN.





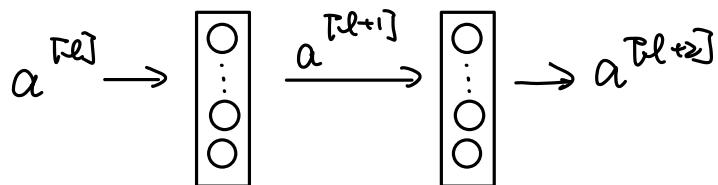
$\sim 1.38 \text{ m}$ params . Pretty large.

The "16" in the name mean there are 16 layers in the network that has weights. Other e.g. VGG-19.

Notice. #filters keep doubling when the network going deeper.

Residual Network . (ResNets) .

Residual Block .



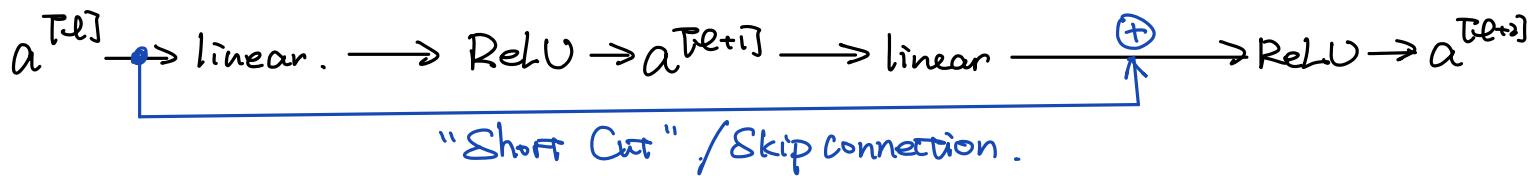
Steps of computation :

$$\begin{aligned}
 a^{Tl} &\xrightarrow{\text{linear.}} \text{ReLU} \rightarrow a^{Tl+1} \xrightarrow{\text{linear}} \text{ReLU} \rightarrow a^{Tl+2} \\
 z^{Tl+1} &= W^{Tl+1} a^{Tl} + b^{Tl+1}, \quad a^{Tl+1} = g(z^{Tl+1}), \quad z^{Tl+2} = W^{Tl+2} a^{Tl+1} + b^{Tl+2}, \quad a^{Tl+2} = g(z^{Tl+2})
 \end{aligned}$$

From a^{Tl} to a^{Tl+2} . it will take all the way.

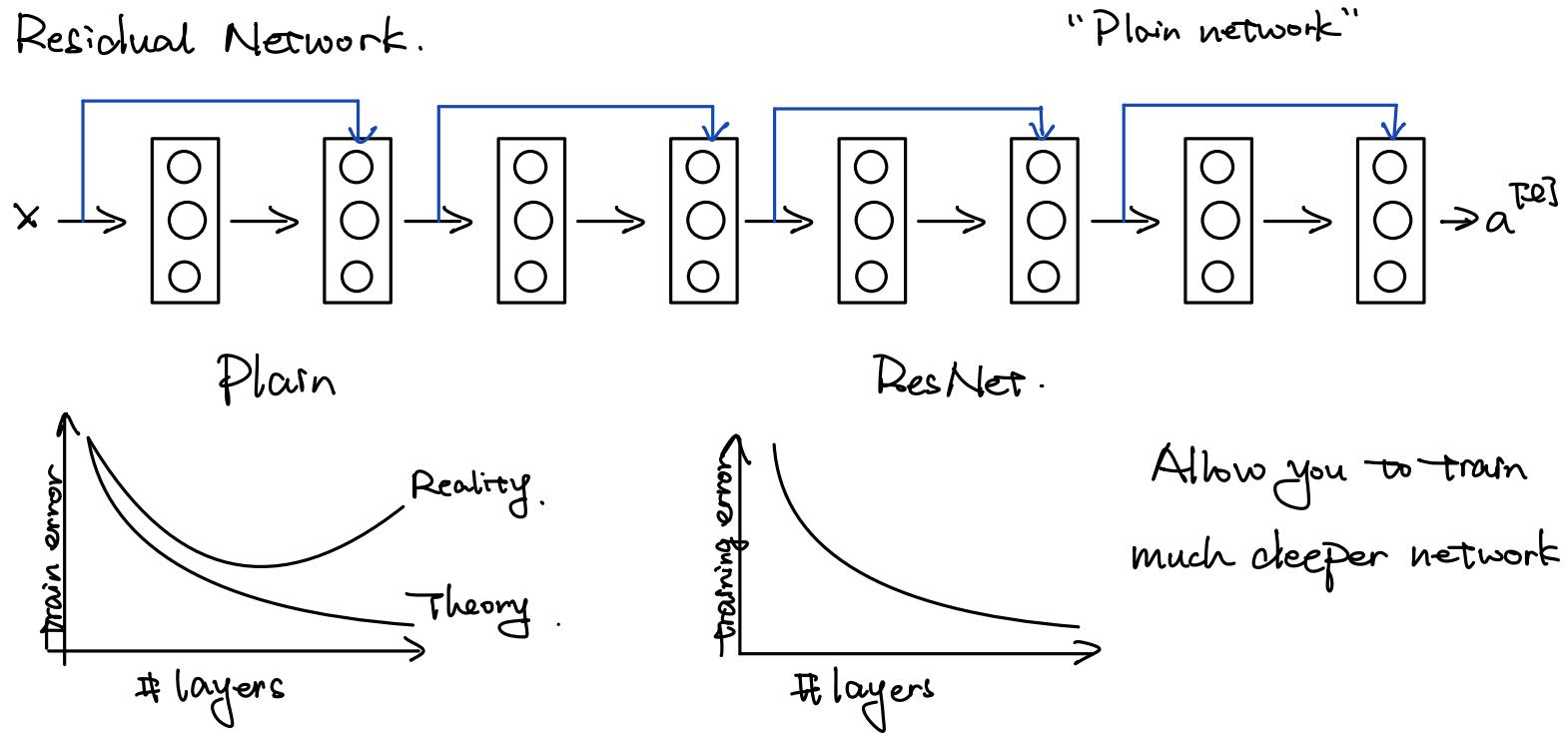
Call this "main path".

