

CONVOLUTIONAL NEURAL NETWORKS

* Computer Vision

Problems: Image classification; object detection,
Neural style transfer

* Edge detection

~~vertical~~ vertical edge detection

Consider a 6×6 grayscale image

$$\begin{array}{|c|c|c|c|c|} \hline
 3 & 0 & 1 & 2 & 7 & 4 \\ \hline
 1 & 5 & 8 & 9 & 3 & 1 \\ \hline
 2 & 7 & 2 & 5 & 1 & 3 \\ \hline
 \end{array}
 \quad *
 \quad
 \begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array}
 = \quad
 \begin{array}{|c|} \hline
 -5 \\ \hline
 \end{array}$$

6×6 3×3 4×4
 filter
 convolution operator
 ✓

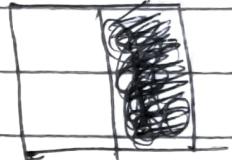
element-wise multiplication with the filter
ie the $(1,1)$ th element of resultant
will be

$$\begin{aligned}
 (3 \times 1) &+ (1 \times 1) + (6 \times 1) + (0 \times 0) + (0 \times 5) + (0 \times 1) \\
 &+ (1 \times -1) + (8 \times -1) + (2 \times -1) \\
 = -5
 \end{aligned}$$

$A =$

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

Consider image A it should look like



$$\text{filter } F = \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$$

$$A * F = \begin{bmatrix} 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \end{bmatrix}$$

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

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

: : :

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

sobel filter:

1	0	-1
2	0	-2
1	0	-1

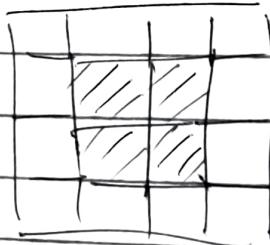
Schorr filter

3	0	3
10	0	10
3	0	3

→ # Padding

$$\begin{array}{ccc} 6 \times 6 & \times & 3 \times 3 \\ n \times n & \times & f \times f \end{array} = \frac{4 \times 4}{(n-f+1) \times (n-f+1)}$$

to keep the dimensions of the resultant image the same as the original image, add a border to it



here padding $p = 1$

resultant will be

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

valid conv: no filter padding

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

$$n + 2p - f + 1 = n$$

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

by convention f is odd

Strided Conv

Consider a 7×7 input, 3×3 filter
 now instead of moving our grid window
 by 1, we move it by 2.

$$\begin{array}{ccccccc}
 2 & 3 & 7 & 4 & 6 & 2 & 9 \\
 6 & 6 & 9 & 7 & 8 & 4 & 3 \\
 3 & 4 & 8 & 3 & 3 & 9 & 7 \\
 \cdot & \cdot & \cdot & & & 2 & 3 & 7 \\
 \cdot & \cdot & \cdot & & & 6 & 6 & 9 \\
 \cdot & \cdot & \cdot & & & 3 & 4 & 8 \\
 \cdot & \cdot & \cdot & & & & &
 \end{array} \quad \text{filter in } 1,1$$

in iteration 2:

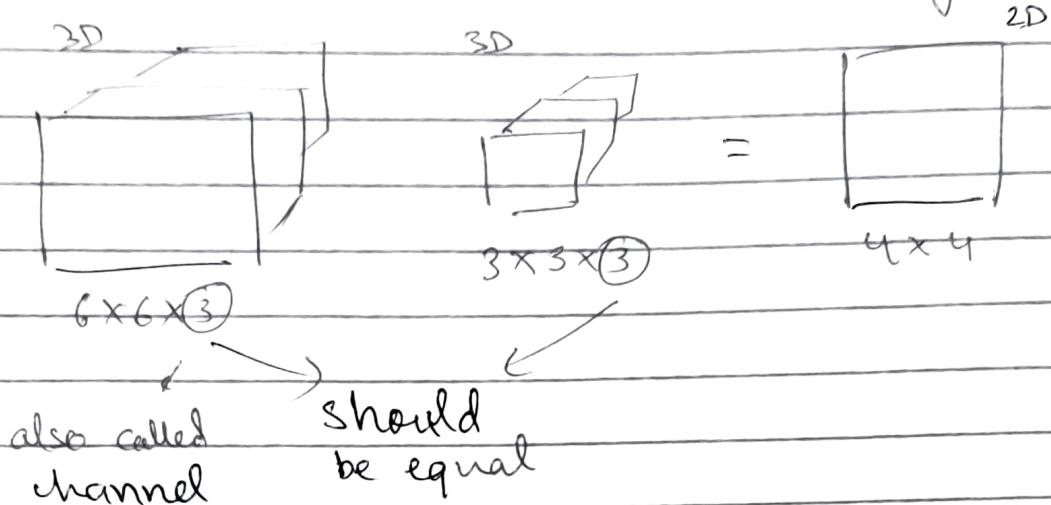
$$\begin{array}{ccc}
 7 & 4 & 6 \\
 9 & 7 & 8 \\
 8 & 3 & 3
 \end{array}$$

Input: 7×7 filter: 3×3 output: 3×3
 Here stride = 2

input $n \times n$ with filter $f \times f$ with
 with padding p & stride s
 output:

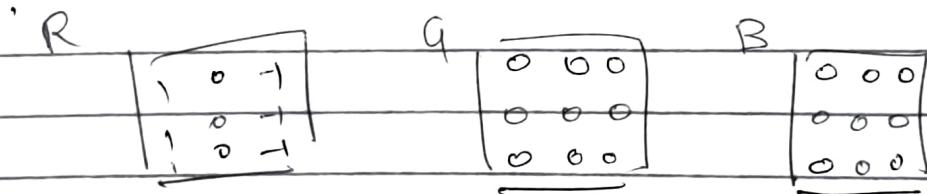
$$\left\lfloor \frac{n+2p-f+1}{s} \right\rfloor \times \left\lfloor \frac{n+2p-f+1}{s} \right\rfloor$$

4 Convolutions over volumes / RGB images

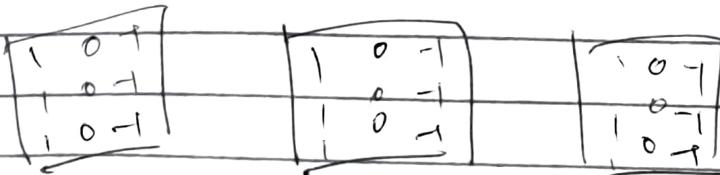


do element wise multiplication on the
cube window

if you want vertical edge detection on only
the red channel



if it ~~do~~ color doesn't matter



If you want to detect both vertical & horizontal
edge

6X6X3

* ^{vertical} $3 \times 3 \times 3 = 4 \times 4$

$\rightarrow 4 \times 4 \times 2$

horizontal

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

\rightarrow

One layer of CNN

$6 \times 6 \times 3$

$$* 3 \times 3 \times 3 \xrightarrow{\text{ReLU}(4 \times 4 + b_1)} \rightarrow 4 \times 4$$

$\rightarrow \text{TR}$

$$* 3 \times 3 \times 3 \xrightarrow{\text{ReLU}(4 \times 4 + b_2)} \rightarrow 4 \times 4$$

Notations:

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

$p^{[l]}$ = padding

$s^{[l]}$ = stride

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

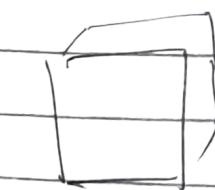
Th

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

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

A simple ConvNet

Consider:



$39 \times 39 \times 3$

$$\begin{aligned} n_H^{[0]} &= n_W^{[0]} = 39 \\ n_C^{[0]} &= 3 \end{aligned}$$

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

$37 \times 37 \times 10$

10 filters

20 filters

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

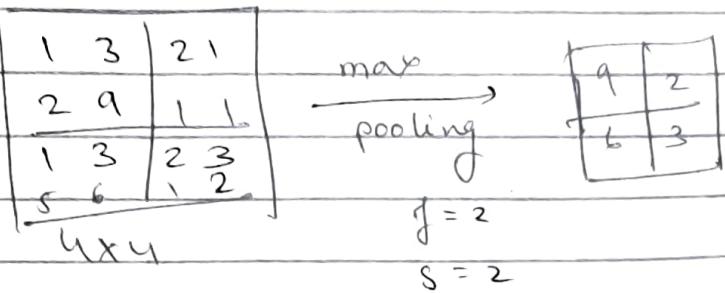
$17 \times 17 \times 20$

$$\begin{aligned} f^{[3]} &= 5 \\ s^{[3]} &= 2 \\ \text{filters} &= 40 \end{aligned}$$

$1 \times 1 \times 40$

logistic softmax

Pooling Layers

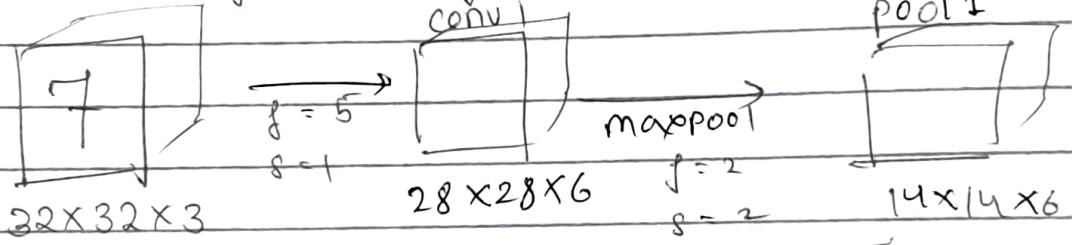


any pooling, min pooling
padding is rarely used in pooling.

