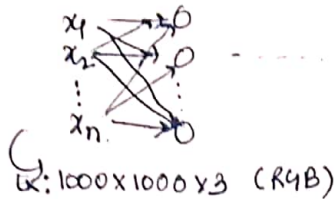


Convolutional Neural Networks:

* Computer vision:

- inputs can be really big.



Ex: $1000 \times 1000 \times 3$ (RGB)

* Edge detection example:

~~edge detection~~

image \rightarrow vertical edges
 \rightarrow horizontal edges

Ex: 6×6 grayscale img

derivative method

3	1	0	1	2	7	4
1	5	0	8	9	3	1
2	7	0	2	5	1	3
0	1	3	1	7	8	
4	2	1	6	2	8	
2	4	5	2	3	9	

6x6

"convolution"
 \ast (overloaded notation)
 $\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$ = 3x3 filter.
 to detect vertical img
 (also called kernel)

-5	-4	0	8
-10	-2	2	3
0	-2	-4	-7
-3	-2	-3	-16

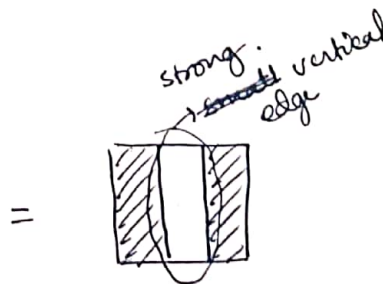
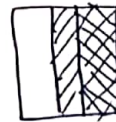
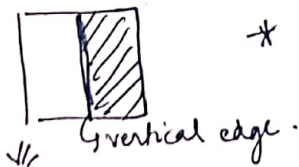
4x4

\rightarrow vertical edge detectors:

Python: conv-forward.

tensorflow: tf.nn.conv2d

Ex:



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

\ast

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

=

0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0

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	10	10	10

\ast

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

=

0	-30	-30	0
0	-30	-30	0
0	-30	-30	0
0	-30	-30	0

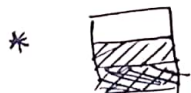
$\begin{matrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{matrix}$ → for horizontal edge.
 filter

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

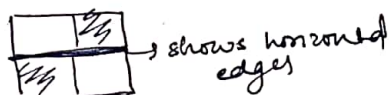
$\begin{matrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{matrix}$

=

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



=



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

Sobel filter.

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

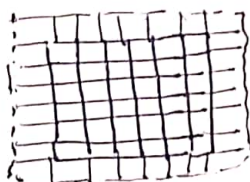
Scharr filter.

flip by 90° to get horizontal detection.
 - we can also learn these parameters.

* Padding:

$n \times n * f \times f = (n-f+1) \times (n-f+1)$
 pixels on the corner are used less as compared to pixels in mid.

soln: "pad" the img



$6 \times 6 \rightarrow 8 \times 8 * 3 \times 3 \rightarrow 6 \times 6$
 ∴ image conserved.

p → padding amount.

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

Valid & same convolutions

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

"Same" : size of i/p = size of o/p.

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

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

→ for same size of p.

f is almost always odd.

* Strided convolutions:

instead of stepping the filter by 1 step, we step 2 steps.

2	3	7	4	6	2	9
6	6	9	8	7	4	3
3	4	8	3	8	9	7
⋮						
0	1	3	9	2	1	7

*

3	4	4
1	0	2
-1	0	3

=

91	100	83
69	91	127
44	72	74

$n \times n * f \times f$ padding p stride ' s '

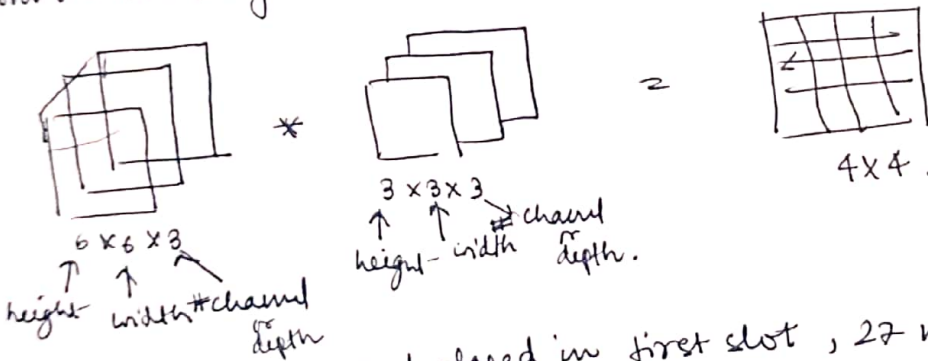
$$\text{of } p: \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor$$

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

if blue box doesn't fit, then don't proceed.

here, we've been using cross correlation (most literatures call this convolution)

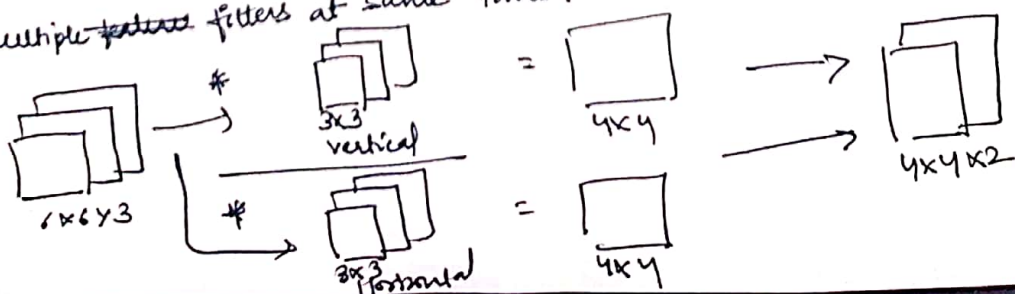
* convolutions over volume: conv on RGB images:



three $3 \times 3 \times 3$ filter taken and placed in first slot, 27 no^s multiplied and added, ...

$$\begin{matrix} R & \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} & G & \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} & B & \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} & \rightarrow \text{red edge detector.} \\ R & \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} & G & \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} & B & \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} & \rightarrow \text{any color edge detected.} \end{matrix}$$

• using multiple filters at same time.