What residual block did is .



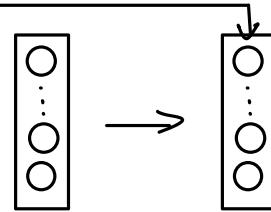
$$\begin{aligned} z^{[l+1]} &= W^{[l+1]} a^{[l]} + b^{[l+1]}, \quad a^{[l+1]} = g(z^{[l+1]}), \quad z^{[l+2]} = W^{[l+2]} a^{[l+1]} + b^{[l+2]}. \quad a^{[l+2]} = g(z^{[l+2]} + a^{[l+1]}) \\ &\quad \text{a}^{[l+2]} = g(z^{[l+2]} + a^{[l+1]}) \end{aligned}$$

Residual Network.



Why ResNet work

$$\begin{aligned} x \rightarrow \text{Big NN} &\rightarrow a^{[l+1]} \\ \Rightarrow x \rightarrow \text{Big NN} &\rightarrow a^{[l+1]} \xrightarrow{\text{ReLU activation } a \geq 0} \end{aligned}$$



$$\begin{aligned} a^{[l+2]} &= g(z^{[l+2]} + a^{[l+1]}) \\ &= g(W^{[l+2]} a^{[l+1]} + b^{[l+2]} + a^{[l+1]}). \\ &\quad \uparrow \text{if } W^{[l+2]} = 0, b^{[l+2]} = 0, \quad \downarrow g(a^{[l+1]}) = a^{[l+1]}. \text{ Since } a^{[l+1]} \geq 0. \end{aligned}$$

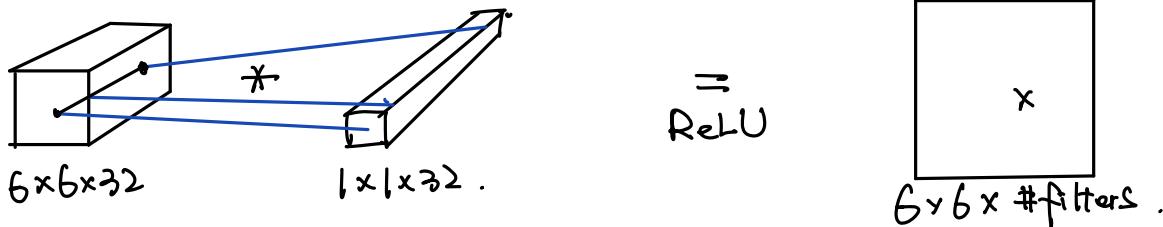
Easier for $a^{[l+2]} = a^{[l+1]}$. So adding this doesn't hurt.

Notice. this $g(z^{[l+2]} + a^{[l+1]})$ need $z^{[l+2]}$ and $a^{[l+1]}$ have same dim.

That's why in ResNet we see many conv. layer preserve dim.
 Otherwise. we may write $g(z^{[l+2]} + W_s a^{[l+1]})$. where W_s makes sure the dimensions are the same.

Network in Network and 1×1 Conv.

$$\begin{array}{r} 1 \ 2 \ 3 \\ 3 \ 5 \ 5 \\ 2 \ 1 \ 3 \end{array} * \begin{array}{r} 2 \\ = \end{array} \begin{array}{r} 2 \ 4 \ 6 \\ 6 \ 10 \ 10 \\ 4 \ 2 \ 6 \end{array}$$

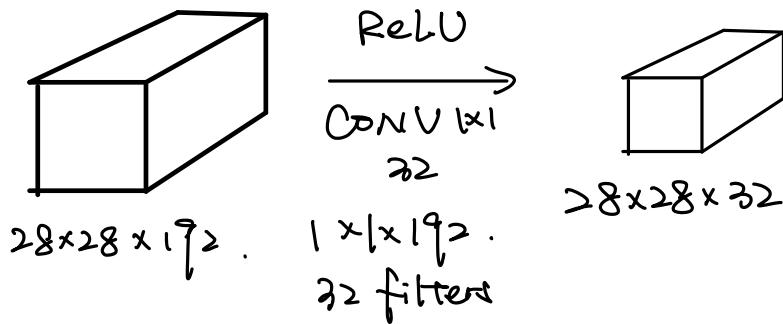


You take one "slice" of numbers, end up with single real number . with ReLU non-linearity.

It's like input 32 numbers with a network. and output to #filters dimension .

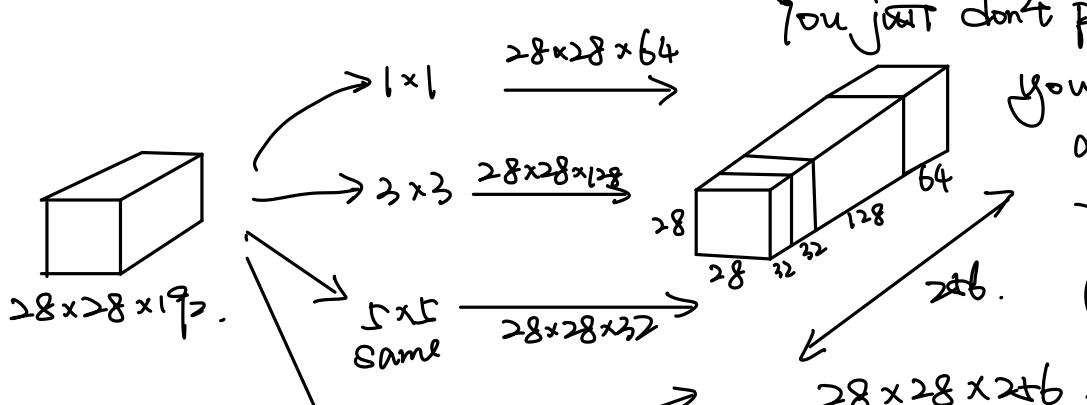


Using 1×1 convolutions .



