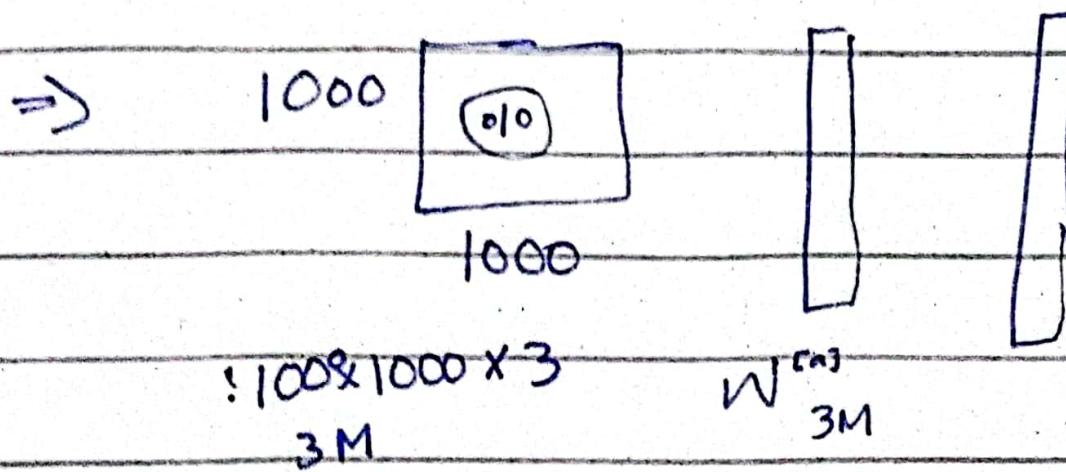


24/4/25  
(After mids)

## CONVOLUTIONAL NEURAL NETWORK



$$w^{[a]} = w^{[1]}, w^{[2]} \quad (1000, 3 \text{ million})$$

$\downarrow$   
3 billion

⇒ Overfitting  $\rightarrow$  Computational  $\Rightarrow$  Flattening  
(less)

⇒ By Default, convolutional neural network nature is to shrink.

$n-f+1$

: Now Everything  
is practical

$\Rightarrow$  Convolutional first layer finds edge , middle layer finds objects small parts ( cat eye, ear etc....) and last layer determine whole object.

$\Rightarrow$  If filter and input are similar then answer will be big otherwise small.

$\Rightarrow$  Sparse Connection and parameter sharing is property of convolutional neural network.

$\Rightarrow n=6, f=3, s=1, p=1$

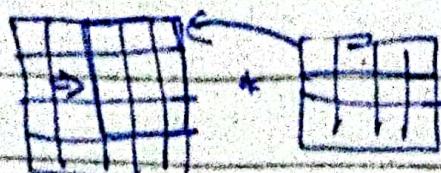
$$n + 2p - f + 1$$

$$: 6 + 2 - 3 + 1$$

$$: 6 \quad (6 \text{ by } 6)$$

$\Rightarrow$  Strided convolution:

$\rightarrow$  Skips the column to reduce image size and get smaller area .



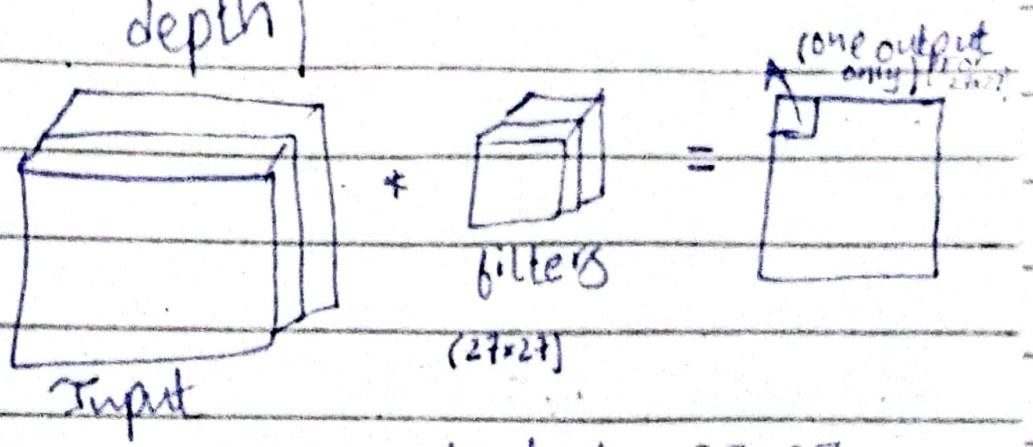
:  $L = 2p - f + 1, n_{\text{op}}$   
S S

$$H \cdot W \Rightarrow n=6, f=3, s=2, p = (6+3-2)/2 = 3.5$$

(Identify Padding for that stride)

