

CNN

Vertical Edge detection

• Convolutional Operation

convolution operation

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

6x6 grayscale image

\times

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

3x3 kernel/ filter

$=$

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

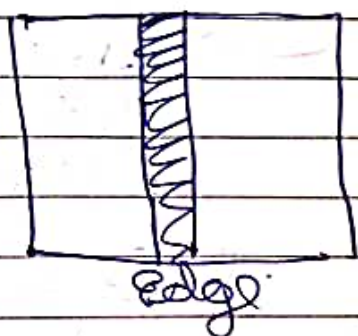
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0

\times

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



Vertical Edge is a $n \times n$ filter which has bright pixels on left and dark on right ^{positive edge}

If transition is light to dark, edge is light
and if " is dark to light, edge is dark
↳ neg edge

Filter for Horizontal Edge

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

• Sobel filter

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

puts more emphasis on centre pixel hence more robust

• | | | | |----|---|-----| | 3 | 0 | -3 | | 10 | 0 | -10 | | -3 | 0 | -3 | → Scharr filter

We can initialize the nine values of 3×3 filter as random values and learn their values by back propagation.

→ Image $\rightarrow n \times n$

Filter / kernel $\rightarrow f \times f$

Output $\rightarrow n - f + 1 \times n - f + 1$

To reduce shrinking of image and loss of data in convolutional operation, use padding.

If p is padding amount then:

Output: $n + 2p - f + 1 \times n + 2p - f + 1$

Generally by convention, padding is of zero and output size same as of image

- Valid Convolution: $h=0$
- Same Convolution: ~~$h=k$~~ $(k \neq 0 \rightarrow h=k)$

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

$$2h - f + 1 = 0$$

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

f is generally odd sized filter i.e. 3×3 or 5×5 etc. so h is an integer and symmetric padding.

If f is even \rightarrow asymmetric padding

In odd sized filter, central pixel can be used to denote location of filter

Input image $\rightarrow n \times n$

kernel $\rightarrow f \times f$

padding $\rightarrow h$

stride $\rightarrow s$

$$\therefore \text{Output} \rightarrow \left[\frac{n + 2h - f + 1}{s} \right] \times \left[\frac{n + 2h - f + 1}{s} \right]$$