You can shrink the network by doing the 1×1 conv.

Inception Network Motivation .

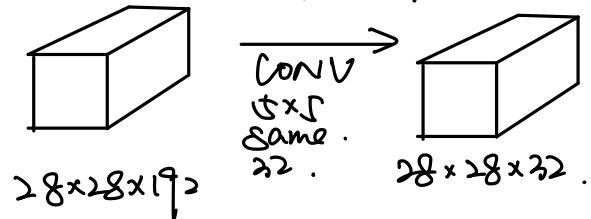


You just don't pick what filter/pooling you want, you just do them all. and stack all of them .

(Only problem: computation)

MaxPool
 $\frac{28 \times 28 \times 32}{4} = 28 \times 28 \times 32$
 Need padding

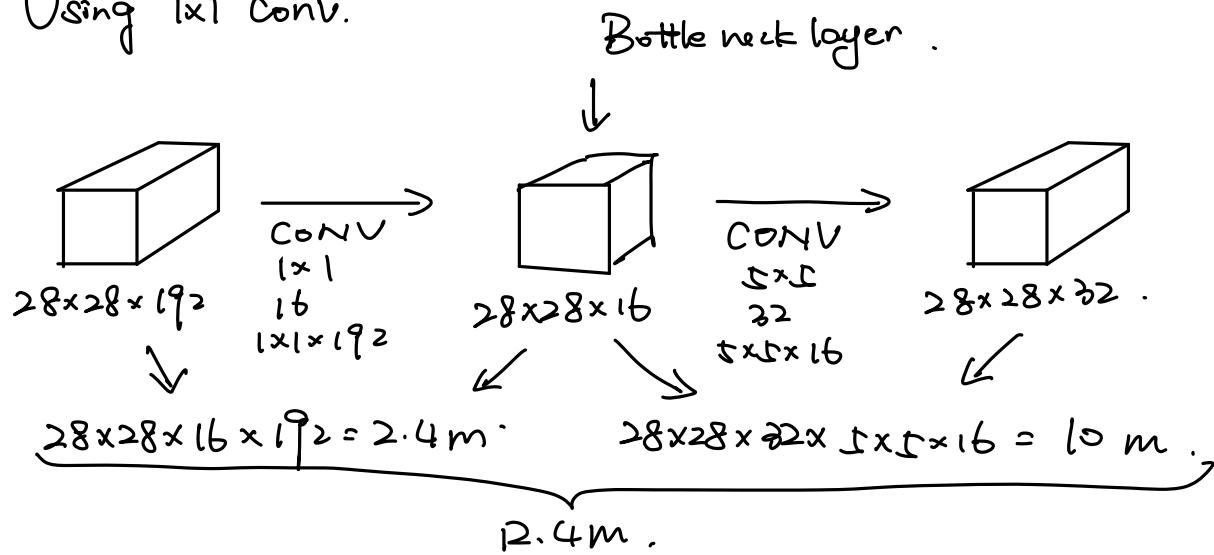
The problem of computation cost.



32 filters. filter are $5 \times 5 \times 192$.

$28 \times 28 \times 32 \times 5 \times 5 \times 192 = 120M$ of computations. (Costly!).

Using 1×1 conv.

