Deep learning computer vision is Relpful in self driving cars, better face recognition.

Applications: Image Classifica^, object Detec^, Neural style transfer
(content image + style image)

Problem ε̄ CNN:

If i/p image is (1000 × 1000) resolu^ image, then i/p features = 1000 × 1000 × 3 = 3 m
i/p features. If hidden layer with 1000 units is used;

$$W^{[1]} = (1000, \; 3\text{million}) \text{ dimensional matrix} \quad \text{i.e } 3 \text{ billion parameters}$$

So, many parameters lead to overfitting in absence of sufficient amount of data. Also, memory and computational requirements of 3 billion parameters is quite infeasible. we use convolu^ computation to solve this problem.

→ Edge detec⁰ example to understand convolu⁰ operation:-

Note: A gray scale image, is just 2-D box no RGB codes. Every pixel has either a value 1 or 0 i·e block or white. [or may be some other nos. that are shades of black and white.]

matrix

The convolu⁰ oper⁰ detects vertical & horizontal lines in the i/p image.

e.g.

| 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 |

(6×6) i/p image.

(a)

"Convolu⁰"

*

| 01 | 0 | -1 |
|---|---|---|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

3×3 filter/kernel

(b)

=

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

4×4 vertical edge ouput.

(c)

$C_{11}$ is computed by doing element wise multiplica⁰ of (a) and (b) and adding all the 9 terms.

(1st window)

To get $C_{12}$, shift window 1 step to the right and repeat above procedure. [i·e· $(3×1)+(1×1)+(2×1)+(0×0)+(5×0)+(7×0)+(1×-1)+(8×-1)+(2×-1)]=-4$

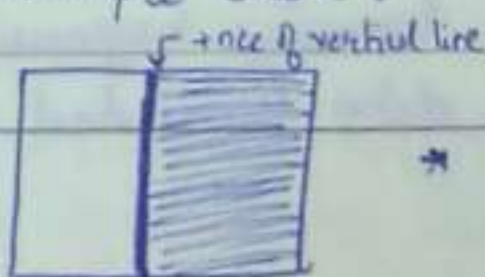Hence, a 6×6 matrix convolve ⊗ 3×3 matrix gives 4×4 matrix.

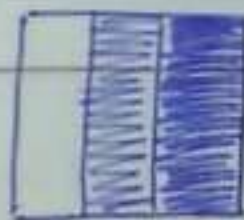To implement the above convolu⁰ fn:
in python we have: conv-forward
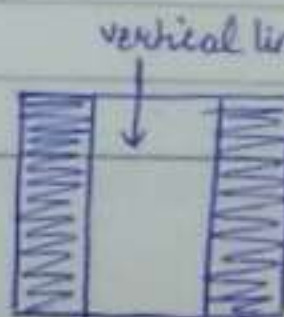in tensorflow " " : tf.nn.conv2d
in Keras " " : Conv2D

Note: Above works as a vertical edge detector. This can be shown through the example below.

+ve of vertical line

vertical line detected



i·e·

| | | | | | |
|---|---|---|---|---|---|
| 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 |

*

| | | |
|---|---|---|
| 1 | 0 | -1 |
| 1 | 0 | -1 |
| 1 | 0 | -1 |
| 1 | 0 | -1 |

(edge detec$^n$ filter)

$\longrightarrow$

| | | | |
|---|---|---|---|
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |

(vertical edge detector)

This indicates the +nce of a strong vertical line in middle of the image. The dimensions seem to be a little bit wrong bcoz the i/p image is just 6×6. If it was of higher dimension, better results would be obtained.

Note: Here, +30 shows that there is a light to dark transition.

If ▨ was the i/p image, we would have got −30 showing transi$^n$ from dark to light. We can take absolute of 30 if we don't care which of the two cases it is.