If this is not integer, round it down

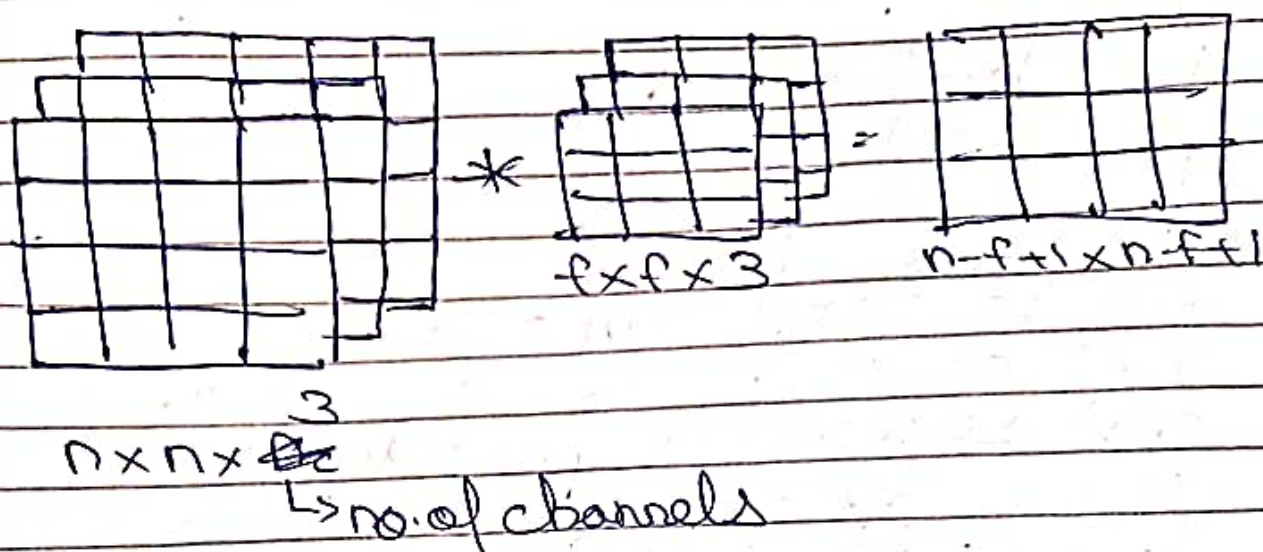
$$\frac{n + 2h - f + 1}{s} = n$$

$$n + 2h - f + s = ns$$

$$2h - f + s = n(s-1) \quad 2h - f + n = s(n-1)$$

$$\therefore s = \frac{2h - f + n}{n-1}$$

In RGB Image



Input $\rightarrow n \times n \times n_c$
 Kernel $\rightarrow f \times f \times n_c$ — same
 Output $\rightarrow n-f+1 \times n-f+1 \times n_f$

no. of filters

If padding = p stride = s then

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

Convolutional Neural net of l layers

$f^{[l]} \rightarrow$ filter size of layer l

$s^{[l]} \rightarrow$ stride of layer l

$p^{[l]} \rightarrow$ padding of layer l