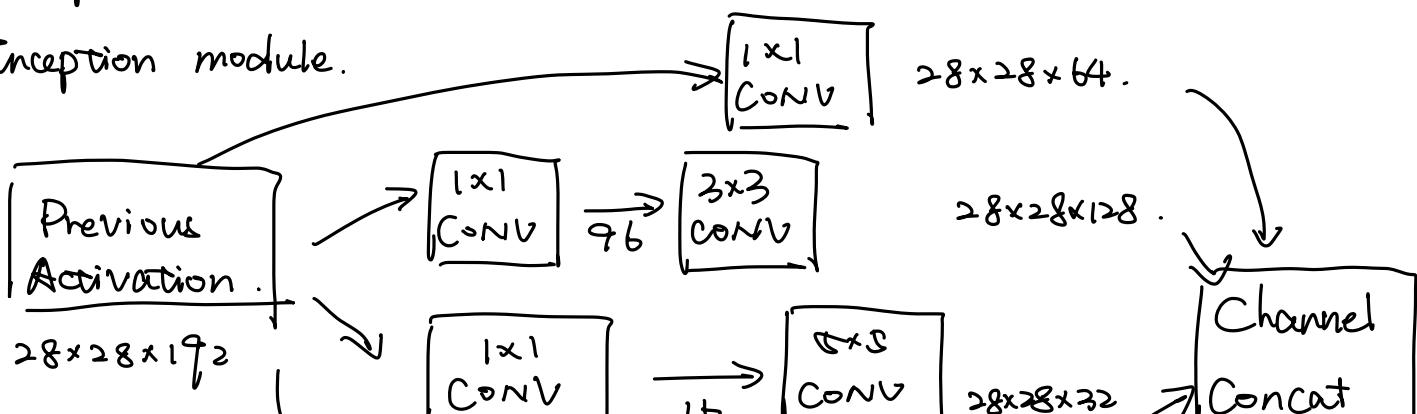


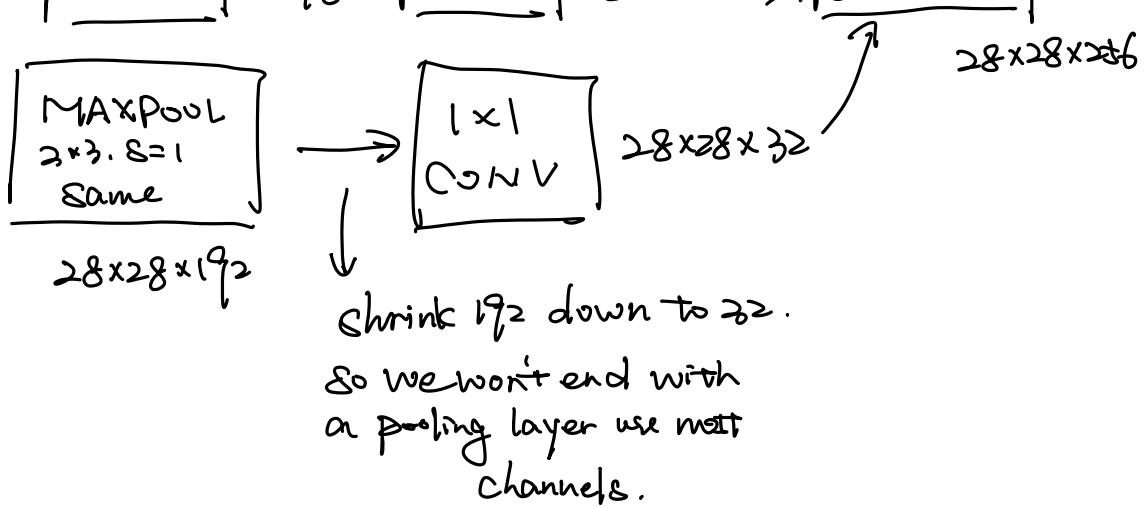
Reduce from $120M$ to $12.4M$. Approx. $\frac{1}{10}$.

Question: Will bottleneck hurt performance? Usually no.

Inception Network.

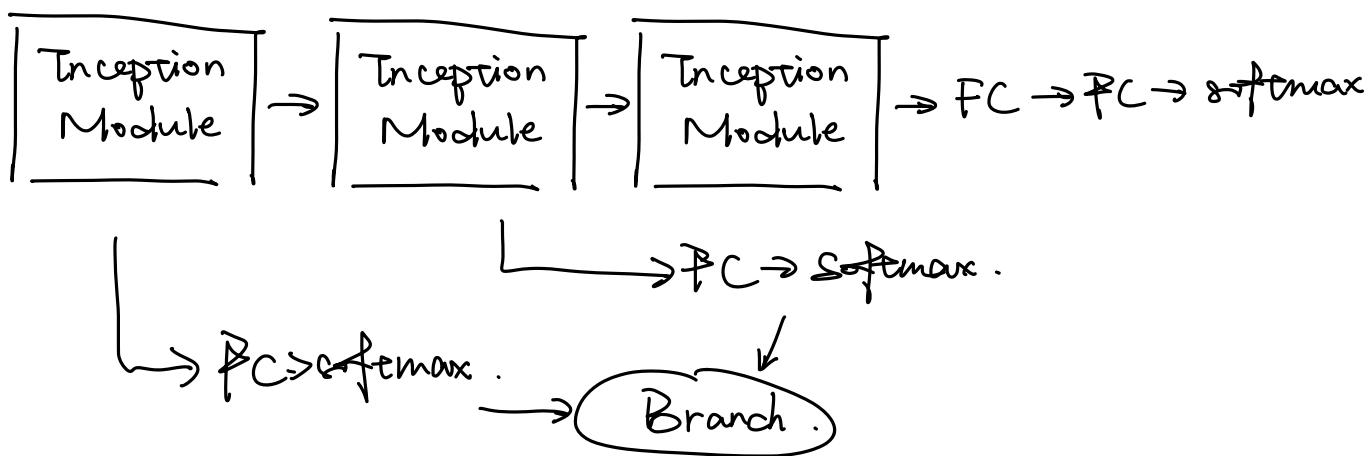
Inception module.





Basically, an inception network is a network with a lot of inception module / block.

Also, there are some hidden layers with branch, connect to FC and softmax to predict result.



Tips for doing well on benchmarks / winning competition.

① Ensembling

- Train several networks independently and average their outputs.
- * But it needs you to run on multiple networks. barely used for production deployed for actual user. (Also waste memory).

② Multi-crop at test time.

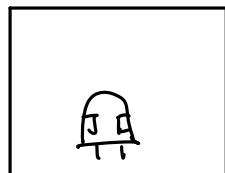
- Run classifier on multiple versions of test images and average results

Use open source code.

- Use architectures of networks published in the literature.
- Use open source code if possible.
- Use pre-trained model and fine-tune on your data.

Object localization.

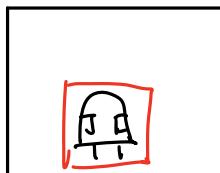
Image Classification .



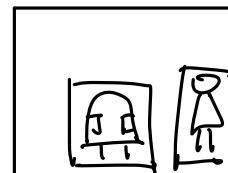
"Car"

1 Obj.

Classification with localization



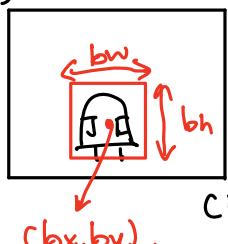
Detection .



Multiple Objs .

Classification with localization .

(b_x, b_y)



C1.1 .