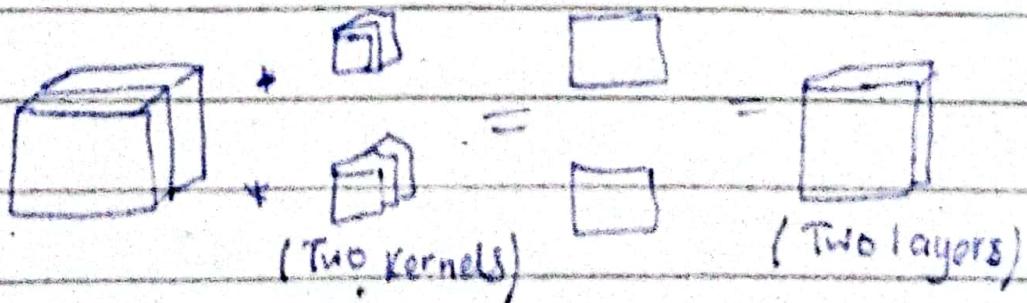
### ⇒ Convolutions over volume:

- For RGB channels (Multiple channels)
- Filter depth should be match channel depth (input depth)



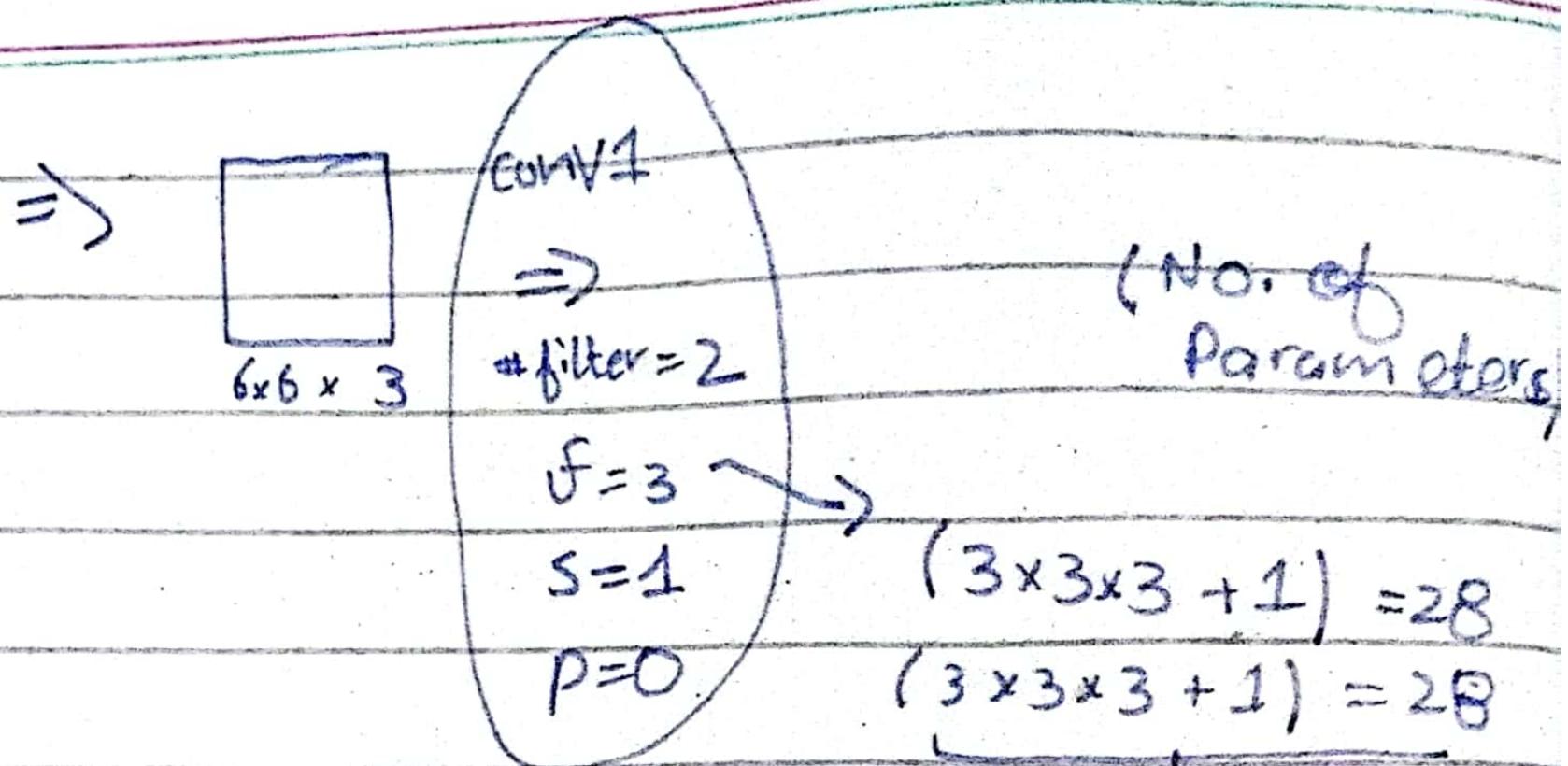
: (Only one output for  $27 \times 27$  input & filter)