$n_c^{[l]} \rightarrow$ no. of filters in layer l

Input $\rightarrow n_h^{[l-1]} \times n_w^{[l-1]} \times n_c^{[l-1]}$

(Input of a layer is output of prev)

Output $\rightarrow n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

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

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

Size of each filter = $f^{[l]} \times f^{[l]} \times n_c^{[l-1]}$

After Act Activation,
 $a^{[l]} \rightarrow n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

In Vectorized form,

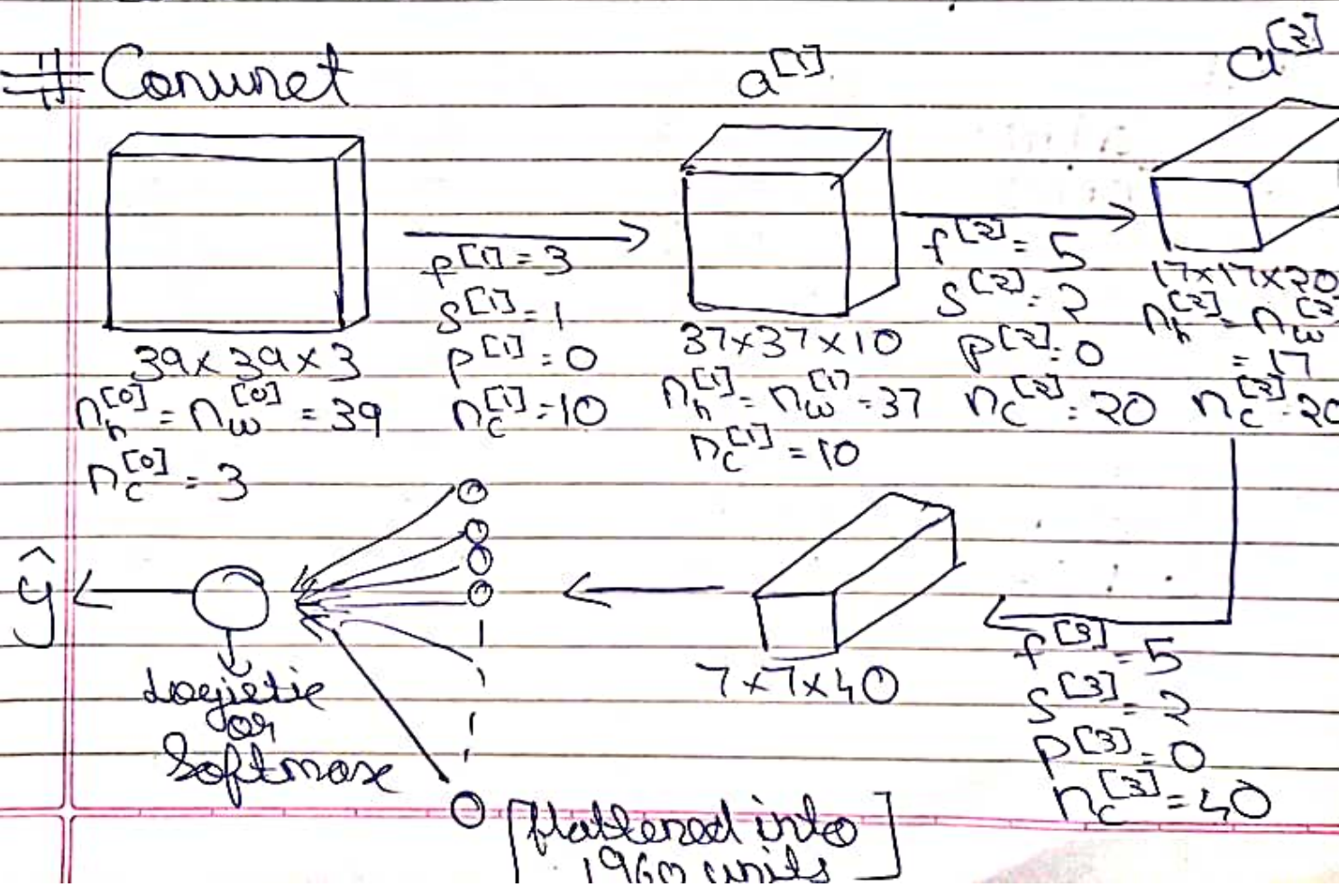
$$A^{[l]} \rightarrow m \times n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$$