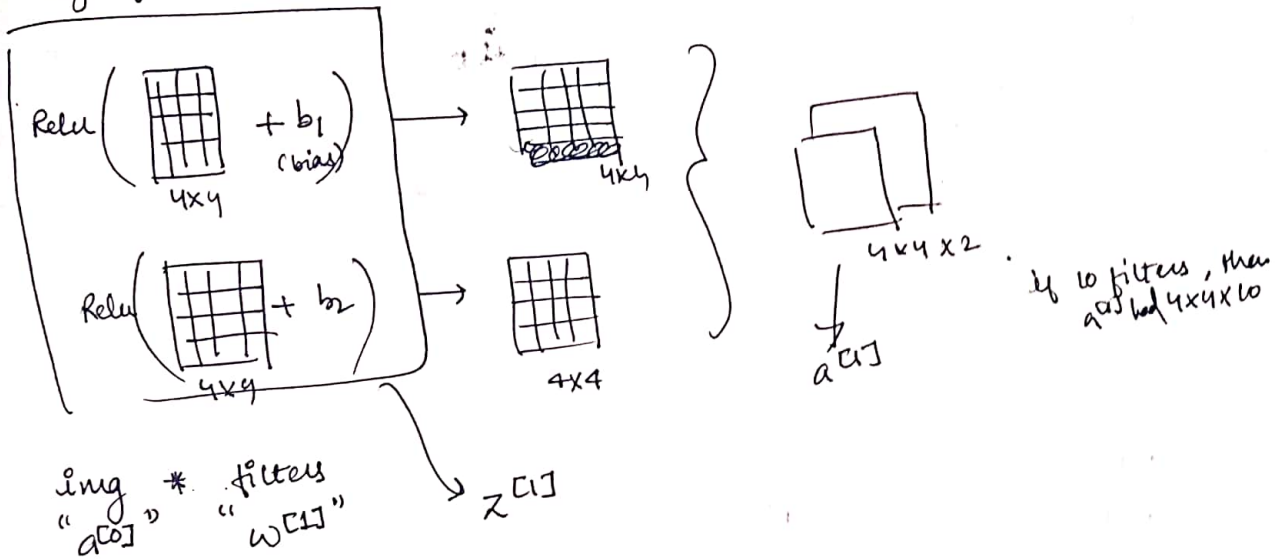


summary :

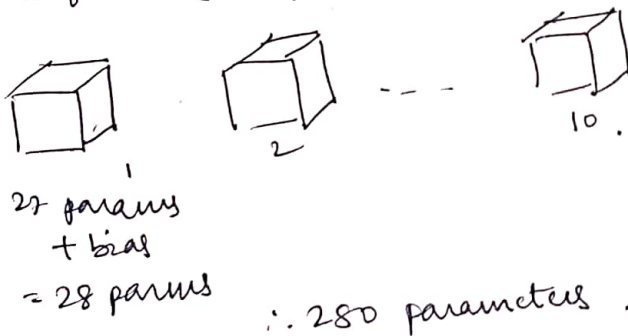
$$n \times n \times n_c * f \times f \times n_c \rightarrow (n-f+1) \times (n-f+1) \times n_c$$

\downarrow
filters.

* one layer of a CNN.



Ex: if 10 filters ($3 \times 3 \times 3$)



if layer l is a convolution layer:

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

$p^{[l]}$ = padding

$s^{[l]}$ = stride

Each filter is

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

activations $a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$

$$A^{[l]} = m \times n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$$

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

bias: $n_c^{[l]} - (1, 1, 1, n_c^{[l]})$ \hookrightarrow # filters in layer 'l'.

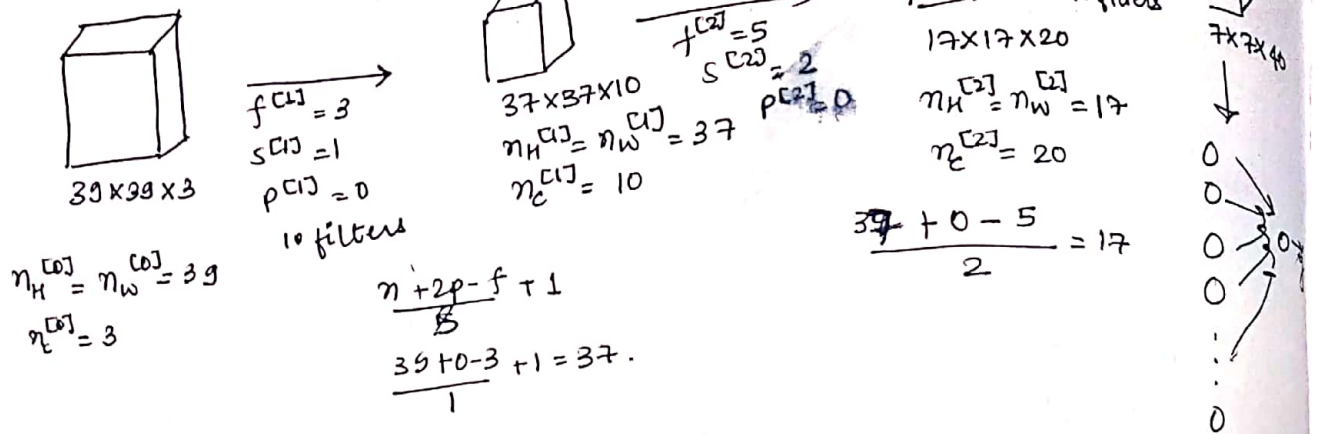
input: $n_H^{[l-1]} \times n_W^{[l-1]} \times n_c^{[l-1]}$

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

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

$H \rightarrow$ height
 $W \rightarrow$ width.

* Simple CNN Example :



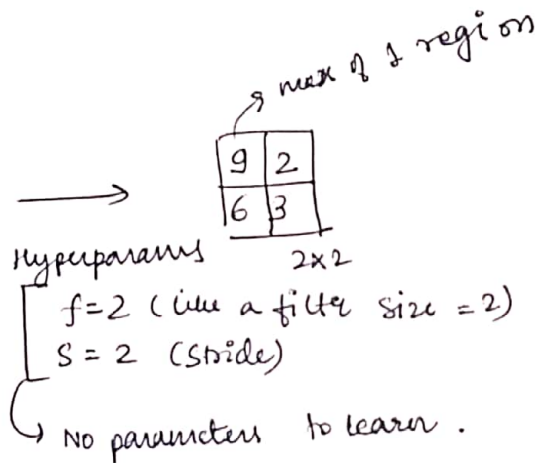
- Types of layers in CNN
- Convolution (CONV)
 - Pooling (POOL)
 - Fully connected (FC)

* Pooling layers :

- Max pooling

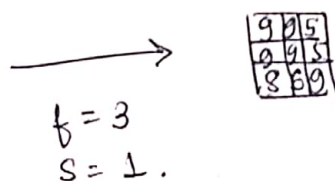
1	3	2	1
2	9	1	1
1	3	2	3
5	6	1	2

4x4



max pool: if feature detected, keep the high number.
 max value for each patch of feature map.

1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
8	3	5	1	0
5	6	1	2	9



if $5 \times 5 \times 2$ of p is $3 \times 3 \times 2$.
 do same computation on previous slice.

• Average pooling

average of ~~pooling~~ values in the patch of feature size.