- If more filters are applied more output will be there:



$$A = \text{ReLU}((W^T \cdot x) + b)$$

12 biases for 2 filters



: If 10 filters, then = 56

No. of parameters (280)

6|5|25

\*



$f=5$

$s=1$

$p=0$

$32 \times 32 \times 3$

# of filters = 6

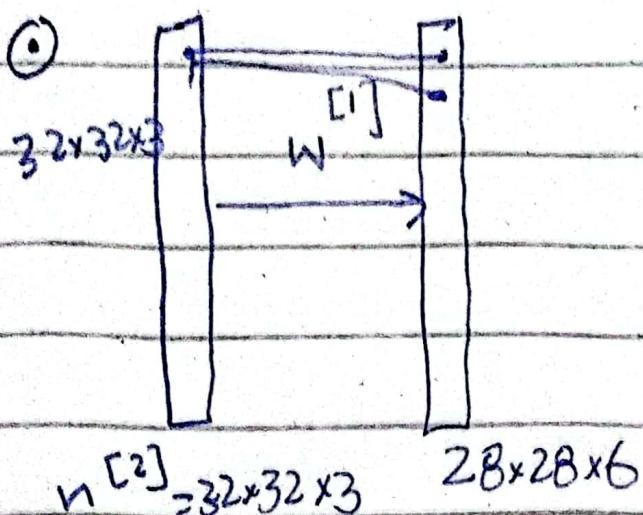
$$\frac{n+2p-f}{s} + 1$$

$$\rightarrow \frac{32+2(0)-5}{1} + 1$$

$$\frac{32-5}{27+1} = 28 \times 28 \times 6$$

Activation size

### ① For CNN



: no. of parameters

$$= 6 \times (5 \times 5 \times 3 + 1)$$

$$= 456$$

(Less Parameters)

(Because of

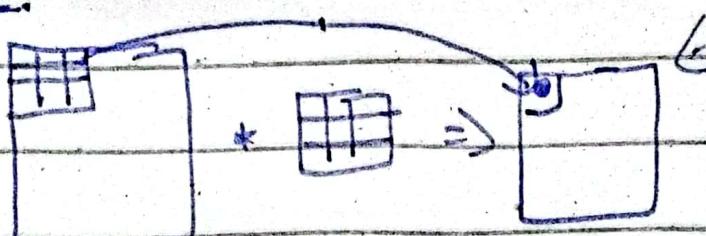
sparse connection & sharing weights)

(No overfitting)

$$w^{[1]} = (28 \times 28 \times 6, 32 \times 32 \times 3)$$

= 14 million (More Parameters)

### ② In CNN



(Less Parameters)

$\Rightarrow$  In CNN we see:
 

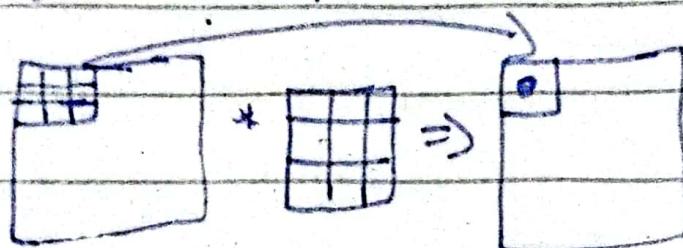
- (i) Convolutional layer
- (ii) Pooling Layer
- (iii) Fully connected layer.

③ Its structure is input layer, convolutional layer, pooling layer, convolutional layer.....

## (Structure of CNN)

$\Rightarrow CL \rightarrow PL \rightarrow CL \rightarrow PL \dots \rightarrow FCL$

④ Receptive Field: How much portion of image a neuron has seen.



(Receptive fields of neuron is 9)

$\Rightarrow$  Receptive Size =  $1 + L \times (f - 1)$

$$\text{E.g.: } f = 1 \times L \times (3-1) = 1 + L \times 2$$

$$6 = L \times 2$$

$$L = 3$$

$$\text{E.g.: } 1000 = 1 \times L(3-1)$$

$$1000 = 1 \times 2L$$

$$\frac{1000}{2} = L$$

$$L \approx 500$$

$\Rightarrow$  For this, we use pooling and strided

$\Rightarrow$  Dense Layer are Learnable & pooling layer are not learnable.

<sup>(For Exam)</sup> In Pooling, variables or parameters will be 3x3

⇒ Max pooling & Avg Pooling  
(Pick Max) (Take Average)

⇒ Convolutional and pooling shrink the input layer.

(As we move forward, it spatial decreases and depth increases)

⇒ In CNN, first layer edges are bind and in intermediate layer we find small objects and last layers find whole object.  
(Bind)

## \* HomeWork : (LeNet - 5)

(Learnable layers)(CL & FC)

(Question for Exam) Find activation size, No. of parameters on each layers



$32 \times 32 \times 3$

Hyper Parameters

Experiment

: Make Correct Tab  
for HW, Quiz, Exam

8/5/25

## ① Architecture Designing:

→ LeNet-5 (for digit recognition)  
: Pooling layers are not included  
in counting because it is  
not learnable layer. (1998)

→ AlexNet (8 layer architecture):

① 2012.