↳ no. of training Examples

$$\text{Weights } w^{[l]} = f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$$

$$\text{Bias } \rightarrow b^{[l]} \rightarrow n_c^{[l]}$$

Convnet



Max and Avg Pooling

Input $\rightarrow N_h \times N_w \times N_c$

Filter size $\rightarrow f \times f \times N_c$

~~Output $\rightarrow N_h - f + 1$~~

Stride $\rightarrow S$

Padding $\rightarrow p$

Output $\rightarrow \left\lfloor \frac{N_h + 2p - f + 1}{S} \right\rfloor \times \left\lfloor \frac{N_w + 2p - f + 1}{S} \right\rfloor \times N_c$

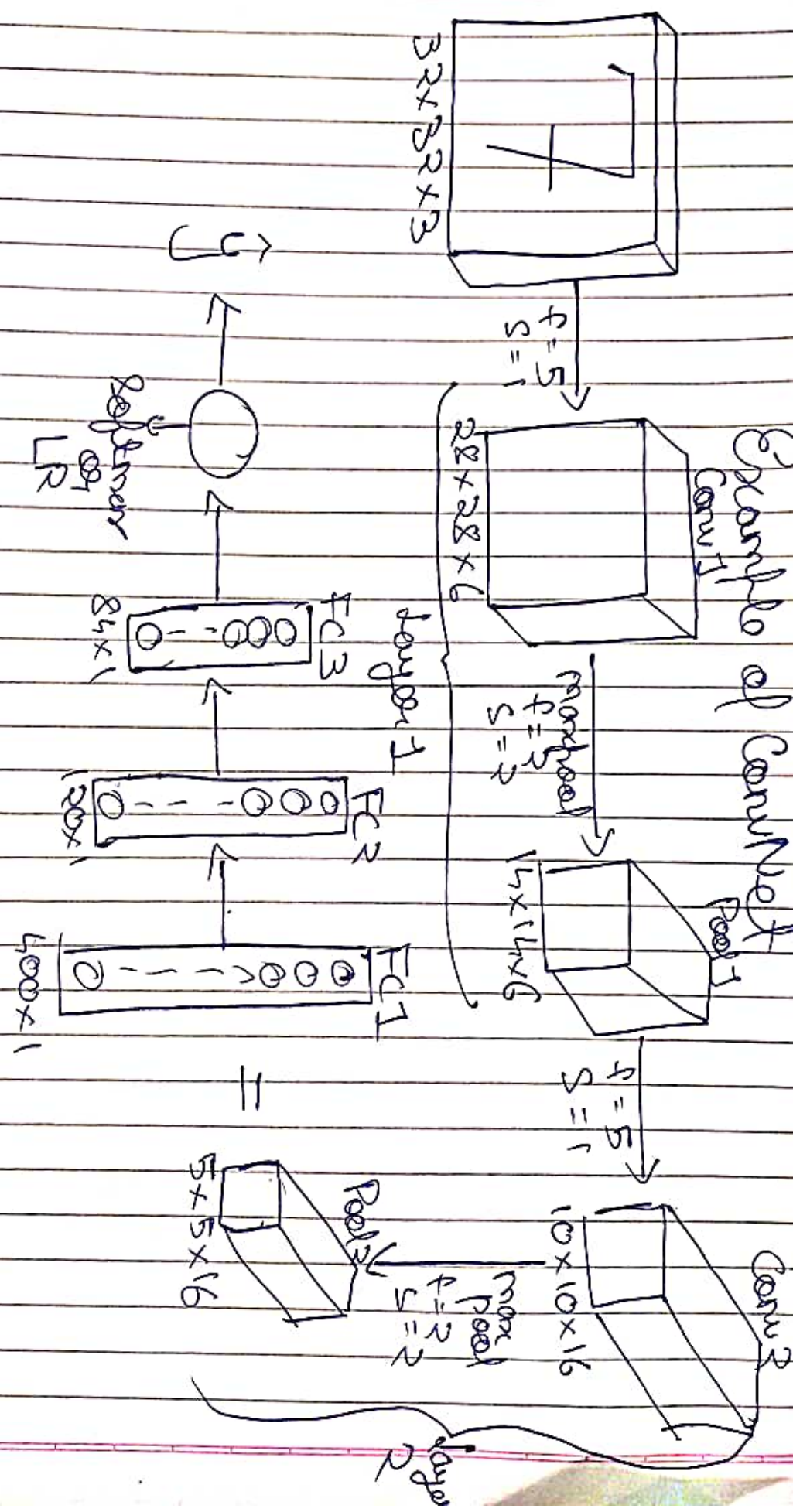
In Each channel, Pooling happens individually

No. of Input channel = No. of output channel

No Parameters to learn, once set, they are constant

Example of Convnet

#		$a^{[L]}$ shape	$a^{[L]}$ size	Params
	Input:	32, 32, 3	3072	0
$f=5, s=1$	Conv 1	28, 28, 16	4704	456
$f=2, s=2$	Pool 1	14, 14, 16	1176	0
$f=5, s=1$	Conv 2	10, 10, 16	1600	2416
$f=2, s=2$	Pool 2	5, 5, 16	400	0
	FC3	120, 1	120	48120
	FC4	84, 1	84	10164
	Softmax	10, 1	10	850
	FC1	400, 1	400	400



Object Detection

Input \rightarrow Random image

Output \rightarrow
 1] Class 1 (Pedestrian)
 2] Class 2 (Car)
 3] Class 3 (Motorcycle)
 4] Background (NOTA)

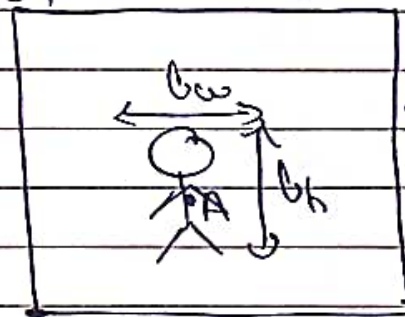
Output Vector \rightarrow

Logistic Regression Loss
 Squared error or something similar

Log loss error

P_c \rightarrow if object is there, otherwise
 b_x b_y } midpoint of object
 b_h b_w } h, w of object
 C_1 C_2 C_3 } which class object belongs