Hyperparameters

f : filter size

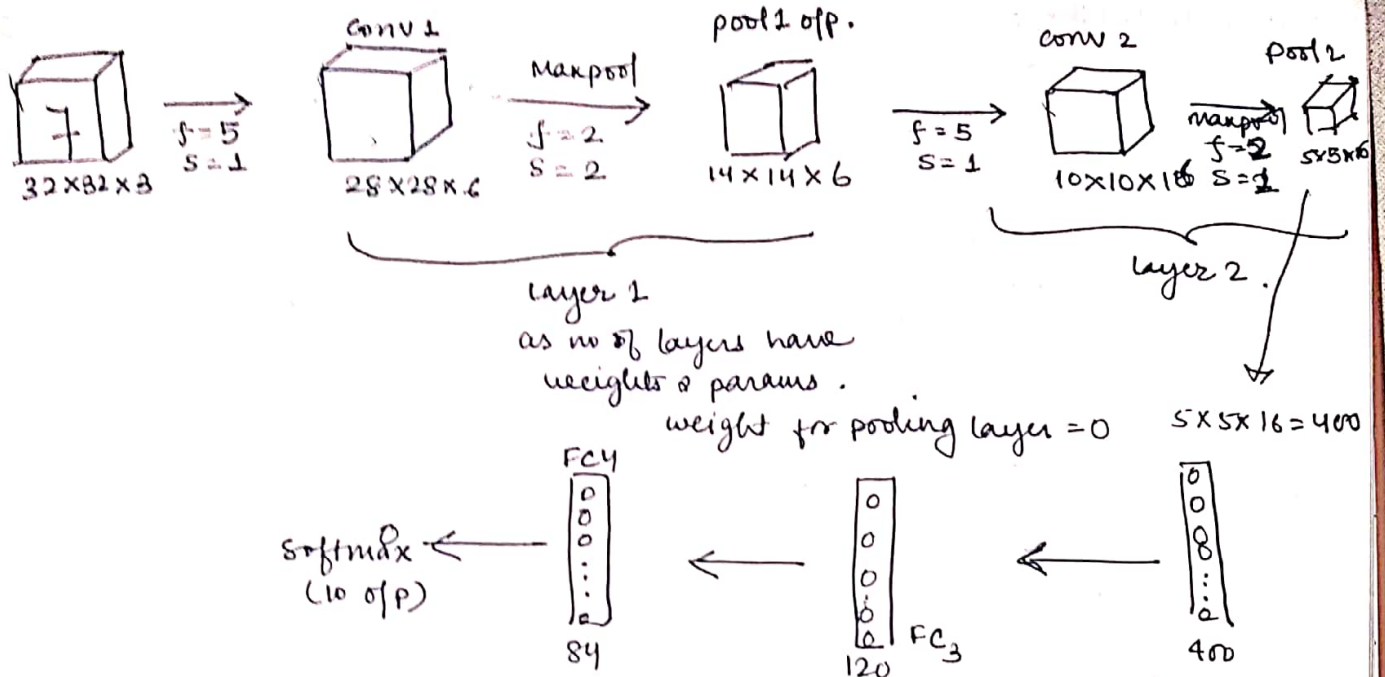
s : stride

max or average pooling.

p : padding.

no of i/p channels = no of o/p channels

* Example of complex CNN (LeNet-5)

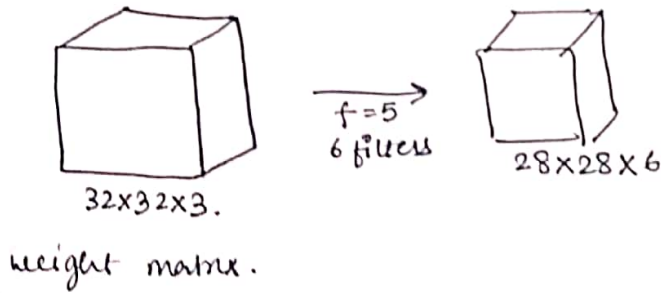


input:	Activ ⁿ shape	Size	# params
	(32 32 3)	3072	0
CONV1 (f=5, s=1)	(28, 28, 6)	6272	208

$$W^{[3]} = (120 \times 400)$$

$$b^{[3]} = (120 \times 1)$$

* why convolutions?



filter

$$5 \times 5 = 25 \text{ + bias}$$

$$= 26 \text{ params}$$

$$6 \times 26 = 156 \text{ params}$$

↓
bits.

$$3072 \times 4704 \approx 14M$$

Parameters sharing: A feature detector (ex vertical edge detector) that's useful in 1 part of img is prolly useful in other part.

Density of connection: in each layer depends only on ~~small~~ a small no of i/p.

Putting all together .



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

use gradient descent to optimise params to reduce 'J'.

⊕ Deep convolutional models : case study .

* why look at case studies :-

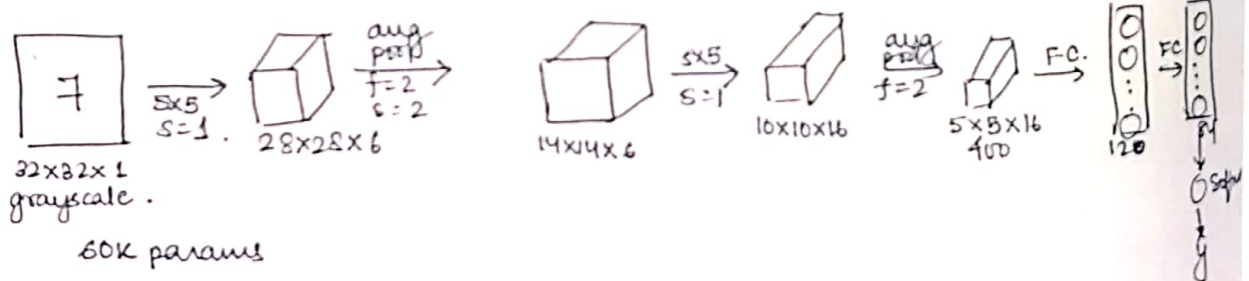
Classic network

- LeNet-5
- AlexNet
- VGG

ResNet (Residual Network)
- 152 layers deep .

* Classic Networks :

- LeNet-5



60K params

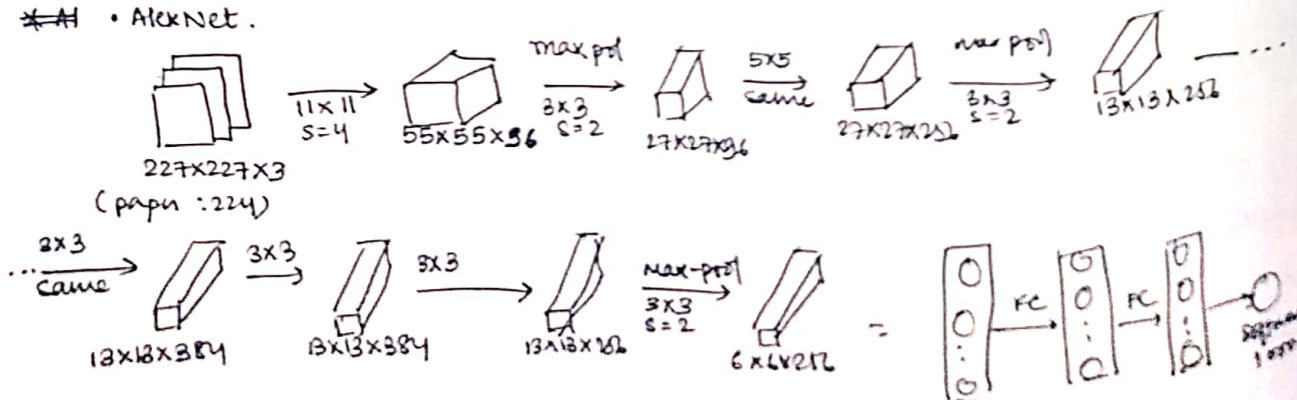
$m_H, m_W \downarrow \quad m_C \uparrow$

conv \rightarrow pool \rightarrow conv \rightarrow pool \rightarrow fc \rightarrow fc \rightarrow output.

Sigmoid/tanh on the research paper .

Le-Net5 had non linearity after pooling. (sections II & III)

* AlexNet .