ConvNet . \rightarrow softmax

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

* Classification part.

e.g. $b_x = 0.5, b_y = 0.7$

$b_h = 0.3, b_w = 0.4$.

↓ New 4 outputs. b_x, b_y, b_h, b_w . for bounding box .

Defining the target label y .

$$y = \begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \\ C_1 \\ C_2 \\ C_3 \end{bmatrix}$$

Bounding box.

Which object it is .

P_c is there an object. 0 no object. 1. has object .

If $P_c = 1$. Other value we just don't care .

The loss function.

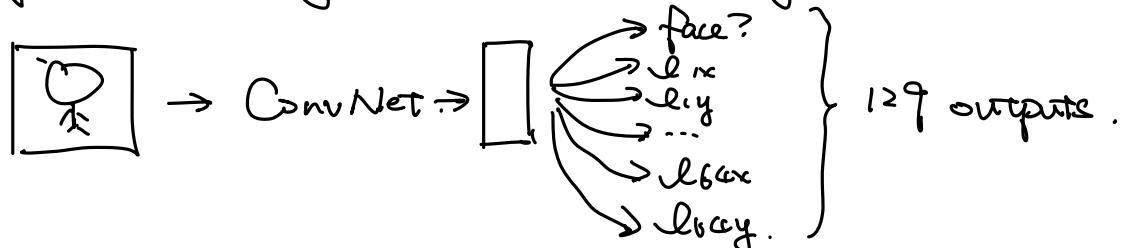
Using mean square : $\sum (\hat{y}_i - y_i)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_8 - y_8)^2$. If $y_i = 1$.
 $(\hat{y}_i - y_i)^2$ Others don't care. If $y=0$.

Of course you may use multiple loss functions for diff. parts.

Landmark detection.

e.g. Given an image of face. output the coord. for the 4 corner of eyes.

e.g. Points help you to land mouth, eyes, ...



Being able to detect these features. e.g. Social media apps can put "hats" "crown" on people head.

Object Detection.

Car detection example.

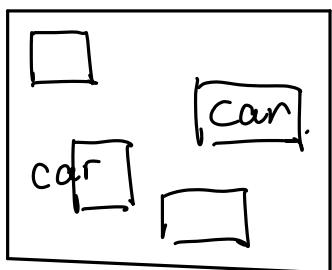
Training set

x	y	
[car]	1	[car] \rightarrow ConvNet \rightarrow y .
;	1	
	:	

Center cropped images.

"The object fill the image".

Then we use Sliding Window detection.



We just feed in the little square region into the ConvNet we trained. → classify whether 0 or 1.

Repeat with resized window (maybe larger).
Do this multiple times.



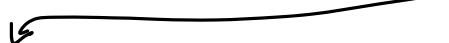
As long as the car is somewhere in the image.

Small square. You can show a square with a car on the image.

* Computational Cost.

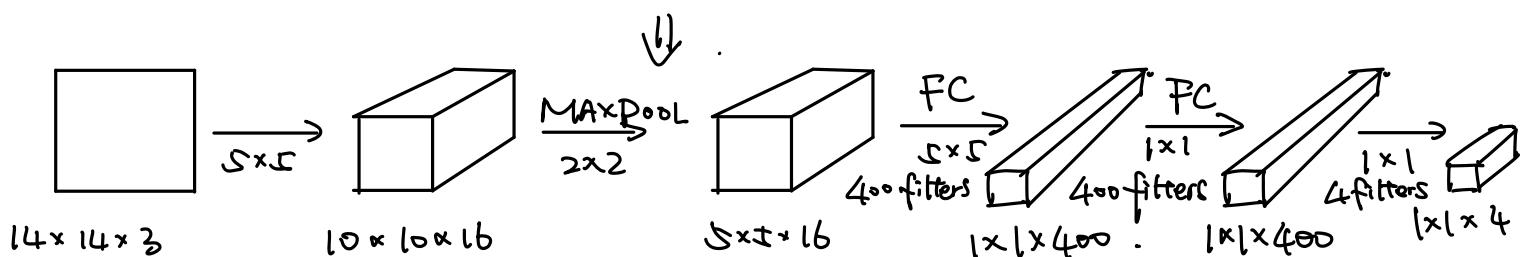
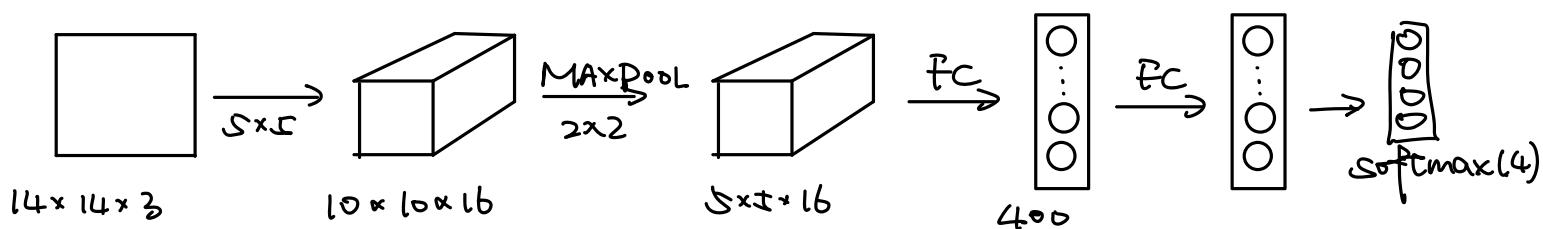
Previously. with simple linear classifier. OK.