② Calculation after filter is  
independent

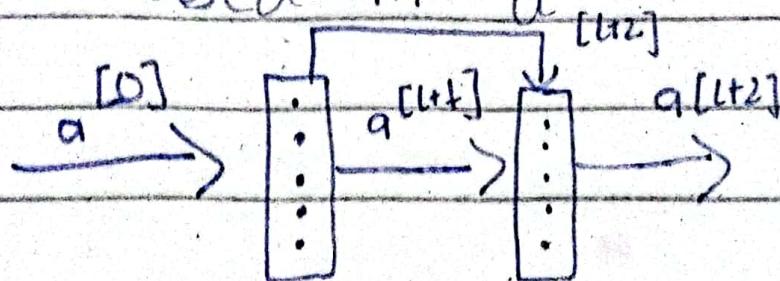
① Hardware and software available =  
in that time period and 3  
more layers added which  
makes network complex  
but pattern is same like  
now ( $\text{CON} \rightarrow \text{PL} \rightarrow \text{CON} \rightarrow \text{PL} \rightarrow \text{FC}$ )

→ VGG-16:

② simplified architecture  
③ Depth becomes double  
but size remains same  
and by pooling dimension  
becomes half.

→ ResNet (Residual Network):

④ skip connection idea is  
used in it



$$a^{[l+2]} = \text{RELU} (w^{[l+2]} * a^{[l+1]} + a^{[l]})$$

$$a^{[l+2]} = \text{ReLU}[a^{[l+1]}]$$

$$= a^{[l+1]}$$

: RELU  
 $(0, a)$   
 $- +$

→ " Scratch Training " (ii) Pretrained

: import torchvision.models

resnet models = models.resnet34

(pretrained=True)

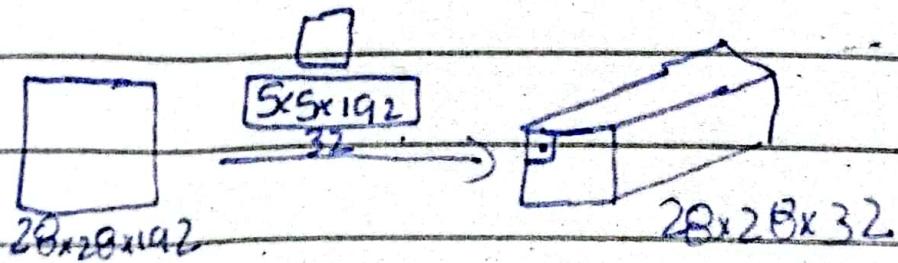
"if false then scratch."

→ Inception network :

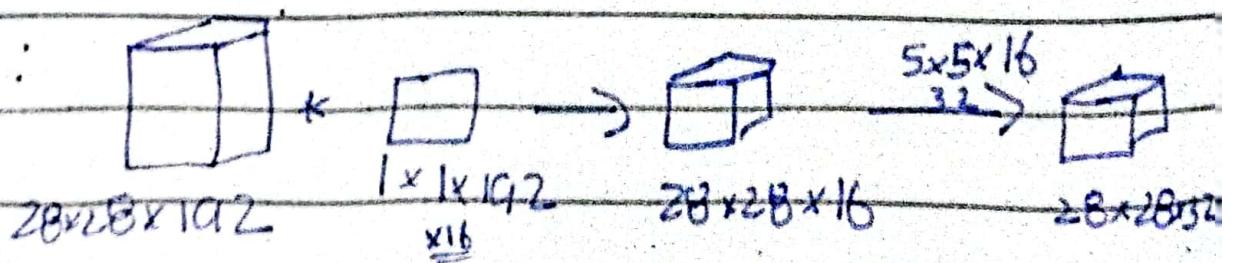
① Convolution operation. 1x1 is used (which is doing scaling).

② Depth wise work is done using it ( Retain, Increase ... )

③ In this , many filters are applied on one layer.



$$\Rightarrow 28 \times 28 \times 32 \times 5 \times 5 \times 192 = 120 \text{ million}$$



$$\Rightarrow 1 \times 1 \times 192 \times 28 \times 28 \times 16 \approx 12$$

: Skip Connection , 1x1 ...

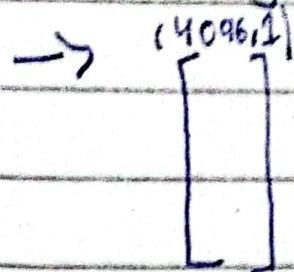
13/5/25

## ④ Object Detection:

→ Classification and localization  
(Not just identify but also locate it)

→ First layer find edges  
and intermediate layer  
find small objects.

→ AlexNet is trained for  
ImageNet (1000 images) (14 billion parameters)



→ 10 images of same object  
and their features are close to each other

(If pixel level is compared  
then patterns of some objects with different background will be different)

→ Feature comparison is better than pixel comparison.

→ Last layer contains complete semantic whole image.