Also

| | | |
|---|---|---|
| 1 | 1 | 1 |
| 0 | 0 | 0 |
| -1 | -1 | -1 |

Horizontal edge detector

Note

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

(Sobel filter)

(another vertical edge detector giving more weightage to the central row making it more robust)

| | | |
|---|---|---|
| 3 | 0 | -3 |
| 10 | 0 | -10 |
| 3 | 0 | -3 |

(Scharr filter)

Note for very big images, these 9 values of the filter are learnt through back propagation and treated as parameters. They can not only be used to detect vertical & horizontal edges but edges at some angles as well like 45° or 73° or so on. This learning happens through data provided as i/p.

→ **Padding:** modified to convolu¹ opera².

$$(n \times n) \underset{\downarrow}{*} (f \times f) \rightarrow ((n-f+1) \times (n-f+1))$$

convolve

- One drawback in what we were doing is that the pixels at the corners and the edges are considered only once as only 1 filter covers it while the pixels at the center are covered by many $(f \times f)$ filters which make us to ignore a lot of info near the edge or corner.

- Other drawback is that the output obtained is a shrinked image, and so if on every hidden layer it will shrink this way, we will get an extremely compressed image as the final output.

To fix these problems, we can pad the image with an addition of border of 'p' pixel all around the image.

In case of $6 \times 6$, $p = 1$ as $(n-f+1) = 6$ i.e same sized image as output
$$(8-3+1)$$
{n = 6+2 , bcoz of padding}.

Convolution ⟨ valid convolu¹ (no padding)
same convolu² (padding such that o/p size = i/p size)

$$\downarrow$$

i.e $n + 2p - f + 1 = n$ → $\boxed{p = \dfrac{f-1}{2}}$

{f - filter size}.

f is almost always add in computer vision conven².

* So, that p is not fractional and we have a central pixel to refer- ▦

→ **Strided convolution:**

Stride = 2 means the $f \times f$ filter that is used is hopped by 2 steps instead of 1 in the $n \times n$ input image.

$$(n \times n) \underset{\downarrow}{*} (f \times f) \longrightarrow \left(\left[\frac{n+2p-f}{s}+1\right] \times \left[\frac{n+2p-f}{s}+1\right]\right)$$

convolve     Padding 'p'         $\hookrightarrow$ floor is taken if it is not

Stride 's'                  an integer.

<u>Note</u> the $f \times f$ filter must lie completely inside the image while hopping over it to do the computations.