Currently. DN. cost a lot. (so)



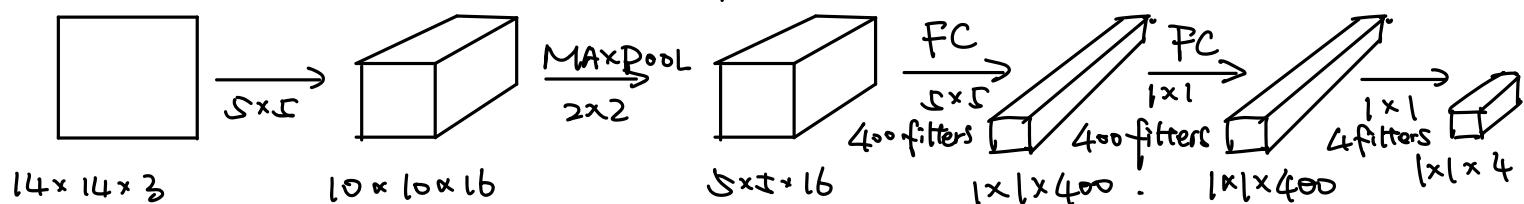
Convolutional implementation of sliding window.

Turning FC layers into convolutional layers.

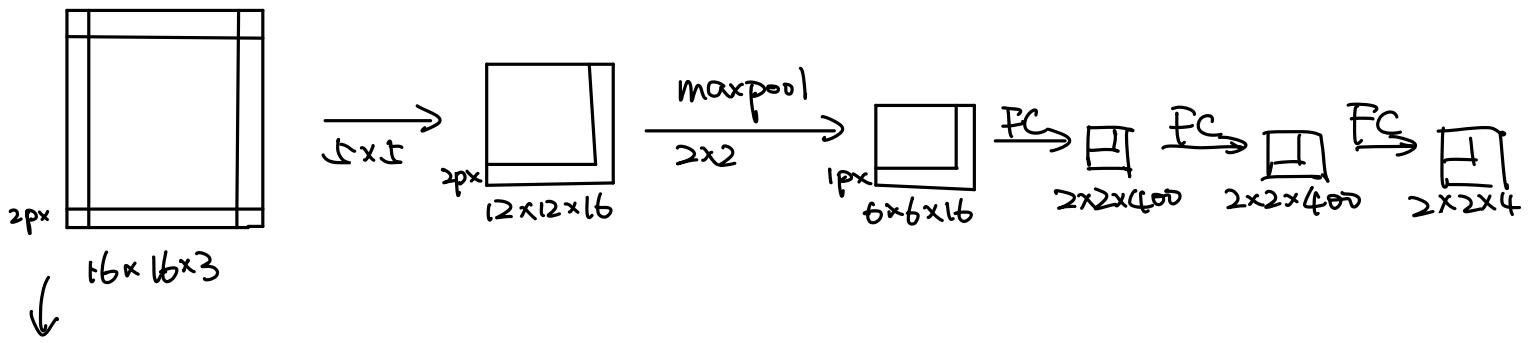


Mathematically. it is the same.

Convolution implementation of sliding window.



Tritting image:



it seems like we need instead. we directly put

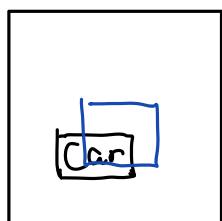
To do 4-times sliding \Rightarrow the image into the network.
window here.

Now consider the last output.

 The upper left gives you the result of upper left $14 \times 14 \times 3$ $2 \times 2 \times 4$. image result. Upper right corner, result of upper right $14 \times 14 \times 3$. The four sliding windows share a lot of computations. We don't need to compute them for multiple times.

Basically. instead of doing the 4 sliding windows sequentially. we just put the entire image into the CNN. and it will give all the results for all the windows.

Intersection Over Union. IoU.



 : Predict box

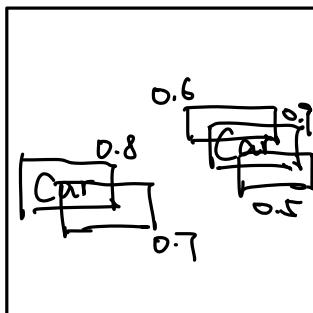
 : Ground truth box .

size of intersection

Intersection over union = size of union.

"Correct" if $\text{IoU} > 0.5$. You can use this to judge whether "May we O.S.O.B.O.T." the localization is correct.

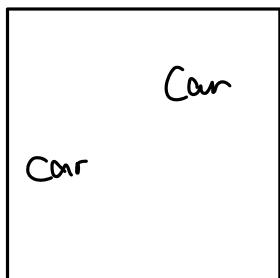
Non-max Supression.



You want each obj only being detect once in your algorithm. But by setting grid cell. Many cells may say they have large probability to contain the center of the object.

Pick the highest probability box, highlight

Probability may comes from IoU.



Each output prediction is:

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

Discard all boxes with $P_C \leq 0.6$.

19×19 .

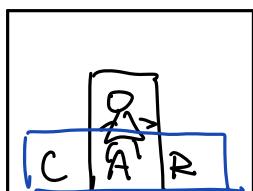
While there are any remaining boxes.

→ Pick the box with the largest P_C output
→ that as a prediction.

→ Discard any remaining box with $\text{IoU} \geq 0.5$ with
the box output in the previous step.

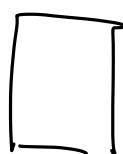
Anchor Boxes.

Overlapping Objects.



$$y = \begin{bmatrix} P_C \\ b_x \\ b_y \\ b_h \\ b_w \\ C_1 \end{bmatrix}$$

Anchor box 1 .



Anchor box 2



$$y = \begin{bmatrix} P_C \\ b_x \\ b_y \\ b_h \\ b_w \\ C_2 \\ C_3 \\ D_1 \end{bmatrix}$$

Anchor box 1 .

$$\begin{bmatrix} c_2 \\ c_3 \end{bmatrix}$$

$$\begin{bmatrix} b_x \\ b_y \\ \vdots \\ c_3 \end{bmatrix}$$

Anchor box 2.