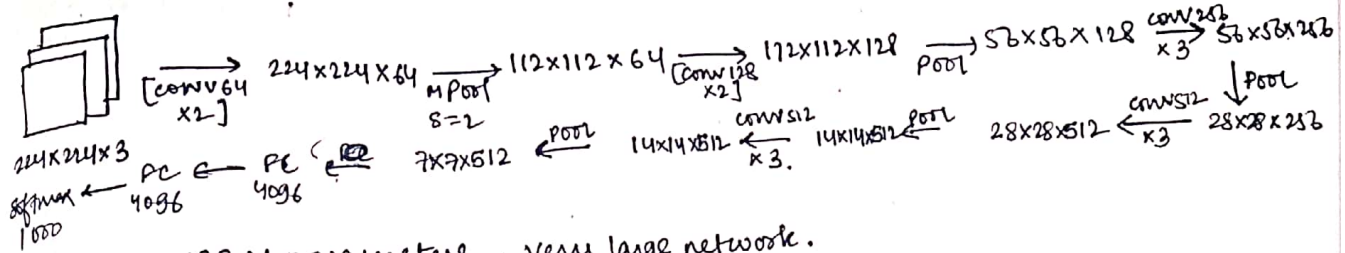


- much bigger
- 60M parameters
- ReLU activation function • when this was written, GPUs were slower.
- Has a local Response normalisation.

~~VGG16~~ • VGG-16 :

conv = 3x3 filter S=1, same

Maxpool = 2x2, S=2

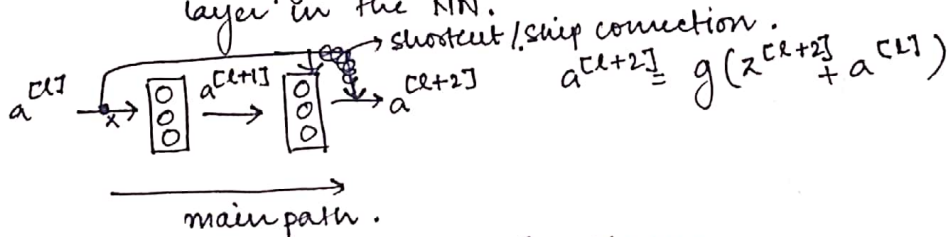


138 M parameters • very large network.

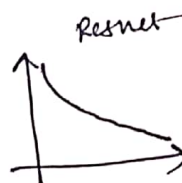
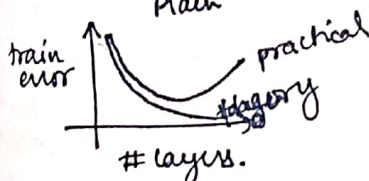
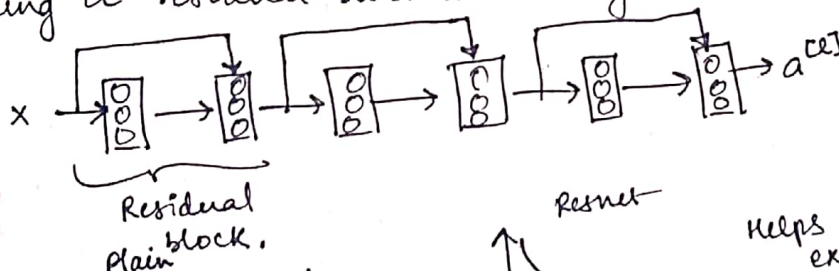
Notes

* ResNets (Residual N/ws)

skip connections : taking a layer & connecting it to a much deeper layer in the NN.

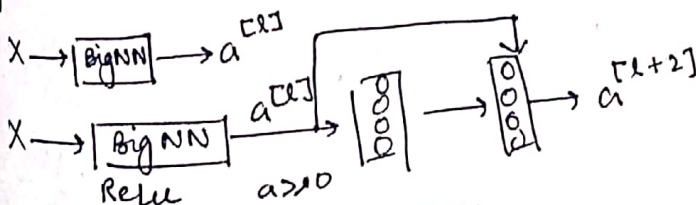


Taking a residual block & stacking them.



Helps with vanishing & exploding gradient problem.

Why ResNets work?



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

$$= g(w^{[L+2]} a^{[L+2]} + b^{[L+2]} + a^{[L]}) = g(a^{[L]}) = a^{[L]} \quad (\text{if } w^{[L+2]} = 0, b^{[L+2]} = 0)$$

identity funcⁿ is easy for residual block to learn.

$z^{[L+2]}, a^{[L]}$ have same dimensions

∴ same convolutions, so that dimensions are preserved.

256 & if p is 128 dim;

- add w_2 (to normalise dimension)

* Networks in Networks & 1×1 convolutions:

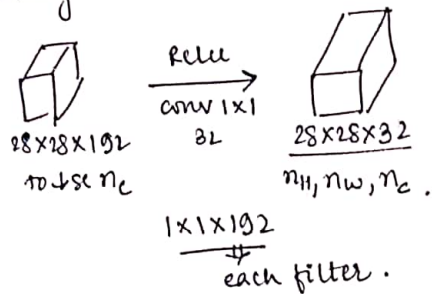
$$6 \times 6 * \boxed{2} = 6 \times 6. \quad (\text{each no doubled})$$

$$6 \times 6 \times 32 * \boxed{1 \times 1 \times 32} = \text{ReLU} \quad 6 \times 6 \times \# \text{filters}.$$

$$32 \rightarrow \# \text{filters} \quad n_c [c+1]$$

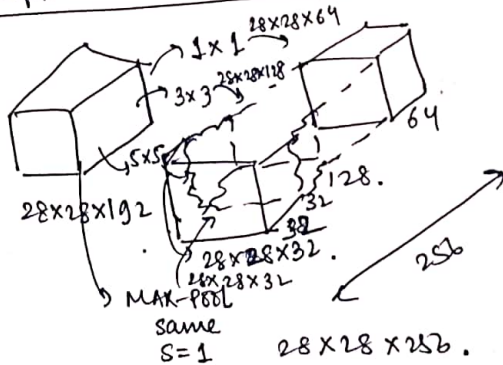
Network in network.

• using 1×1 convolutions.

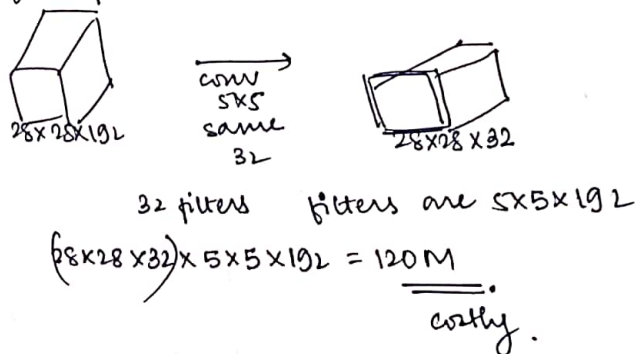


∴ no of channels can be used by 1×1 conv. ∴ reducing computation.