Eg: (0,0)



A

$A = (b_x, b_y)$
 $(0.5, 0.5)$

(1,1)

1
0.5
0.5
0.3
0.6
1
0
0

If $P_c = 0$, other all elements are don't care cases

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

$y_1 = P_c$

Face Recognition

- Verification 1%1
Input img, name / ID
Output whether input image is desired person

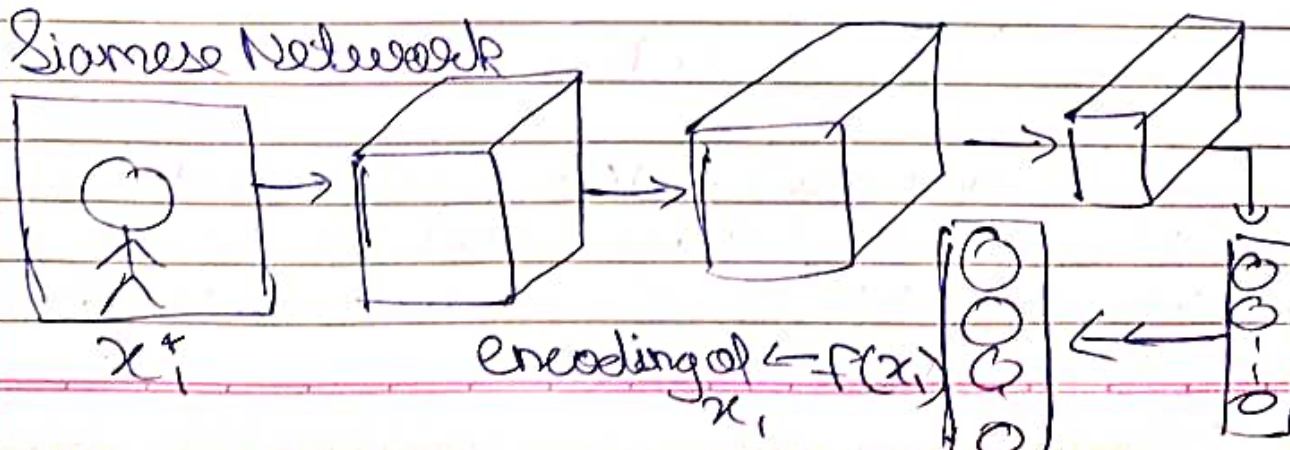
- Recognition 1%K
Has a database of K person
Given input image
Output ID if the image is any of the K persons

≠ One shot learning
Learning from one example to recognize the person again

- Similarity func
 $d(\text{img1}, \text{img2}) \rightarrow$ degree of difference between images
for Verify { $d(\text{img1}, \text{img2}) \leq \tau$ (Same person)
 $> \tau$ (Diff person)
 $\tau \rightarrow$ Higher parameter

For face recognition
Compute d for input image with all img in database

≠ Siamese Network



Calculate $f(x_1)$ and $f(x_2)$ for two input images x_1 and x_2 then
 $d(x_1, x_2) = ||f(x_1) - f(x_2)||^2$
 Train NN such that $||f(x_1) - f(x_2)||^2 \leq \alpha$
 if images are same else greater
 NN has same parameters for all images

≠ Triplet Loss func

Img 1	Img 2	Img 1	Img 3
Anchor	Positive	Anchor	Negative
A	P	A	N

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

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

$\alpha \rightarrow$ margin
 α is used so that NN doesn't output zero or same value encoding for each image