<u>Note</u> <u>Cross correlation v/s convolution</u>

what we have doing till now was cross correla?. Convolu? in actual involves one more step i.e flipping of the filter matrix.
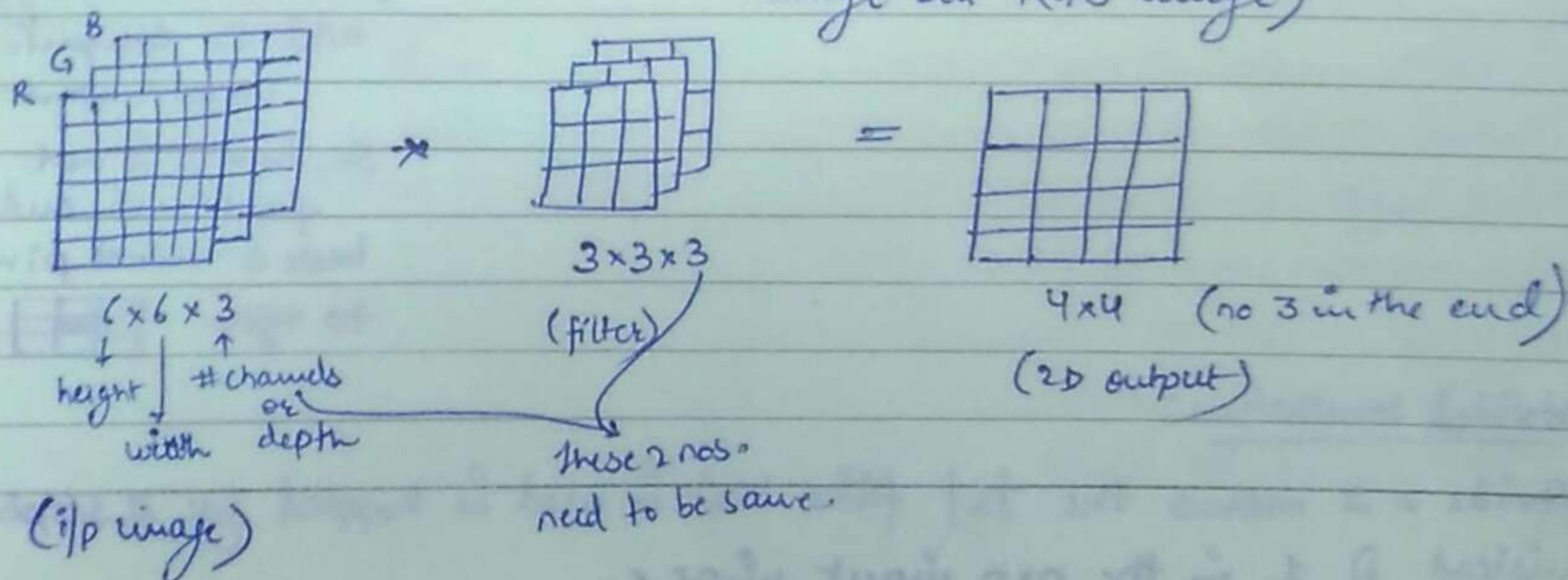
| 3 | 4 | 5 |
|---|---|---|
| 1 | 0 | 2 |
| -1 | 9 | 7 |

$\longrightarrow$

| 7 | 2 | 5 |
|---|---|---|
| 9 | 0 | 4 |
| -1 | 1 | 3 |

(Narrowing of filter both on vertical and horizontal axis)

the element wise product and summing is actually done after this step but generally we skip this step and call that also as convolu? instead of cross correla?.

$\Rightarrow (A * B) * C = A * (B * C)$      Associativity of convolu?.

$\rightarrow$ <u>Convolutions over volumes:</u> (to detect feature not in a grey scale image but RGB image)



$6 \times 6 \times 3$

height | #channels

    or

width   depth

(i/p image)

$3 \times 3 \times 3$

(filter)

these 2 nos.
need to be same.

$4 \times 4$    (no 3 in the end)

(2D output)

This $3 \times 3 \times 3$ cube is swiped over the i/p image and this time 27 nos. are corresponding multiplied and added to give a number of the $4 \times 4$ o/p matrix. Similar to the previous grey scale image convolu^, only difference is that it is in 3D.

To detect vertical edges only in the red channel, $3 \times 3 \times 3$ filter is like:



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

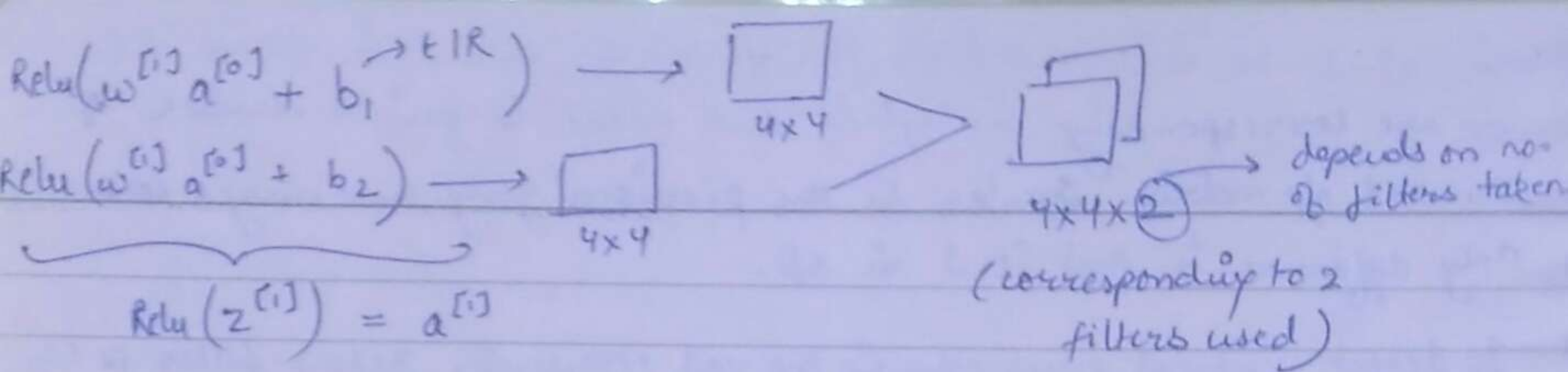$\qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \quad \hookrightarrow$ no. of filters used.