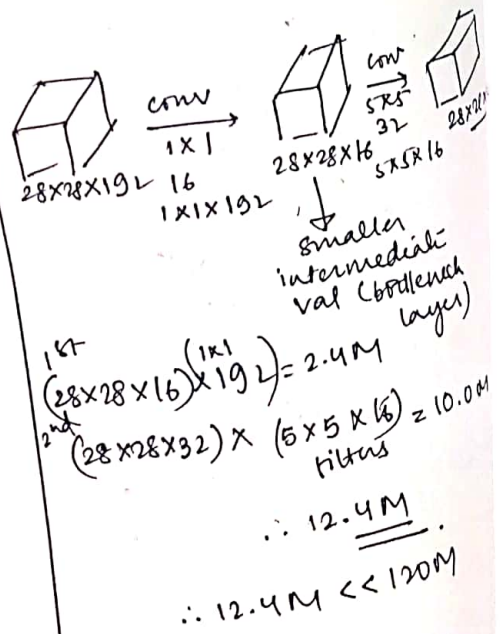
* Inception N/w motivation:



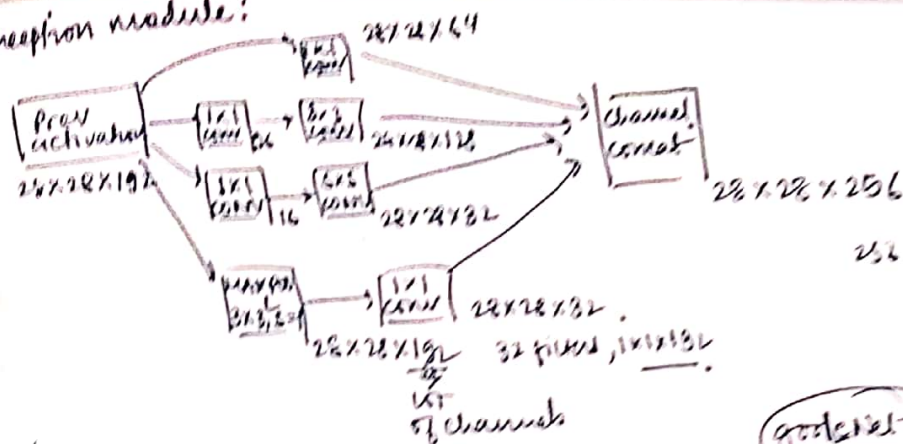
Problem of comp. cost:



alternative architecture (saves computation)



* Deception module:



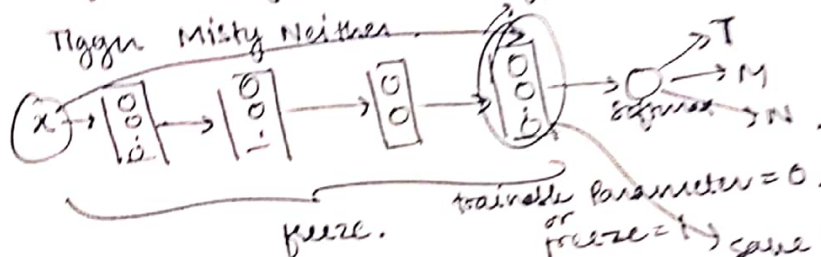
$$256 = 64 + 128 + 32 + 32$$

* Using open source implementation:

- github.

* Transfer learning :-

Taggy Misty Neither



→ if smaller data set, freeze more layers.

if larger data set, freeze ~~less~~ ^{when} layers, train later layers.

if ^{more} larger data set, train whole NN, as we take all weights to initialize rather than random initialization.

* Data Augmentation:

method:

- mostly used, {
- mirroring
 - random cropping (isn't perfect, but works well)
 - rotation
 - shearing
 - local warping

- color shifting: add different distortions to RGB channels.

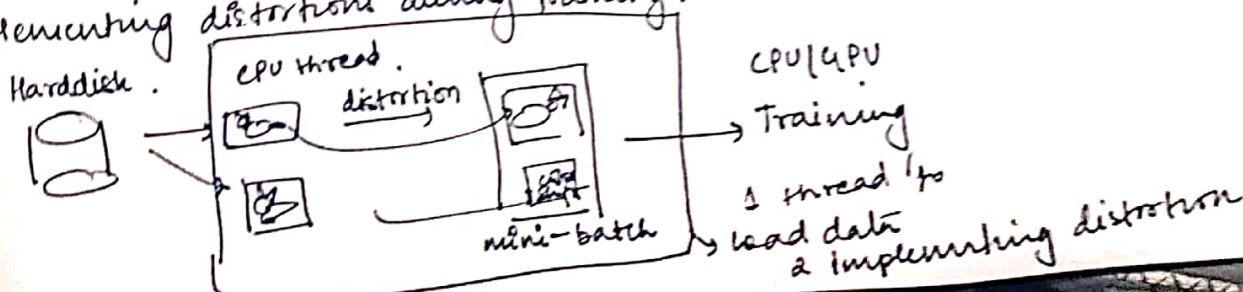
$$\begin{matrix} \text{Ex: } R & G & B \\ +10 & -20 & +20 \end{matrix}$$

* PCA (Principle Component Analysis) color Augmentations.

↓
AlexNet paper

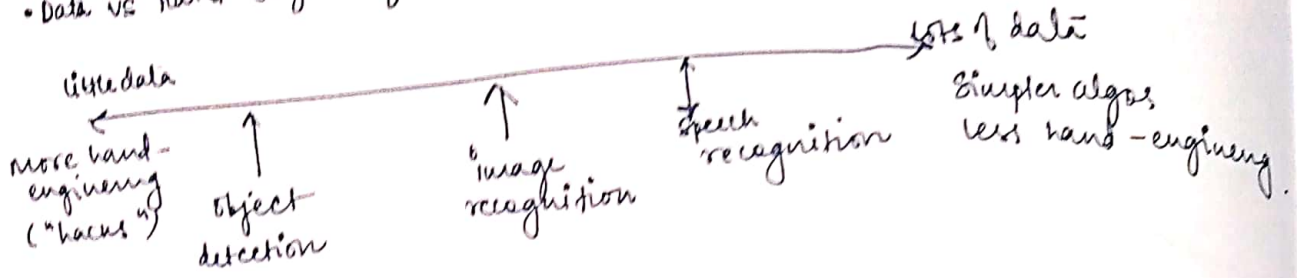
RB G
↓ ↓ ↓
subtracts 60% to RB subtracts less for G.

* Implementing distortions during training.



* State of Computer Vision:

- Data vs Hand engineering.



Transfer learning:

- Tips for doing well on benchmarks / winning competitions

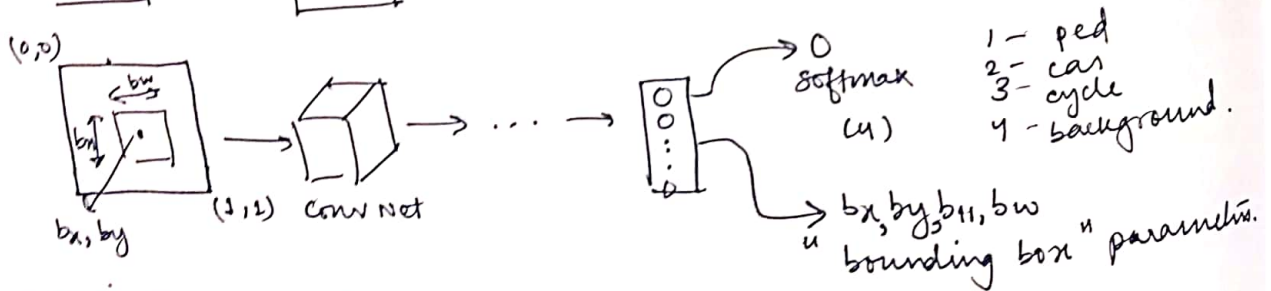
- Ensembling
 - Train several n/w & average their o/p.
- Multicrop at test time
 - Run classifier on multiple version of testing & average result.
 - Ex: 10-crop.
- Use open source code.
 - use architecture of n/w published in literature
 - use pretrained models & fine-tune on your dataset.