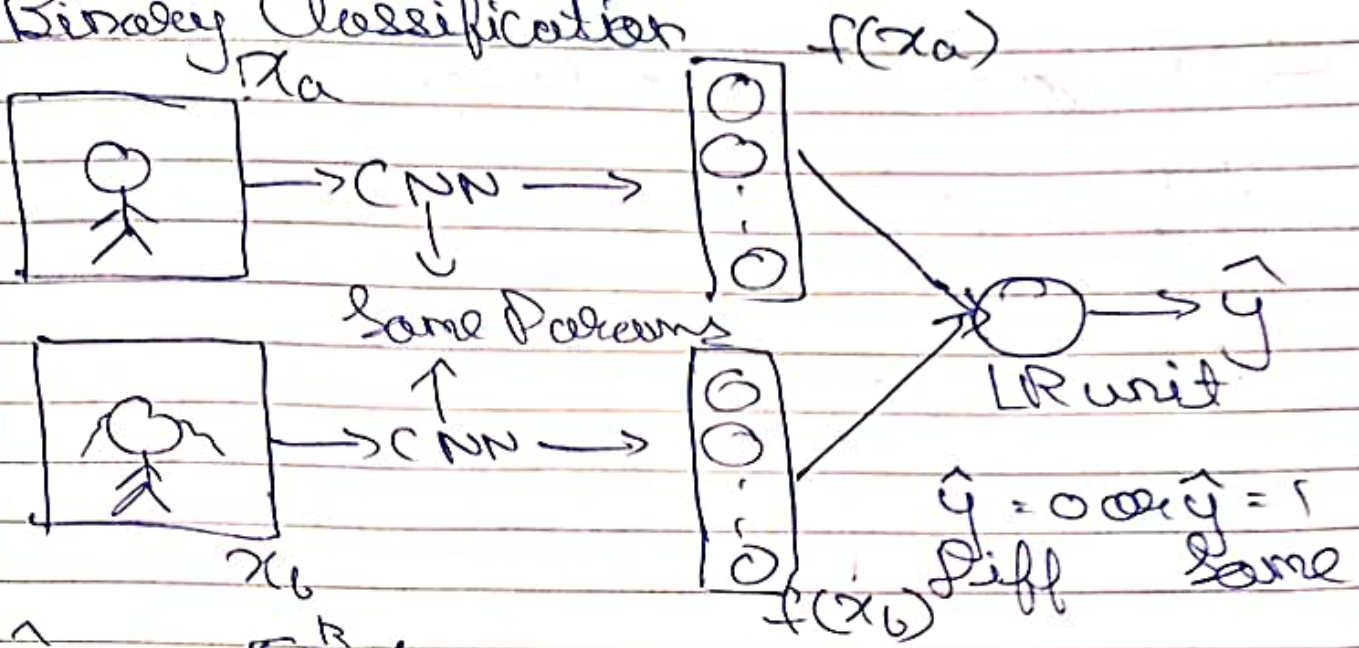
α Pushes (A, P) pair and (A, N) pair further away from each other

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

$$J = \sum_{i=1}^R L(A_i, P_i, N_i)$$

If training set is 10k pics of 1k diff persons - take 10k pics and generate triplets and train NN for gradient descent on Cost func

Binary Classification



~~$$\hat{y} = 0 \rightarrow \sum_{i=1}^R (w_i | f(x_a)_i - f(x_b)_i |) = f$$~~

$$\hat{y} = \sigma \left[\sum_{i=1}^R (w_i | f(x_a)_i - f(x_b)_i | + b) \right]$$

Then do gradient descent

Neural Style Transfer

Content (C) + Style (S)

↓
Generated image (G)

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

First randomly generate G then

$$G = G - \frac{dJ(G)}{dG}$$

Object Detection

Performance Evaluation Metrics

- 1] How good is location (IOU)
- 2] How good is classification (mAP)

• IOU \rightarrow how close is predicted bounding box to ground truth

$$\text{IOU} \rightarrow \frac{\text{Area}(A \cap B)}{\text{Area}(A \cup B)} \quad [0, 1]$$