→ Block ( $CL \rightarrow PL$ ) ( $[L \rightarrow CL-PL]$ )

## ① Transfer Learning:

→ To overcome overfitting,  
dataset size increases.

→ If dataset size is small  
then use transfer  
Learning.

→ Three ways:

(i) ① If dataset size  $< 10000$ ,  
then last layer will  
be fine tuned (trained)  
and all other <sup>ALONET layer</sup> will be  
frozen instead of last.  
(3 neurons in last layer only)

(ii) ② If dataset size  $> 10000$ ,  
last layer will be  
removed, first layer  
freeze, last four  
layers will be  
fine tuned and  
transfer it.

## : CelebriDataSet

(iii) ⚡ If dataset is 50000 or 100000 and classes 3, last layer will be removed and train all remaining layers.

→ If Model is there, but not trained and we have three classes and dataset small, then train it on related large dataset and then train it on your own data.

→ Assignment: (ResNet)

① Train last Layer sign Dataset  
(Last layer contains 10 previous classes)

→ Image → classification → Detection  
classification with localization [centralization]

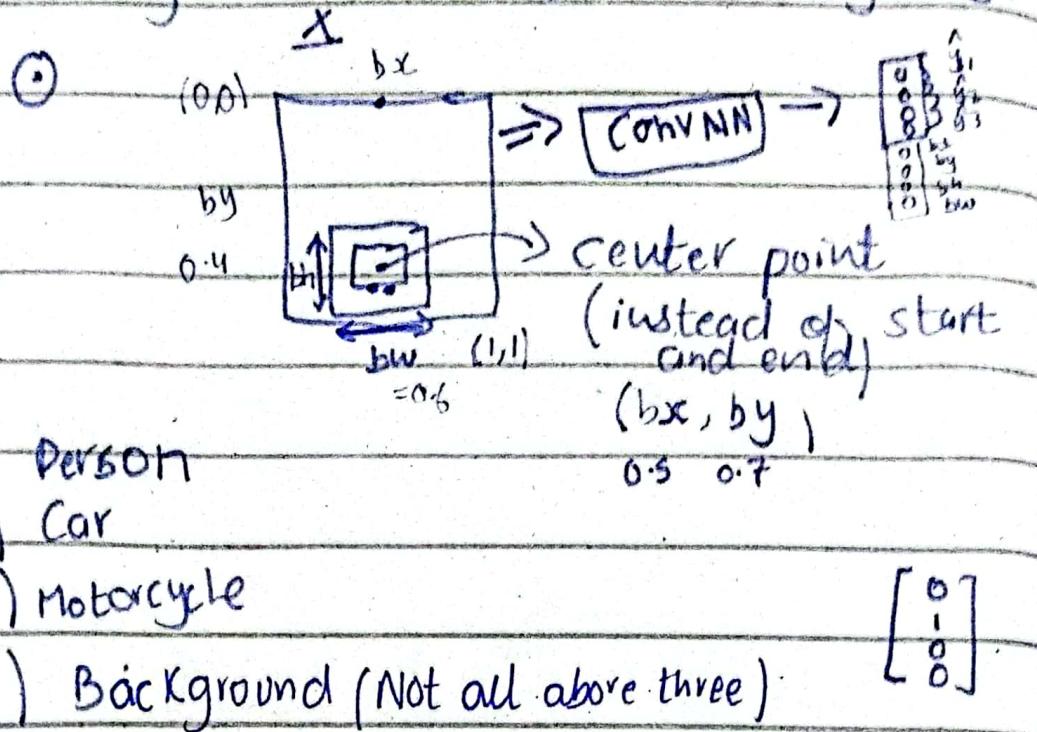
# 15/5/25

\* Self-Driving Car uses Object Detection

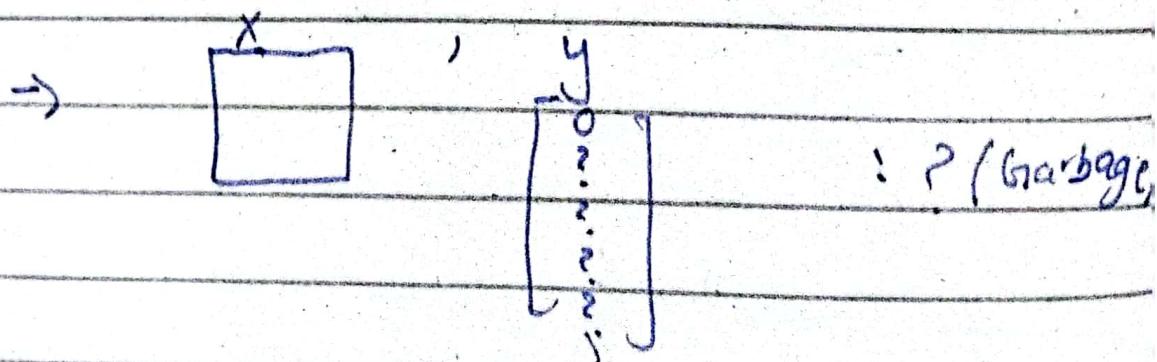
\* Classification → object detection  
(one object) (one image but multiple images)