Object Detection:

* Object localisation:



Classification prob: 1 ~~prob~~ obj
Detection prob: multiple obj.



Defining Target label 'y'.

- 1 - ped
- 2 - car
- 3 - motorcycle
- 4 - background.

Need to output b_x, b_y, b_h, b_w , class label (1-4).

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \quad \text{is there an obj?} \quad (0 \text{ for background only})$$

if no obj,
 $p_c = 0$
rest are don't care

$$\text{loss} = \frac{1}{N} \sum_{i=1}^N \left[(y_i - \hat{y}_i)^2 + (y_i - \hat{y}_i)^2 + (y_i - \hat{y}_i)^2 \right] \quad \text{if } y_i = 0$$

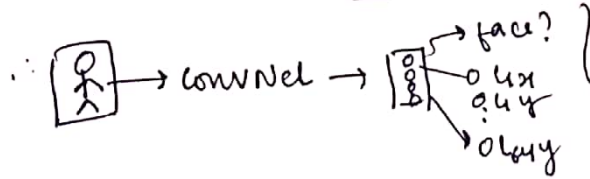
* landmark detection

l_1x, l_1y, l_2x, l_2y



maybe 64 landmarks on face.

l_1x, l_1y
 l_2x, l_2y
 \vdots
 $l_{64}x, l_{64}y$



128 units.
~~parameters~~
+ for face or not
128 landmarks (64x, 64y)

Ex: snapchat filters.

Ex: car detection.

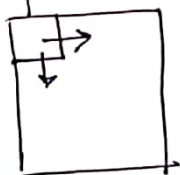
T.S.
Car 1
Car 2 1
No car 0
 \vdots



data set has cropped car images.

* sliding window detection.

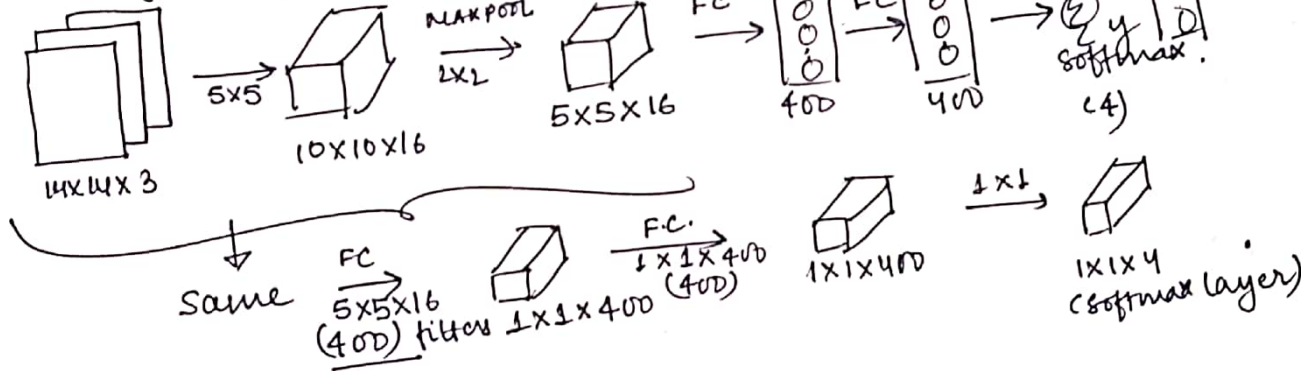
convNet $\rightarrow 0$



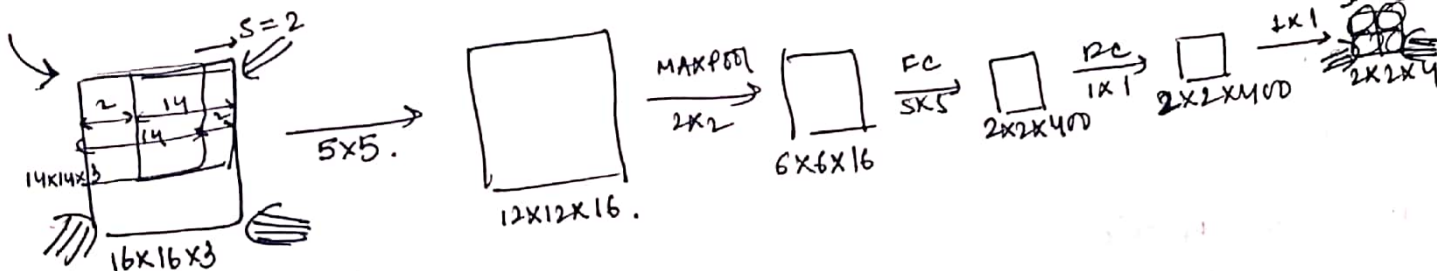
sliding window. (faster if large slides)

repeat using a larger window, may do 3rd time, even larger window there will be eventually a window which contains a car.

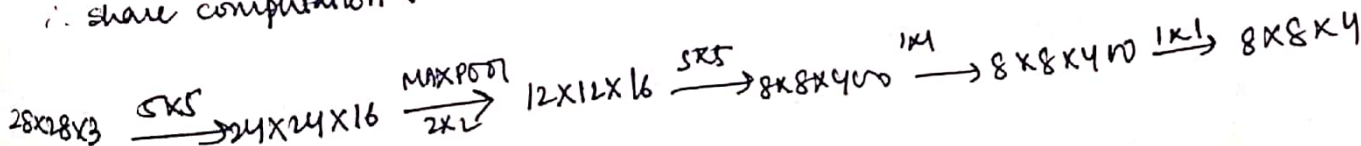
* convolutional implementation of sliding window:-



\rightarrow if test set has 16x16x3 img.

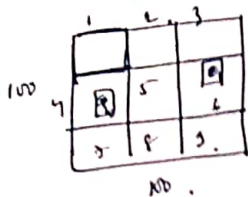


\therefore share computation.



* Bounding box predictions:
 - As convNet, sliding window pred. doesn't predict bounding boxes accurately.

• YOLO algorithm (You only look Once)



3x3 grids.

labels for training

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

YOLO research papers.

for 1, 2, 3, 5, 7, 8, 9.

$$y = \begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \\ 8 \\ 9 \end{bmatrix}$$

- width of object determines grid cell.

$$4: y = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

$$6: y = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Tgt of p: $\frac{3 \times 3 \times 8}{\text{grids}} \times 8 \text{ units in op.}$

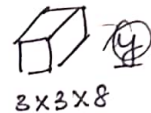


3x3x8.