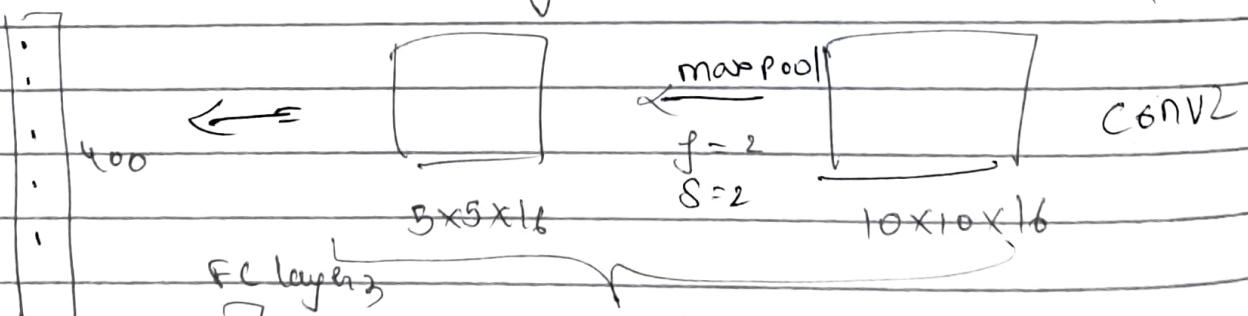
CNN Example

LeNet - 5

to recognize digits



layer 1

 $f = 5$
 $s = 1$ 

FC layers

fully connected layer

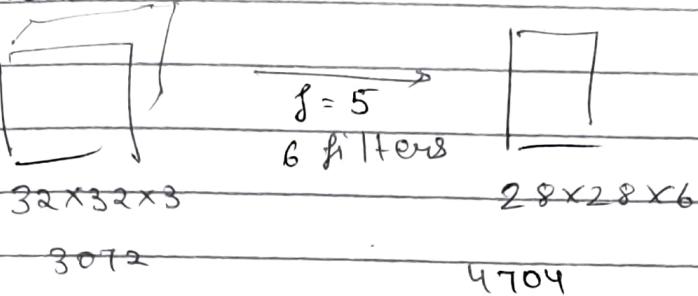
layer 2

Fc4

Softmax
(10 outputs)

Why Convolutions

Consider



if it were a fully connected layer, we would need $3072 \times 4704 \approx 14 \text{ million parameters}$

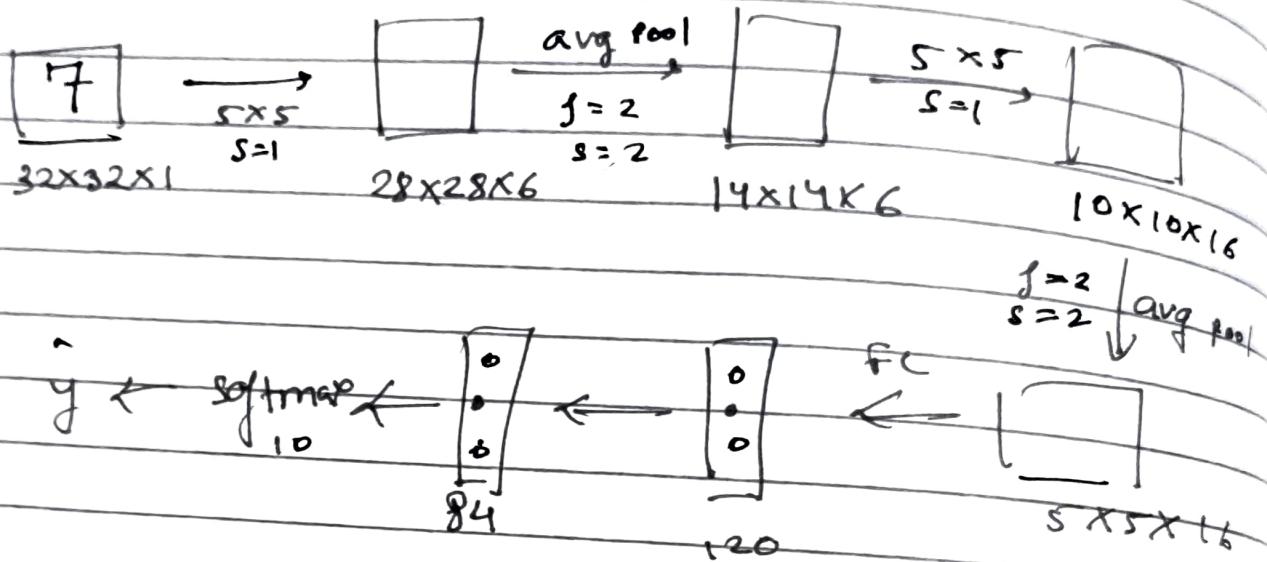
but from conv we have $(5 \times 5 + 1) \times 6 = 156$ params

Why convs get away w so few parameters?