## → YOLO:

- ① You only look once
- ② Self Driving uses this algorithm to detect objects.



$$y = \begin{bmatrix} pc \\ bx \\ by \\ bw \\ bh \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 0.5 \\ 0.7 \\ 0.4 \\ 0.6 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$



If we use  
sigmoid, then  
we use BCE

## ① Loss Function:

$$y_i (y_i = \pm 1)$$

$$L(y, \hat{y}) = \text{BCE} + (1-y) \log(1-\hat{y}) + (y - \hat{y})^2$$

$$\text{else : } || (y_i = 0)$$

$$L(y_i, \hat{y}_i) = (y_i - \hat{y}_i)^2$$

→ For First node apply BCE

after sigmoid

→ For other nodes <sup>(coordinates)</sup> we apply

MSE.

→ For last 3 nodes, classes probability  
apply softmax and use

## ② Landmarks Detection

### ③ 68 Landmarks on human faces

which is used by different filters  
etc... (glasses) to fit on face.



For smile  
Detection  
check  
Lips landmark  
Distance

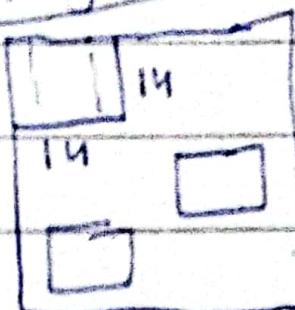
For  
Death  
Detection  
Check Eye  
Landmark  
Distance

⇒ We can decide pose  
by landmarks.

## ① Car detection:



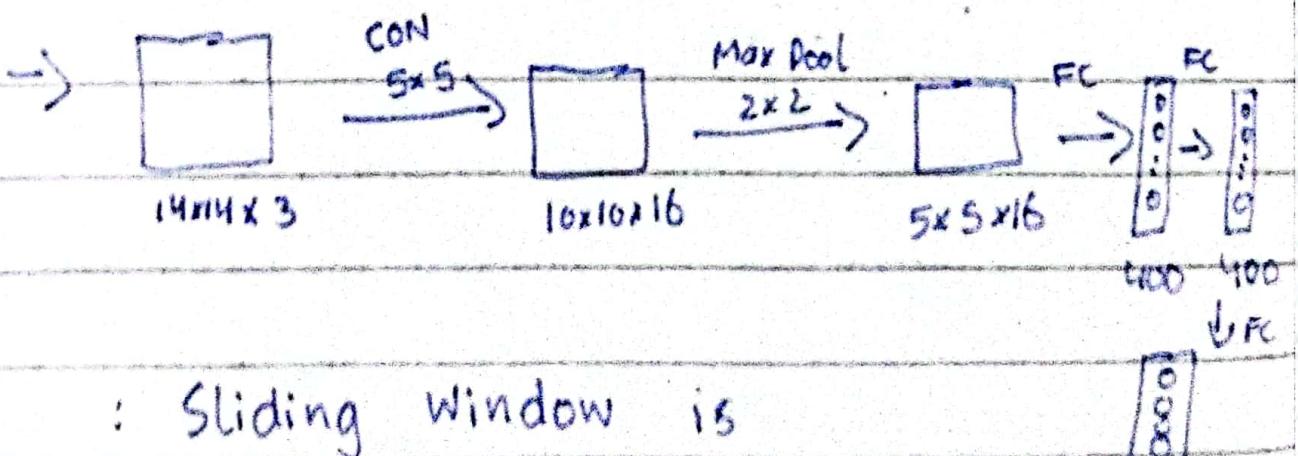
→ Closely cropped image (only images in box frame)  
sliding window