100x100x3

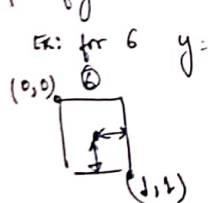
conv → conv layers → MAXPOOL →



3x3x8

↑ use 3x3 so that 1 grid cell has obj.

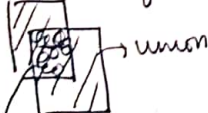
• specify bounding boxes:



$$y = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix} \begin{matrix} 0.4 \\ 0.3 \\ 0.5 \end{matrix} \left. \begin{matrix} \text{blue line} \therefore 0.9 \\ \text{50% of total grid} \end{matrix} \right\} \text{could be } > 1$$

* Intersection over Union:

(Iou) bounding boxes



intersection.

$$Iou = \frac{\text{size of intersection}}{\text{size of union}}$$

"correct" if $Iou \geq 0.5$ (threshold)

accuracy $\propto Iou$.

* Non-max suppression.

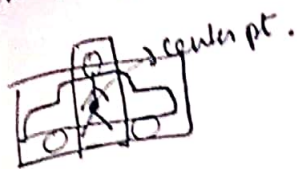
makes sure that an object is detection only once.
 - cleans up multiple detections

- rectangles with highest Iou part are highlighted, rest suppressed.
 - Discard all boxes with $p_c \leq 0.6$

- while any remaining boxes
 pick the one with highest p_c , discard the rest.

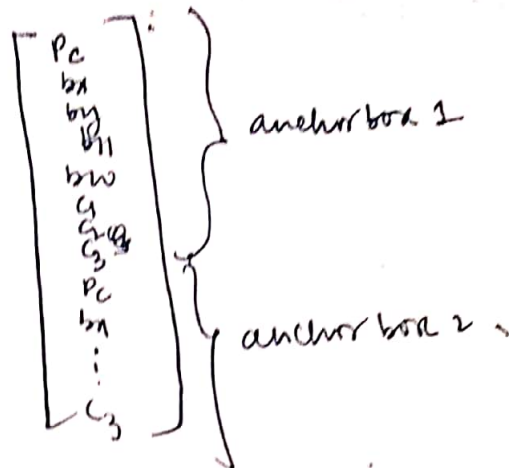
* Anchor Boxes:

overlapping objects:



Anchor box 1
Anchor box 2

$y =$



Now, of p 'y' becomes $3 \times 3 \times 16$

with 2 anchor boxes:

- Each obj in Tr. image is assigned to grid cell that contains object's midpt
- 2 anchor box for grid cell with highest IoU.

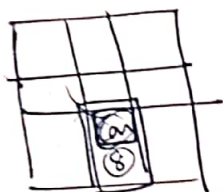
If only car:

- for anchor box 2: specify parameters
- for anchor box 1: $p_c = 0$, rest one don't care.

KYOLO Algorithm:

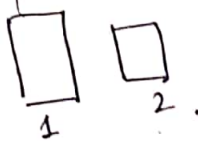
Training

- 1 - Ped
- 2 - car
- 3 - cycle.

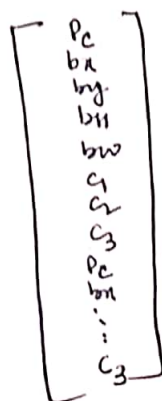


y is $3 \times 3 \times 2 \times 8$

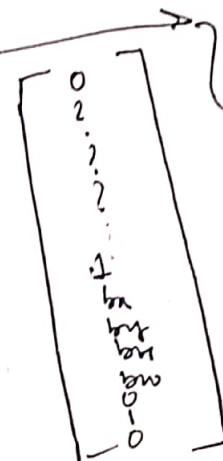
#anchors $\rightarrow 5 + \# \text{classes}$



$y =$



for '8'.



That's why
2 is empty, as
1 is begin

Outputting non max suppression of p :

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

* Region proposals: R-CNN.

- pick just a few regions with objects.
- runs a segmentation algo.
- Output: label + bounding box.

Fast R-CNN

- use convolution implementation of sliding windows to classify all proposed regions.

⊕ Facial Recognition :-

* Face recog:

* What is face recognition:

• face verification vs face recognition

Verification

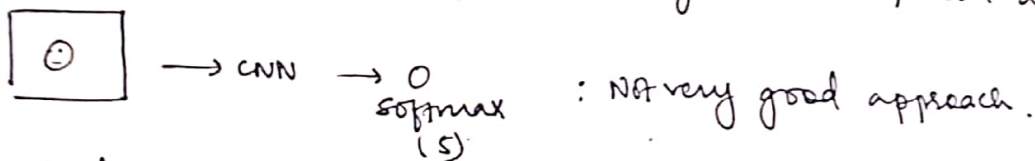
- I/p name ID
- output whether i/p img is that person or not.

Recognition

- Database of k persons
- get an i/p img
- Output ID if img is of any of the k persons (or not recognised)

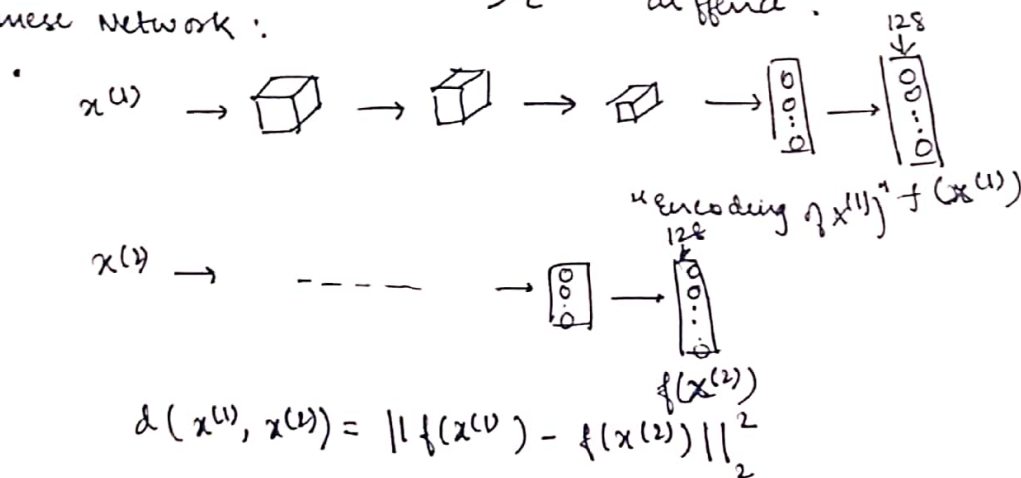
* One shot learning:

• learn from one example to recognise the person again



- good approach: learn similarity f^m
 $\rightarrow d(\text{img1}, \text{img2}) = \text{degree of difference b/w images}$
 if $d(\text{img1}, \text{img2}) \leq \tau$ same
 $> \tau$ difference.

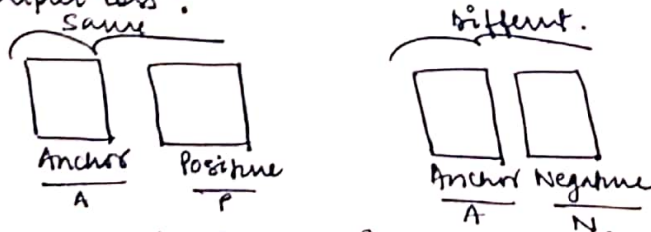
* Siamese Network:



• Deepface closing the gap to human level performance

if $x^{(1)}$ & $x^{(2)}$ are same person: $||f(x^{(1)}) - f(x^{(2)})||^2$ small
 else large.