• Average Precision

1] Confusion Matrix

Predicted	True +ve	False +ve
	True neg	False neg
True		

$$2] \text{ Precision} = \frac{TP}{TP + FP}$$

$$3] \text{ Recall} = \frac{TP}{TP + FN}$$

Average Precision \rightarrow Area under Precision Recall curve

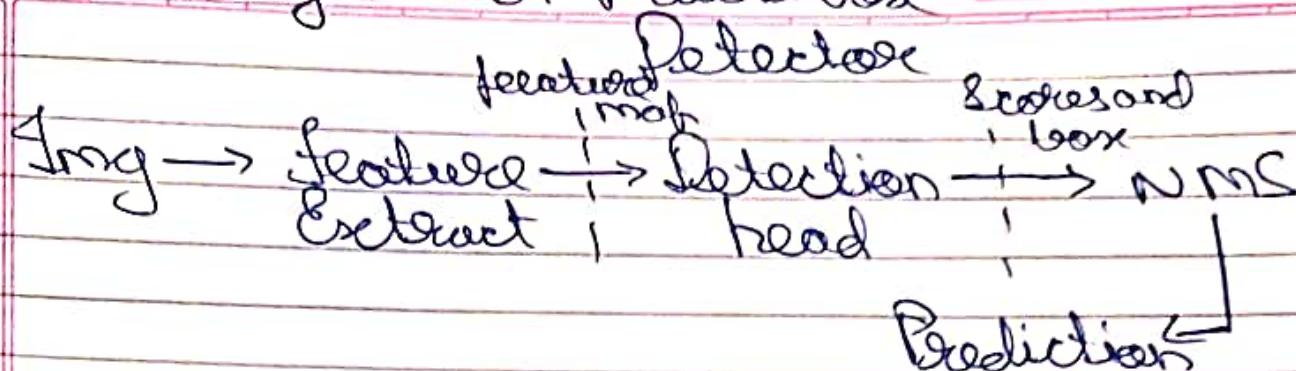


Average Precision \rightarrow Single Class
Mean Avg Precision \rightarrow Multiple Class
mAP

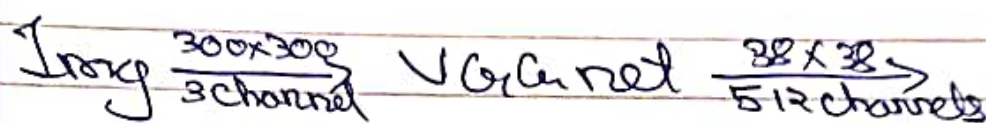
Single Shot Multibox

Page No.

Date

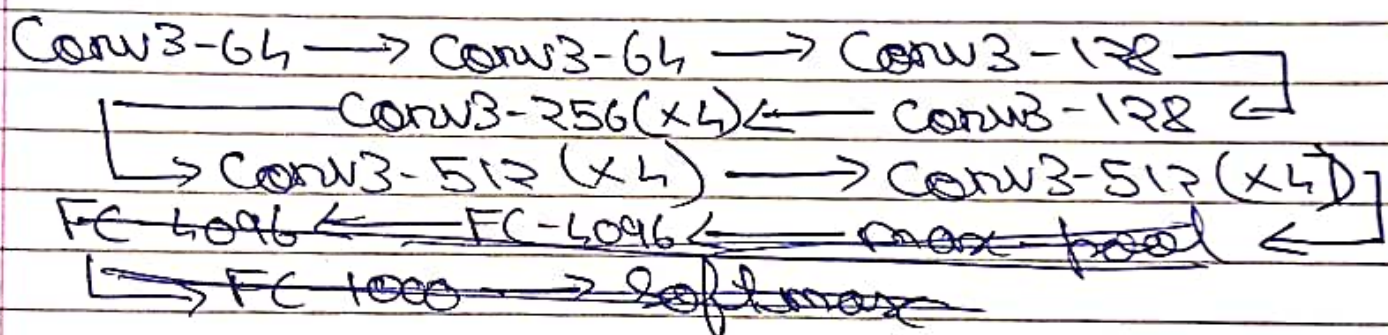


For Feature Extraction



feature maps

VGG net:



Feature maps