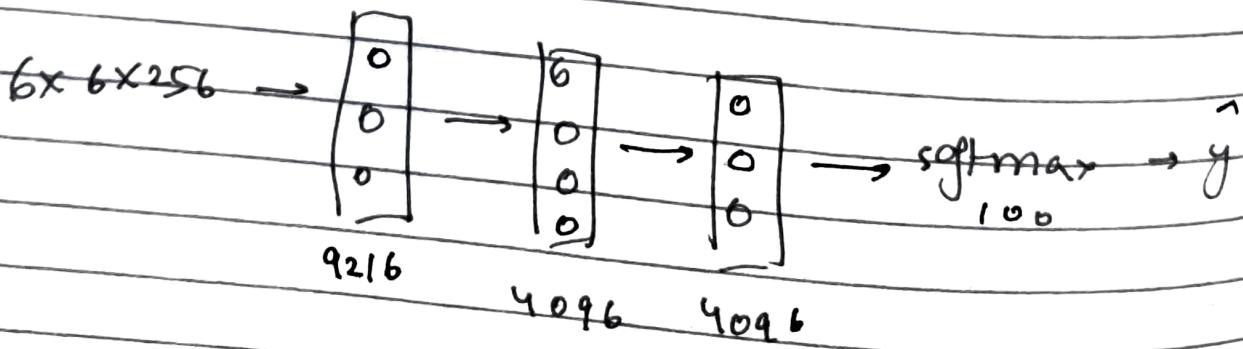
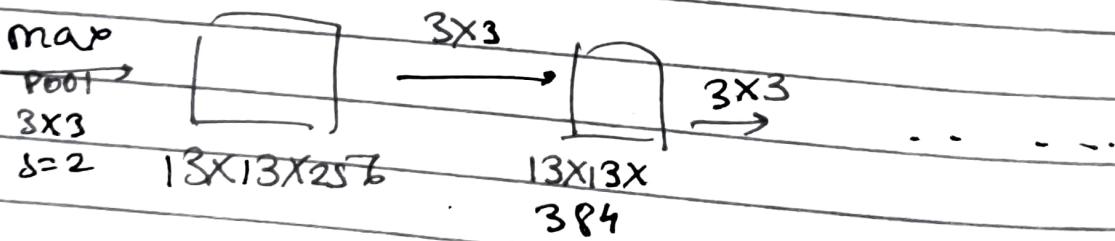
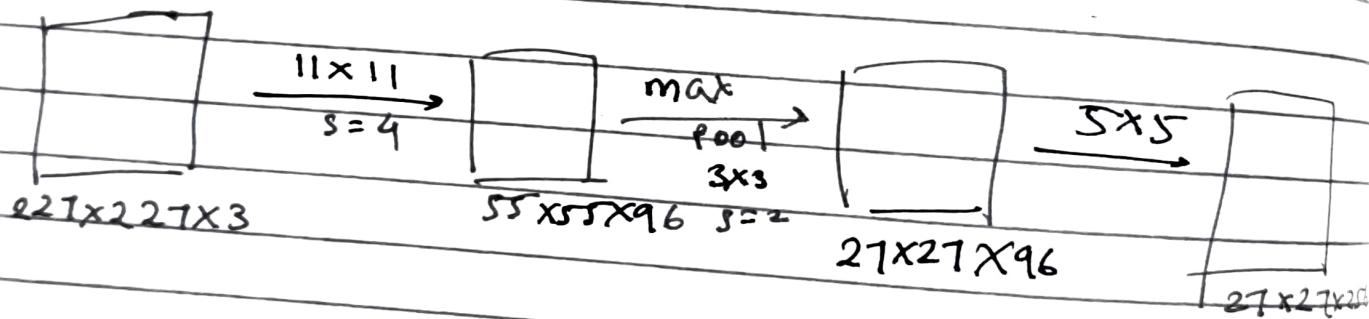
- 1) Parameter Sharing: A feature detector (e.g. vertical edge detector) will mostly be usable in other part of the image.
- 2) Sparsity of connections: In each layer, an output is dependent on a small no. of input features.

Classical Networks

* LeNet 5



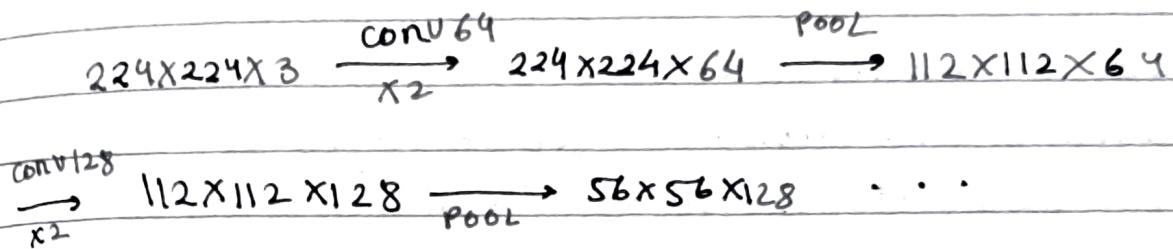
* AlexNet



* VGG16

conv: 3x3 filters s=1

maxpool: 2x2, s=2



Residual Networks (ResNets)

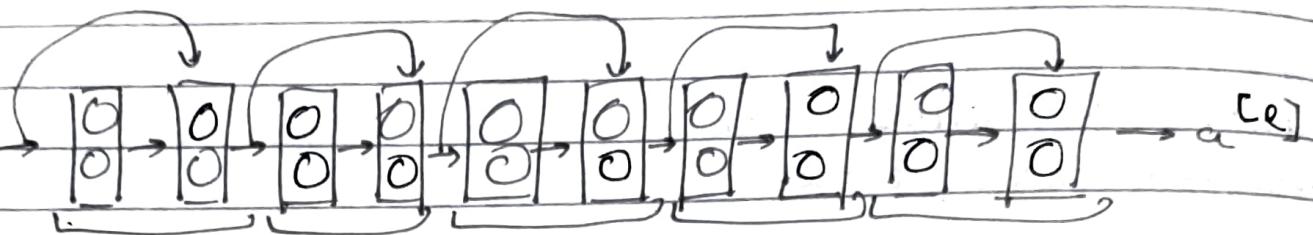
* residual block

The diagram shows a residual block. The forward pass consists of three main stages: 1) An input $a^{[l]}$ is processed by a linear layer and then a ReLU activation to produce $a^{[l+1]}$. 2) $a^{[l+1]}$ is processed by another linear layer and then a ReLU activation to produce $a^{[l+2]}$. 3) A skip connection branches off before the first ReLU and adds the original input $a^{[l]}$ to the output of the second stage, $a^{[l+2]}$.