Previously :

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

Output : $3 \times 3 \times 8$.

With two anchor boxes.

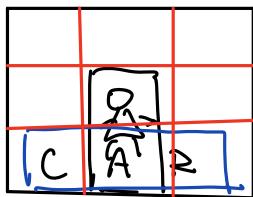
Each object in training

image is assigned to grid cell that contains that

obj's midpoint. and anchor box for the grid cell with highest IoU.

Output : $3 \times 3 \times 16 / 3 \times 3 \times 2 \times 8$.

example .



$y =$

P_c	1	Car only.	
b_x	b_x	0	
b_y	b_y	?	
b_h	b_h	?	
b_w	b_w	?	
C_1	1	?	
C_2	0	?	
C_3	0	?	
P_c	1	1	Anchor box 1
b_x	b_x	b_x	
b_y	b_y	b_y	
b_h	b_h	b_h	
b_w	b_w	b_w	
C_1	0	0	
C_2	1	1	Anchor box 2.
C_3	0	0	

of anchor boxes allows you how many objects can stack together.

Put it Together : YOLO Alg.

Region proposal : R-CNN .

Rather than running sliding window on every window, it use some image segmentation algorithm . And only run windows on some part of your image.

Face recognition

Face verification vs. face recog.

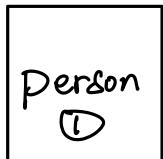
Verification:

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

Recognition:

- Has a database of K persons.
- Get an input image.
- Output ID if the image is any of the K persons
(Or "not recognized")

One Shot Learning



We need to learn from one example to recognize the person again.



Learning a "similarity" function.

$d(\text{img}_1, \text{img}_2)$ = degree of diff. between images.

If $d(\text{img}_1, \text{img}_2) \leq T$. (some threshold). Predict: same?

"diff." } -tion.

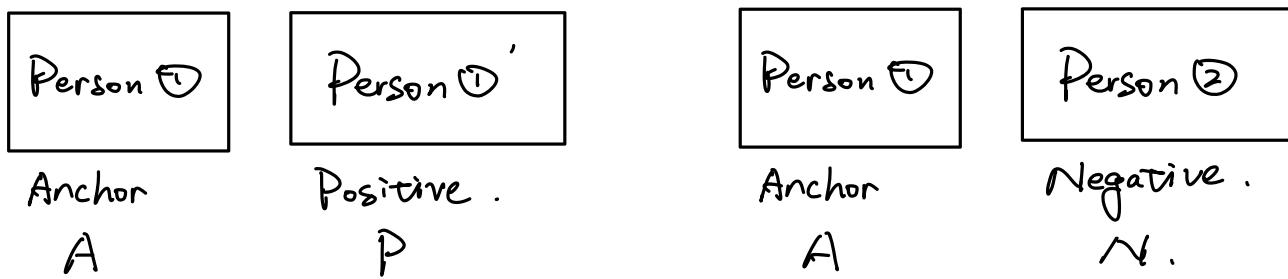
for recognition problem.

Just do pairwise comparison.

Pick the highest one larger than the threshold.

Triplet Loss

Learning Objective.



$$\text{Want: } \|f(A) - f(P)\|^2 \leq \|f(A) - f(N)\|^2 \\ d(A, P) \quad d(A, N).$$

$$\Rightarrow \|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 \leq 0.$$

But: $f(\text{img}) = \vec{0}$ if for all images. it also satisfy the equation.

So we need to modify: \rightarrow margin. (Like margin in SVM).
 $\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha \leq 0.$

Loss function.

Given three images A. P. N. (Triplet loss).

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

① As long it is smaller than 0. we get loss 0.

② As long it greater than 0. loss is pos.

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

Training set: 10k pictures of 1k persons. (Avg. 10 for each).

If only 1 for each maybe it's hard to train.

But after training, for new added person, may only have one pic

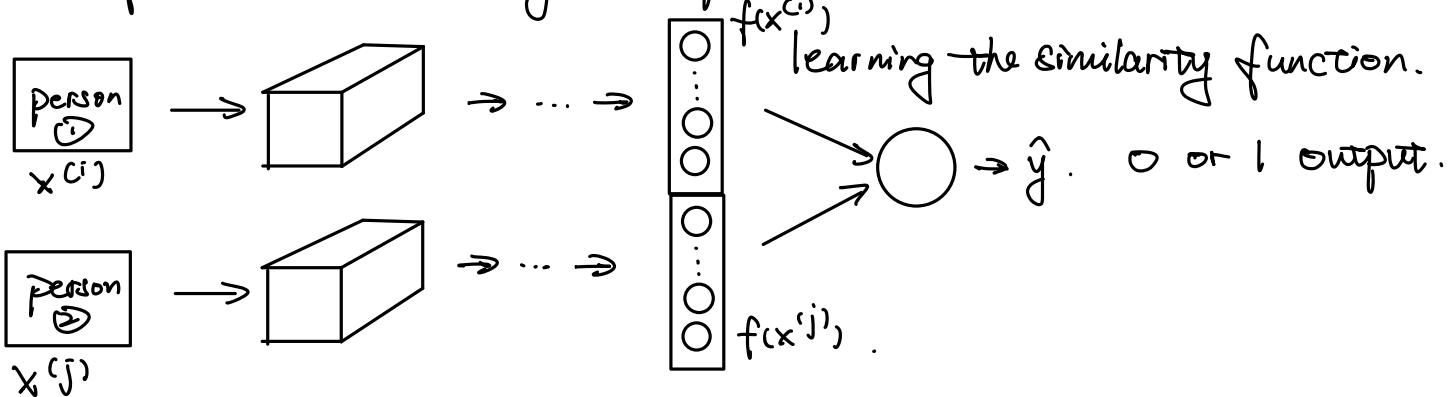
Choosing the triplet A.P.N.