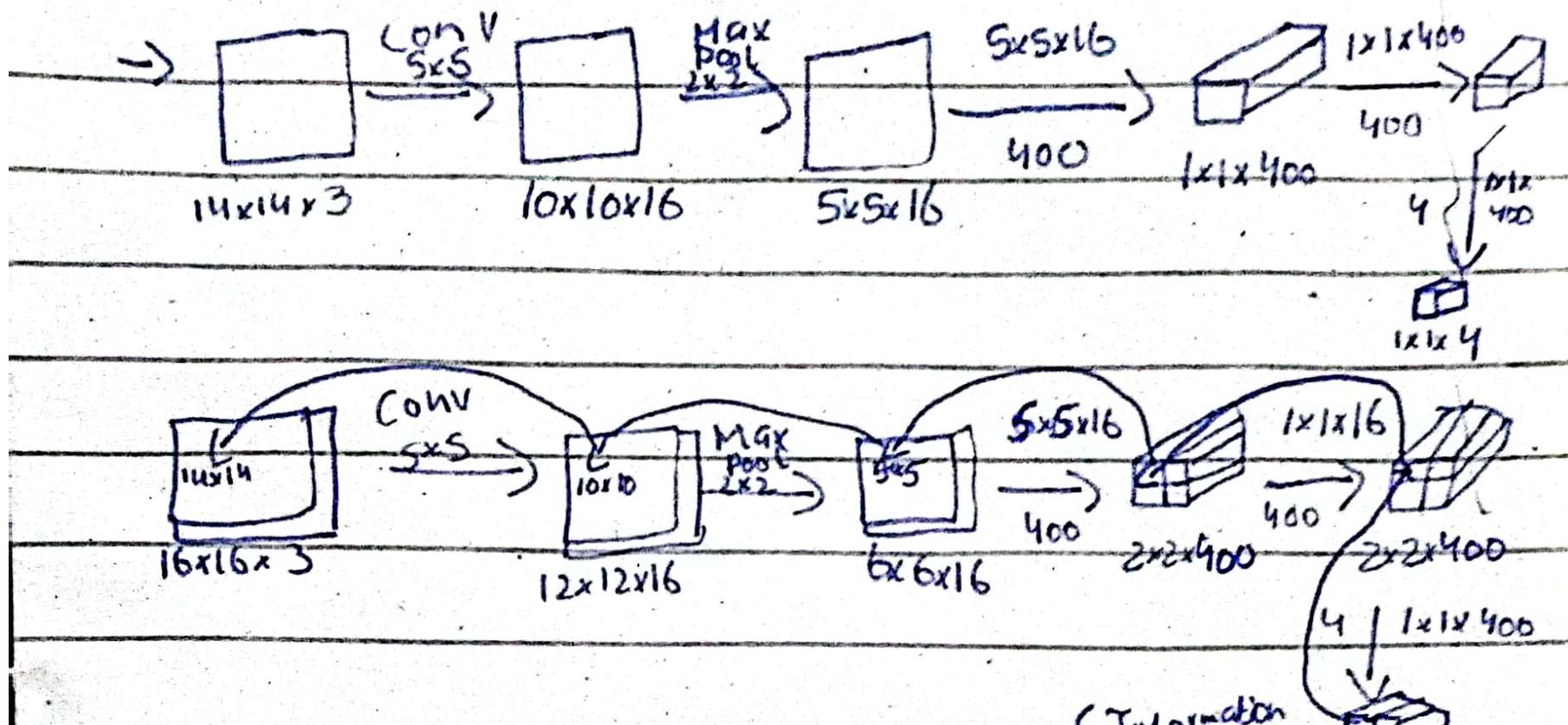
- : If image of car size is greater than 14x14, then image size will be lesser to detect properly.
- : Big Image (contains multiple objects)

: (Computations Increases)  $\frac{[n+2p-f]}{s} \cdot 1$

→ Solution is YOLO which looks <sup>one</sup>.



: Sliding Window is replaced by Fully Convolutional Neural Network!



Now Check for  
receptive field

(Information  
about  
4 region)

$2 \times 2 \times 4$

(You Only Look Once) (Classify done) (Box Done)

→ It is refined by max suppression etc...

20/5/25

① Evaluation:

→ Label and prediction

intercept region

→ Sample IOU.

$$\text{IOU} = \frac{\text{size of intercepted area}}{\text{size of union area}}$$

;  $\text{IOU} \geq 0.5$  (satisfactory)

$\text{IOU} \geq 0.7$  (Good)

$\text{IOU} \geq 0.9$  (Best)

## ④ Non-Maximum Suppression:

- Drawing tangent
- We see the region where center is there
- The values which are not of big
- The closest or neighborhood region should be greater.
  - If values is 0.7, suppress it.
  - If values is 0.17, then finish it and leave
- Select maximum and suppress others.

: Exam  
Imp Question

## ④ Example:

: Intercept  
over Union

- There are 4 rectangles  
 $0.9, 0.8, 0.75, 0.7$

- 1) Select highest one 0.9
- 2) Find, its IOU with others  
Supress angles with 0.7 and output will be 0.9.

## ④ Anchor Boxes:

- If center of gravity of two images is

same, so it's difficult to handle, so we use anchors:

PC
b <sub>x</sub>
b <sub>y</sub>
b <sub>h</sub>
b <sub>w</sub>
c <sub>1</sub>
c <sub>2</sub>
c <sub>3</sub>

anchor box

①

PC
b <sub>x</sub>
b <sub>y</sub>
b <sub>h</sub>
b <sub>w</sub>
c <sub>2</sub>
c <sub>3</sub>

Car

PC
b <sub>x</sub>
b <sub>y</sub>
b <sub>h</sub>
b <sub>w</sub>
c <sub>1</sub>
c <sub>2</sub>
c <sub>3</sub>

Person

1
b <sub>x</sub>
b <sub>y</sub>
b <sub>h</sub>
b <sub>w</sub>
0

1
b <sub>x</sub>
b <sub>y</sub>
b <sub>h</sub>
b <sub>w</sub>
1

: PC  
(Probability)

If there  
is no  
image,  
then undefined  
and  $PC = 0$

② Objects can be 3 but anchors will be 2

③ Call NMS for Person object and then for Car.