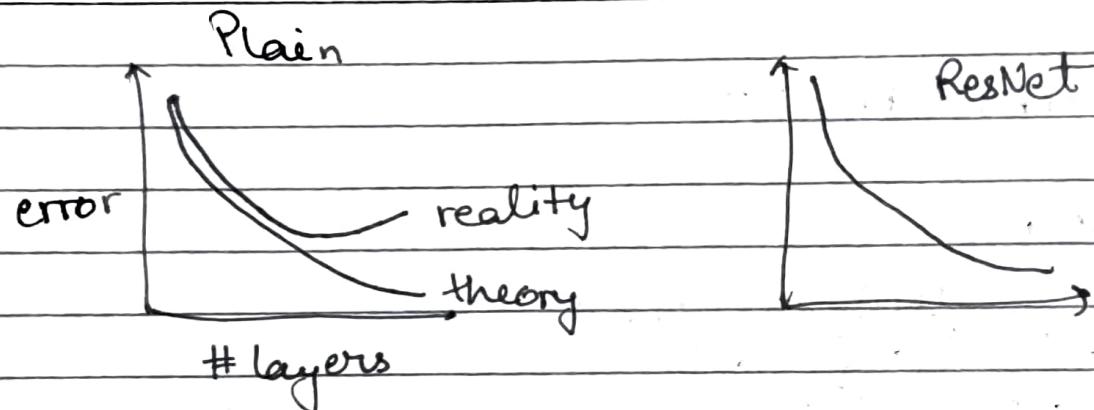
$$z^{[l+2]} = w_a^{[l+1]} z^{[l]} + b^{[l+2]}$$

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

$$g(a^{[l+2]}) = \underbrace{g(z^{[l+2]})}_{\times \times} + \underbrace{g(z^{[l+2]}) + \tilde{a}^{[l]}}_{\text{residual block}}$$

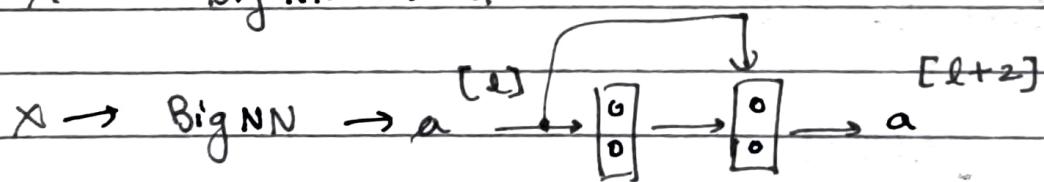


has 5 residual blocks
 without the residual connections,
 the network is called "plain"



Why ResNets work

$$x \rightarrow \text{Big NN} \rightarrow a^{[e]}$$



activation : ReLU $a > 0$.

$$a^{[l+2]} = g(a^{[l+2]} + a^{[l]})$$

$$= g(w^{[l+2]} \underbrace{a^{[l+1]}}_x + b^{[l+2]} + a^{[l]}) = g(a^{[l]})$$

$$\text{if } w^{[l+2]} = b^{[l+2]} = 0$$

$$= a^{[l]}$$

1x1 Convolutions

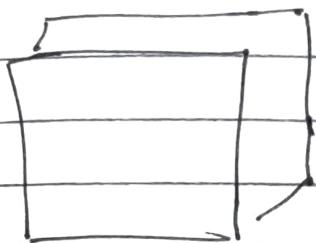
$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$$

$3 \times 3 \times 1$

\star

$$\begin{bmatrix} 2 & 4 & 6 \\ 8 & 10 & 12 \\ 14 & 16 & 18 \end{bmatrix}$$

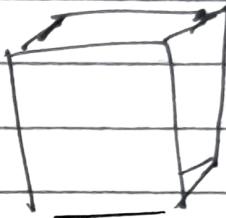
is basically
element wise
multiplication
in this case



\star



=
ReLU



$6 \times 6 \times 32$

$1 \times 1 \times 32$

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

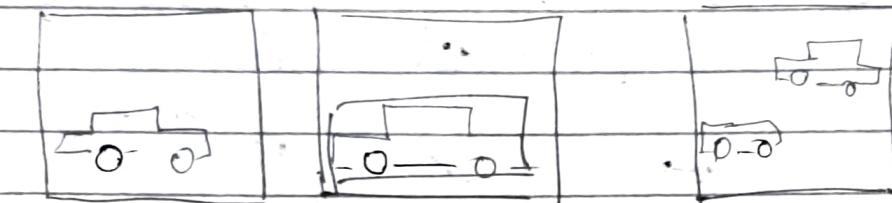
← inception network →

Practical Advice

- * Use open source CNN implementations available on github
- * Use weights trained by other people for your network. (transfer learning)
- * Augment your data to increase the dataset size and hence improve performance
 - mirroring
 - random cropping
 - rotating
 - color shifting ($R+20, G+30, B+20$)