If A.P.N chosen randomly, $d(A.P) + \alpha \leq d(A.N)$ is easily satisfied.

So we need to choose triplets that are hard to train on.

e.g. $d(A.P) \approx d(A.N)$.

Face verification and binary classification.

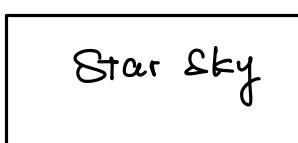
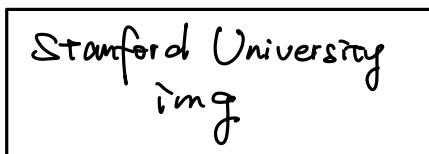


$$\hat{y} = \text{Sigmoid} \left(\sum_{k=1}^{128} w_k |f(x^{(i)})_k - f(x^{(j)})_k| + b \right).$$

Take the elementwise diff between two encodings.

Or. χ^2 diff:
$$\frac{(f(x^{(i)})_k - f(x^{(j)})_k)^2}{f(x^{(i)})_k + f(x^{(j)})_k}.$$

Neural Style Transfer.



Content Style



Stanford img
with Star sky

Generated image G .

Cost Function.

Neural style transfer cost func.

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

where C is the content img. S is the style img. G is the generated img.

Find the generated img G .

① Initiate G randomly. (White noise img).

$$G = 100 \times 100 \times 3$$

② Use G to minimize $J(G)$.

$$G := G - \frac{\partial}{\partial G} J(G) \quad \text{More and more like both imgs.}$$

Content Cost func.

Overall function : $J(G) = \alpha J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$.

→ Say you use hidden layer l to compute content cost.

→ Use pre-trained ConvNet (e.g. VGG).

→ Let $a^{[l]}(C)$ and $a^{[l]}(G)$ be the activation of layer l on the imgs.

→ If $a^{[l]}(C)$ and $a^{[l]}(G)$ are similar, both images have similar content.

$$J_{\text{content}}(C, G) = \frac{1}{2} \| a^{[l]}(C) - a^{[l]}(G) \|^2$$

Style Cost Func.

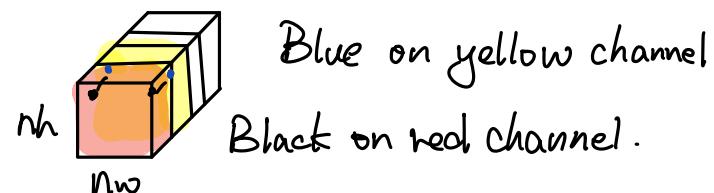
Meaning of "Style" of an image.

Saying you use layer l to measure "style".

Define style as correlation between activations across channels.

"How correlated are the activations across different channels?".

e.g. Suppose we have 5 channels.



Blue on yellow channel

Black on red channel.

We take each pair dots on different channel.

Look across all the pairs on nh and nw. look how correlated they are.

The correlation gives you one way to measure how often diff features occur or doesn't occur together.

Style Matrix:

let $a_{i,j,k}^{[e]} = \text{activation at } (i, j, k)$. $G^{[e] \times [e]}$ is $n_h^{[e]} \times n_w^{[e]}$.

$$G_{kk'}^{[e] \times [e]} = \sum_i^{n_h^{[e]}} \sum_j^{n_w^{[e]}} a_{ijk}^{[e]} a_{ijk'}^{[e]}$$

Two different activation on different channels but in same pos.

$$G_{kk'}^{[e] \times [e]} = \sum_i^{n_h^{[e]}} \sum_j^{n_w^{[e]}} a_{ijk}^{[e]} a_{ijk'}^{[e]}$$

e.g. if we have 5 channels. it will be a 5×5 symmetric matrix. (Gram matrix).

$$\text{Jstyle}(S, G) = \sum_k \sum_{k'} \left(\frac{\|G_{kk'}^{[e]} - G_{kk'}^{[e] \times [e]}\|_F^2}{\|G_{kk'}^{[e] \times [e]}\|_F^2} \right)^2$$

Also can add norm term $\frac{1}{(\sum h^{[e]} w^{[e]} c^{[e]})^2}$

$$\text{Jstyle}(S, G) = \sum_e \lambda^{[e]} \text{Jstyle}(S, G)$$

$\lambda^{(j)}$ to allow you to pick which layer you want to use. And how to use it.

To sum up.

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

1D and 3D generalizations of model.

Now we learned how to process 2D images. (May have multi channels)
e.g. 1D signal.

$$\begin{matrix} \text{14x1} \\ \text{e.g. 14 dim} \end{matrix} * \begin{matrix} \text{5x1} \\ \text{5 dim.} \end{matrix} \Rightarrow \begin{matrix} \text{10x1} \\ \text{Also if you have multiple channel.} \end{matrix}$$

3D data. like CT of brain. (Multiple slices).

$$\begin{matrix} \text{14x14x14} \\ \text{3D volume} \end{matrix} * \begin{matrix} \text{5x5x1} \\ \text{3D filter.} \end{matrix} \Rightarrow \begin{matrix} \text{10x10x10} \\ \text{This is one channel. May multi channel.} \end{matrix}$$

like. If we have 16 filter. : $10 \times 10 \times 10 \times 16$.

$* 5 \times 5 \times 5 \times 16$. 32 filters.

$$\Rightarrow 6 \times 6 \times 6 \times 32$$

e.g. 3D volume: Movie data. Medical data.