④ Mean Average Precision:

→ Important Metric

→ It follows classification and IOU and then gives value.  
classification  
and  
localization

→ Model is trained and now we wanted to test it using (mAP)

→ Now check for each image  $\Rightarrow$  confidence, TP or FP

(i) TP if ( $P_c > 0.5$  & IoU between label and prediction is  $> 0.5$ )

→ (i) Write all images score together.

(ii) Sort in descending order

(iii) Precision or Recall:

		+ GT	
Predict		TP	FP
+ GT	-	FN	TN

$$\therefore \text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\therefore \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\therefore \text{Accumulated} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

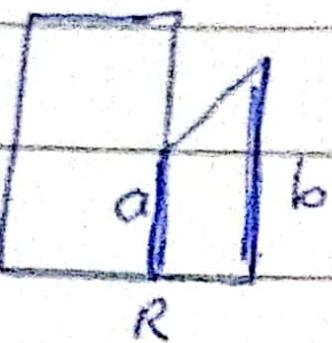
$$\therefore \text{Accumulated} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\therefore \text{Precision} = \frac{\text{AccFP}}{\text{AccTP} + \text{AccFP}} = \frac{1}{1+1} = \frac{1}{2}$$

$$\therefore \text{GT} = \text{U}_{\text{objects, images}}$$

$$\therefore \text{Recall} = \frac{\text{TP}}{\frac{\text{TP} + \text{FN}}{\text{GT}}}$$

(iv) Plot the Precision-Recall graph  
(Any side of precision or recall  
on graph)



: calculate  
from graph

22/5/25

## \* Face Recognition:

- One to one match
- Checks face data from database (Mobile Face Recognition)

## (\*) Face Verification:

- Multiple faces or data present in database. (<sup>multiple</sup><sub>faces</sub>)

: Link for face detection project present in slides.

## \* Face Recognition System:

- Employee data of company (4 employees and wanted to recognize their face and apply neural network).

- If employees increases, then again training and testing is required. If data is less, training underfit.  
So more data is required to train. (One-shot learning)
- The solution to one-shot Learning is to learn similarity function.
- If one image is passed one feature is made and if again that image is passed, then feature comparison is done and check for less difference.
- Just trained the model over big dataset and if new person comes, then compare it on feature vector and check for difference.

: code

nn.Sequential (Run layers in sequence)

## ① Triplet loss:

→ To train Siamese network,  
we use set of 3 images:

: Anchor + Positive + Negative  
original image      same image      different image

$$: \underbrace{d(A, P) + \alpha}_{0.5} - d(A, N) \leq 0 \quad 0.2 \quad 0.51$$

$$0.2 \leftarrow 0$$

→ Two Branches in Siamese, in  
which we are training model

→ Loss function:

$$L(A, P, N) = \max(d(A, P) + \alpha - d(A, N), 0)$$

(positive) 0.5      0.2      0.51  
(negative) 0.5      0.2      0.9

→ Selection of triplet:

① Anchor and Positive should be  
different in look to train  
network hard

② Anchor and Negative look  
same to train hard.

$$\textcircled{1} \quad z = w \boxed{x} + b$$

$$z = w \left( \frac{f(x)^{(i)} - f(x)^j}{\|f(x)^{(i)} - f(x)^j\|^2} \right) + b$$

• Feature Learning  
By Autoencoder:  $\rightarrow$  Image retrieval

## ① Deep Face:

$\Rightarrow$  If matched 1, otherwise 0

so deep face is much better than face Net.

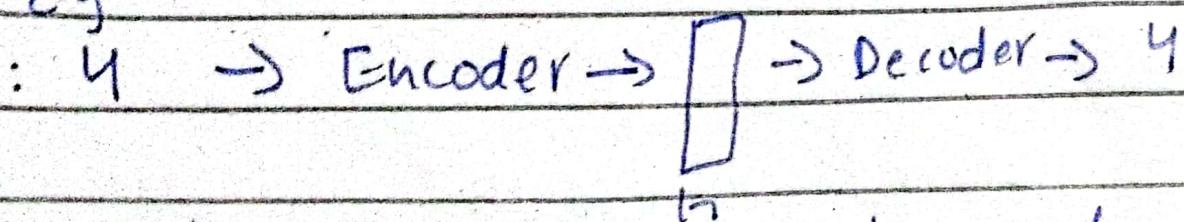
- ① FaceNet Pre computed and pass through neural network and check difference for comparison.

## ② Autoencoder Architecture:

$\rightarrow$  It has two parts:

(i) Encoder      (ii) Decoder  
(compress info in small vector)      (generate original image)

E.g.



- ③ CNN can also be used and layers of convolution is used.

## ④ Application:

$\rightarrow$  Noisy image given or good given, then give noisy image to encoder and it makes original image.

$\rightarrow$  Image colorization (old images into RGB)

$\rightarrow$  Anomaly detection (Input and Output difference will be more then anomaly. image is same just one same image)