Object Detection

* Object Localization



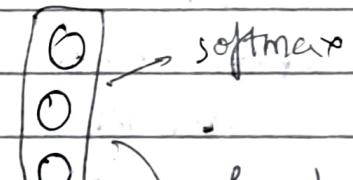
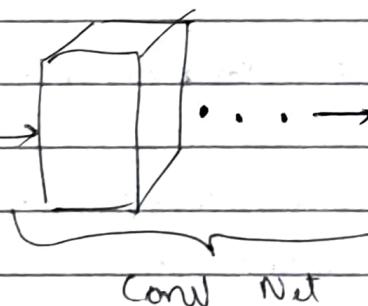
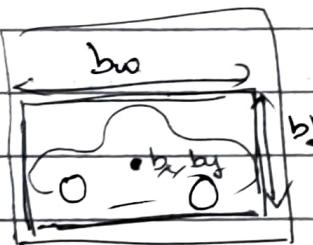
"car"
image
classification

"car"
image
classification
with localization

Detection

object

multiple objects



softmax
for local
bbox by b_x b_y b_w
bounding
box

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

$y =$	p_c	is there any object?
	b _x	
	b _y	
	b _w	
	C ₁	
	C ₂	
	C ₃	
	...	

$$f(\hat{y}, y) = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 \dots (\hat{y}_8 - y_8)^2$$

if $y = 1$

$$= (\hat{y}_1 - y_1)^2 \quad \text{if } y_1 = 0$$

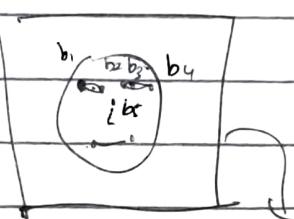
* Landmark Detection

"face detect"

has 64 "landmarks"

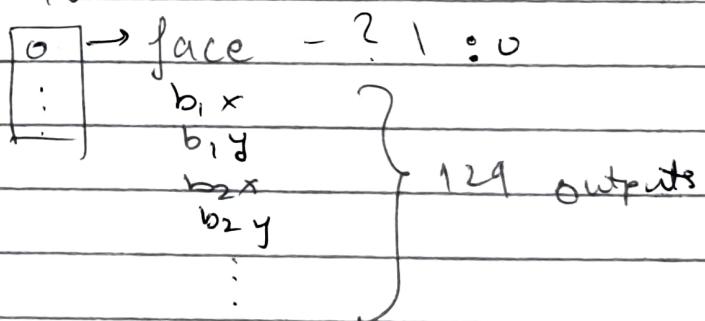
(single points)

corners of the eyes,
nose locators, lips etc



convnet

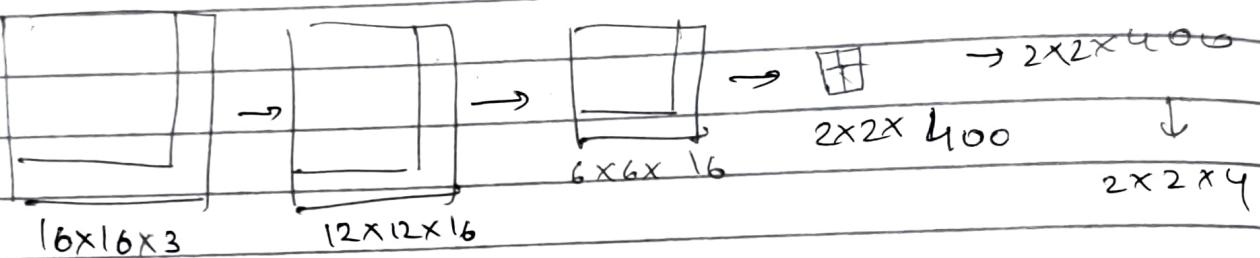
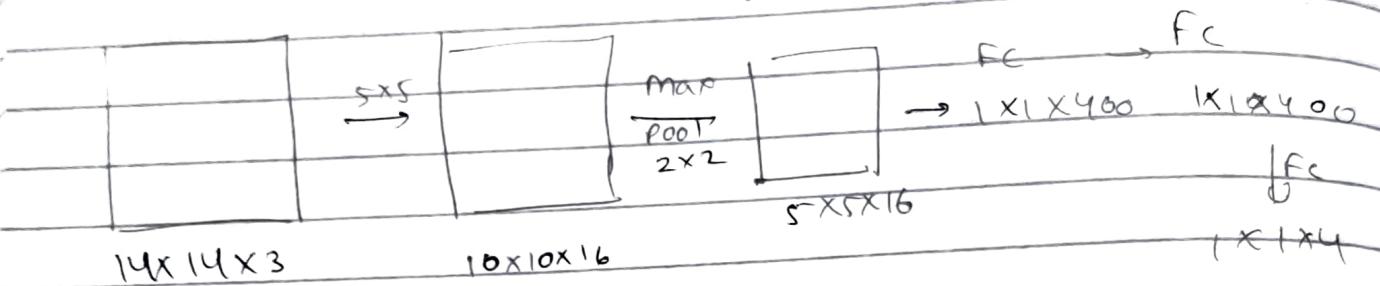
Used in AR



* Object Detect "

- Train the model using a dataset of closely cropped object images / bgs.
- Use a sliding window grid over the larger test image to localize the object.
- change the size of window & iterate again
- has high computational cost.