* Triplet loss:



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

want: $||f(A) - f(P)||^2 \leq ||f(A) - f(N)||^2$
 $d(A, P) \leq d(A, N)$

hyperparameter that prevents NN from outputting trivial solⁿ.

Given 3 images A, P, N

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

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

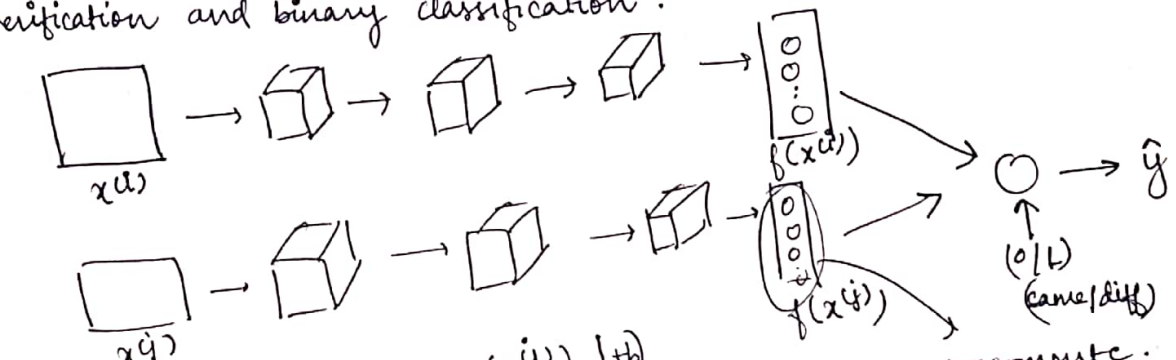
Paper: FaceNet

Training set: 10K pics of 1K persons

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

• those triplets that are "hard" to train on
 $d(A, P) + \alpha \leq d(A, N)$
 $d(A, P) \approx d(A, N)$

* Face verification and binary classification:



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

or

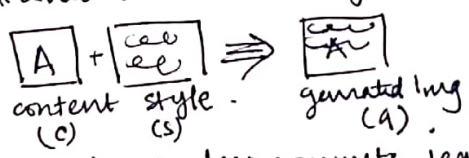
$$\chi^2 = \frac{(f(x(i))_k - f(x(j))_k)^2}{f(x(i))_k + f(x(j))_k}$$

precompute.
 upper convnet computes encoding & precomputed is compared with.

χ^2 (chi square formula)

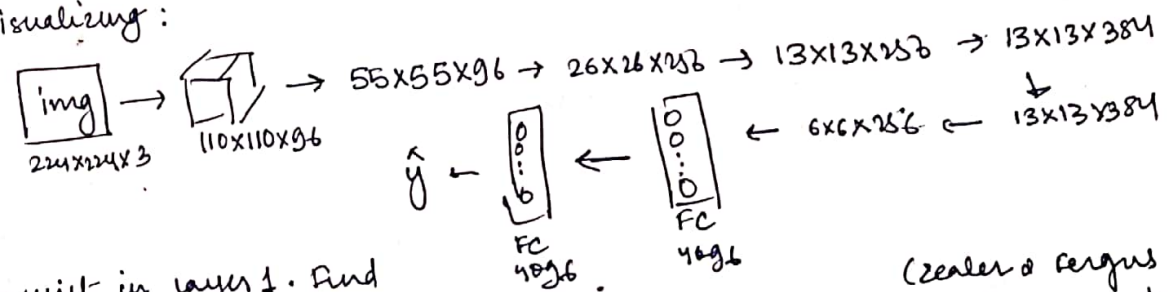
* Neural Style Transfer.

* What is Neural Style Transfer?



* What are deep convnets learning?

• visualizing:



- Pick a unit in layer 1. Find the nine img patches that maximise unit's activation

- Repeat.

(Zeiler &ergus)
 Visualizing & understanding convolutional networks.

* Cost function:

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

How similar is content How similar is cost

(A neural algo of artistic style gets el at it)

α, β - hyperparameters.

• finding generated 'image' G .

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

2. Use G.D. to minimise $J(G)$.

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

* Content cost function: ($J(C, G)$)

• say you use 'l' to compute content cost.

• use pre-trained convNet (Eg: VGG N/w)

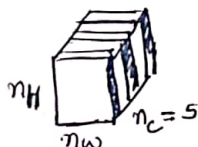
• let $a^{[l]}(C)$ and $a^{[l]}(G)$ be activation of layer l on imgs.

• if $a^{[l]}(C)$ & $a^{[l]}(G)$ are similar, both have similar content.

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

* Style cost funcⁿ: ($J(S, G)$)

• 'l' to measure "style".



How correlated are activations across channels?

let channel 1 corresponds to sth, channel 2 corresponds to sth, then check how much they are correlated.

• style matrix

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

$$G_{kk'}^{[l]}(S) = \sum_{i=1}^{n_H^{[l]}} \sum_{j=1}^{n_W^{[l]}} a_{ijk}^{[l]}(S) a_{ijk'}^{[l]}(S)$$

n_c
 $G^{[l]}$
 k, k'
 $k=1 \dots n_c^{[l]}$
 (no of channels)

$$G_{kk'}^{[l]}(G) = \sum_{i=1}^{n_H^{[l]}} \sum_{j=1}^{n_W^{[l]}} a_{ijk}^{[l]}(G) a_{ijk'}^{[l]}(G)$$

do this for every val of k' & k to complete style matrix
 "gram matrix"

$$J_{\text{style}}^{[l]}(S, G) = \| G^{[l]}(S) - G^{[l]}(G) \|_F^2$$

$$\sum_{i=1}^{n_H^{[l]}} \sum_{j=1}^{n_W^{[l]}} \sum_{k=1}^{n_c^{[l]}} \sum_{k'=1}^{n_c^{[l]}} (G_{kk'}^{[l]}(S) - G_{kk'}^{[l]}(G))^2$$

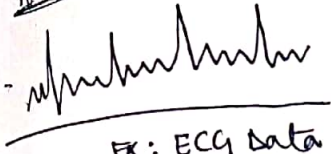
$$J_{\text{style}}^{[l]}(S, G) = \frac{1}{(2 n_H^{[l]} n_W^{[l]} n_c^{[l]})^2} \sum_k \sum_{k'} (G_{kk'}^{[l]}(S) - G_{kk'}^{[l]}(G))^2$$

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

↑
hyperparameter

* 1D and 3D generalisations:
 • conv in 2D & 1D.

11	20	15	3	18	2	4
----	----	----	---	----	---	---



ex: ECG data

14

*



1	3	10	3	1
---	---	----	---	---

5

= 10 dimensional

3D data:

ex: CT scan.



3D vol

*



3D filter

$14 \times 14 \times 14 \times 1$
 $H \quad W \quad D$

#chans

$5 \times 5 \times 5 \times 1$

#chans.

16 files $\rightarrow 10 \times 10 \times 10 \times 16$
 $\rightarrow 10 \times 10 \times 16$ volume