* Conv implementation of Sliding window.



user strides

* Bounding box predictions

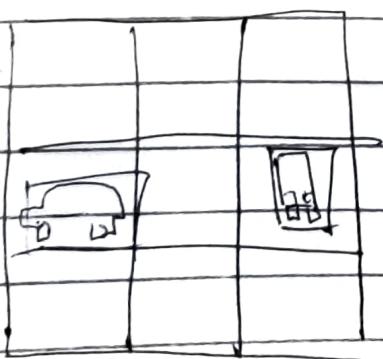
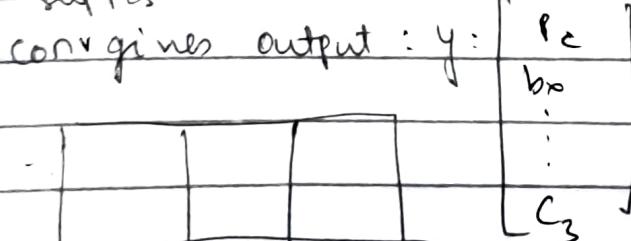
the ~~size~~ sliding windows from the convolution strides may not

fit the objects accurately

- Use You only Look Once algo (YOLO) algorithm

- Divide the test image to $n \times n$ grids
apply conv to each grid

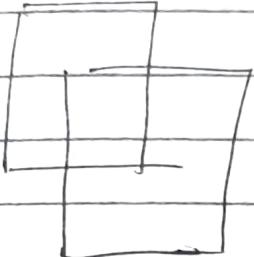
Suppose



→ now returns a $3 \times 3 \times 8$ output

→ Intersection Over Union

Let's say there are 2 boundary regions who which localized the object



Intersection over Union (IoU)

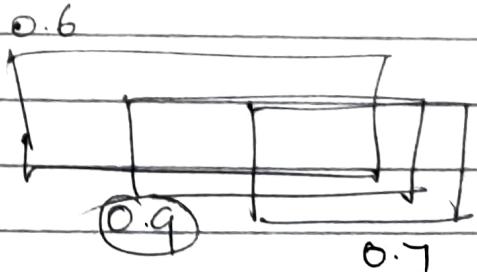
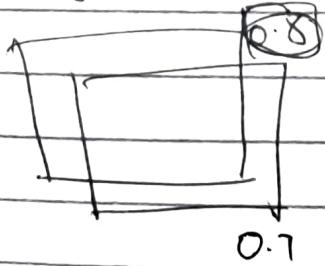
$$= \frac{\text{Size (area) of intersect region}}{\text{Size (area) of union region}}$$

correct if $\text{IoU} \geq 0.5$ usually

more generally, it is the measure of overlap of two boxes

* Non Max Suppression:

What if the object appears on more than one grid. There should be
The count of object should not be incremented for each appearance
of the object



the higher pc grid is kept

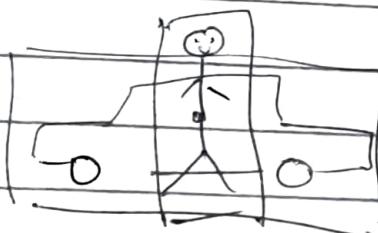
Consider the image divided into 19×19 grid
each output prediction gives

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

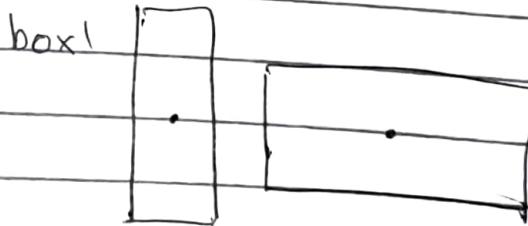
- discard all boxes with $p_c \leq 0.6$
- pick the box with highest p_c , output that as predict
- discard any box with $\log \frac{b_w}{100} \geq 0.5$
with the box from the previous step

* Anchor box

What if we have multiple overlapping objects.



we have two anchor boxes:



Output earlier: $3 \times 3 \times 8$

now: $3 \times 3 \times 16$

* Putting Together Yolo

= Face Recognition

Face recognition vs

Face verification

input: image

image, name/ID

output id of person in image/
not recognized

whether the input image
is of the claimed person

* One Shot learning

algo gets only 1 image per person

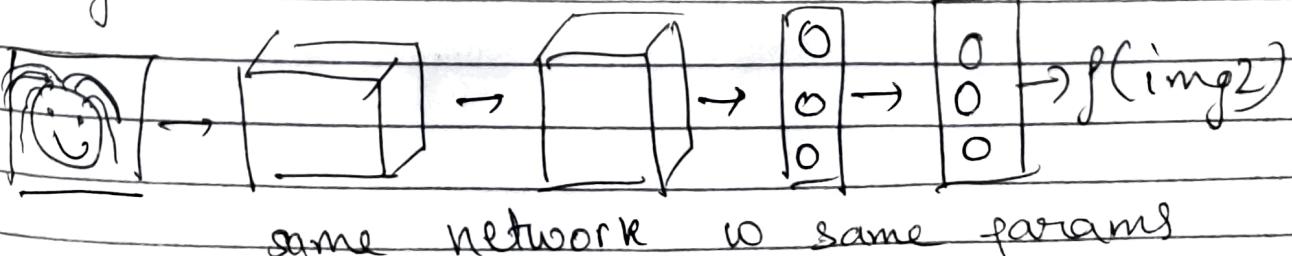
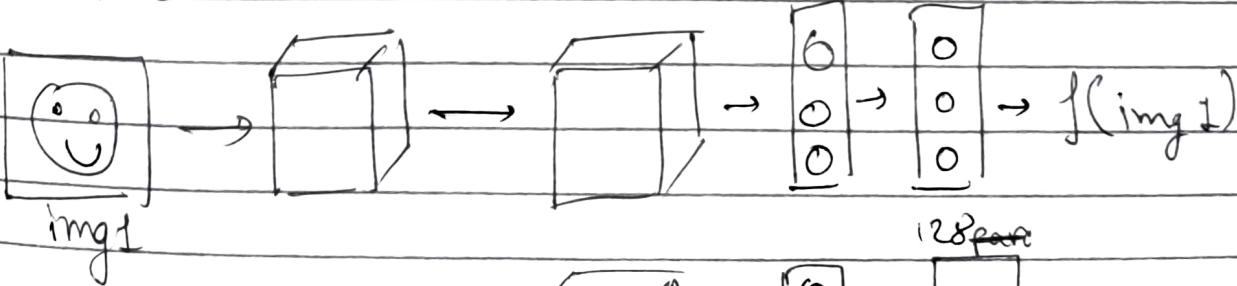
conv softmax doesn't work very well in
this situation

The model has to learn a 'similarity' fn
 $d(img_1, img_2) = \text{degree of diff b/w images}$

$d(img_1, img_2) \leq \tau$ "same"

$d(img_1, img_2) > \tau$ "not same"

* Siamese Network



$$d(\text{img1}, \text{img2}) = \|f(\text{img1}) - f(\text{img2})\|_2^2$$

Learn parameters s.t.

if img1 & img2 are same person,

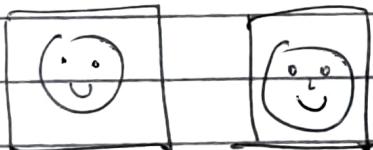
$\|f(\text{img1}) - f(\text{img2})\|$ is small

if img1 & img2 are diff person,

$\|f(\text{img1}) - f(\text{img2})\|$ is large

* Triplet Loss fn

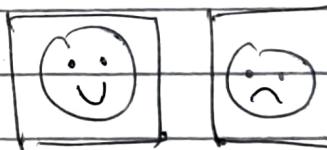
Learning objective:



anchor positive

~~A~~ P

same person



anchor negative

A N

diff person

$$\text{want } \|f(A) - f(P)\| \leq \|f(A) - f(N)\|$$

Or often written

$$\|f(A) - f(P)\| - \|f(A) - f(N)\| + \underline{\text{margin}} \leq 0$$

Given 3 images A, P, N

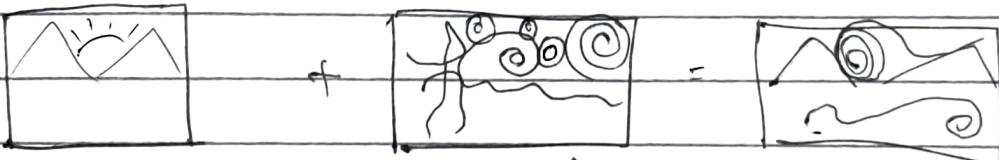
$$L(A, P, N) = \max \left(\|f(A) - f(P)\|_1 + \|f(A) - f(N)\|_1 + \alpha, 0 \right)$$

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

Choose A & corresponding P don't choose N randomly for training

Choose triplets that are hard to train on
i.e N should can be the image A
closely resembles to.

Neural Style Transfer



Landscape

Content(C)

Van gogh

Style(S)

generated image

(a)

kinda like what prima did back in the day

* What are deep convnets learning here
As we move towards deeper layers,
they detect more complex structures ~~structures~~

* Cost fn

define a cost fn J st.

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

To find generated image G

1. Initialize G randomly

$$G: 100 \times 100 \times 3$$

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

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

* The Content cost fn
i.e.

$$J_{content}(G) = \alpha J_{content}$$

- Say you use hidden layer l to compute content cost

- Use pretrained Convnet (VGG network)

- Let $a^{(c)}$ & $a^{(G)}$ be the activation of layer l on the images

- If $a^{(c)}$ & $a^{(G)}$ are similar, both have similar content.

$$J_{content}(C, G) = \frac{1}{2} \| a^{(c)} - a^{(G)} \|_2^2$$

* Style Cost fn

Any Take any intermediate conv layer volume from the net

It has multiple channels > 3

They store more complex properties of the image like vertical, edges, circular patterns etc

(let

$a_{ijk}^{(l)}$ activation at (i, j, k) . $g^{(l)}$ is $n_l \times n_k$

$$G_{kk}^{(l)} = \sum_i \sum_j a_{ijk}^{(l)} a_{ijk}^{(l)}$$

$$G_{kk}^{(l)} = \sum_i \sum_j a_{ijk}^{(l)} a_{ijk}^{(l)}$$

$$J_{\text{style}}^{(l)}(s, a) = \ln G_{kk}^{(l)} - G_{kk}^{(l)} \pi_f^2$$

$$= \frac{1}{(\frac{n_l n_k}{2})} \sum_i \sum_j (G_{kk}^{(l)} - G_{kk}^{(l)})$$