Note Adv. of convolution volume is that now we can use multiple filters thus detecting multiple features in the input image.

no. of filters used = no. of channels in o/p image.

(e.g. horizontal filter + vertical filter, then $n_c = 2$) 1st channel giving horizontal edges and 2nd channel giving vertical edges detected in the i/p image.

→ **Convolution NN layer :**



$(6 \times 6 \times 3)$

$a^{[0]}$

$(3 \times 3 \times 3)$

$(4 \times 4)$

$(3 \times 3 \times 3)$

$\underbrace{\qquad}_{w^{[1]}}$

$(4 \times 4)$

$\underbrace{\qquad}_{w^{[1]} a^{[0]}}$

$Relu\left(w^{[1]}a^{[0]} + b_1 \xrightarrow{} \in IR\right) \longrightarrow \square$ $4 \times 4$

$Relu\left(w^{[1]}a^{[0]} \pm b_2\right) \longrightarrow \square$ $4 \times 4$

$\underbrace{\phantom{Relu\left(w^{[1]}a^{[0]} \pm b_2\right)}}$

$Relu\left(z^{[1]}\right) = a^{[1]}$

$\square$ $> $ $\square$ $4 \times 4 \times ②$ $\xrightarrow{}$ depends on no. of filters taken

(corresponding to 2 filters used)

Hence, activaⁿ for next layer obtained.

Convoluⁿ step was basically the linear operaⁿ step i.e $w^{[1]}X$, then bias was added giving $z^{[1]}$ and then $\sigma(z')$ gave $a^{[1]}$.

**Q.** Calculate no. of parameters if we have 10 filters that are $3 \times 3 \times 3$ in 1 layer of a NN.

**Ans** $\left((3 \times 3 \times 3) + \underset{\underset{bias}{\downarrow}}{1}\right) \times 10 = 280$ parameters.

**Note** No matter how big the i/p image is, no. of parameters $= 280$ is fixed for 10 $3 \times 3 \times 3$ filters. This property saves CNN from overfitting.
$\phantom{xxxxxxxxxxxxxxxxxxxx}\hookrightarrow$ to extract
$\phantom{xxxxxxxxxxxxxxxxxxxxxxxx}$ different features
$\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxx}$ from images.

**Notations:**

$f^{[l]}$ = filter size (i.e $f \times f$ filter used for layer $l$ of NN)

$p^{[l]}$ = padding in layer $l$

$s^{[l]}$ = stride in layer $l$.

Input : $n_H^{[l-1]} \times n_w^{[l-1]} \times n_c^{[l-1]}$ {i.e activaⁿ from the previous layer}

Output : $n_H^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

$$n_H^{[l]} = \left[\frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1\right] \quad \text{(Similarly for } n_w^{[l]}\text{)}$$

$n_c^{[l]}$ = no. of filters. and each filter is $f^{[l]} \times f^{[l]} \times n_c^{[l-1]}$

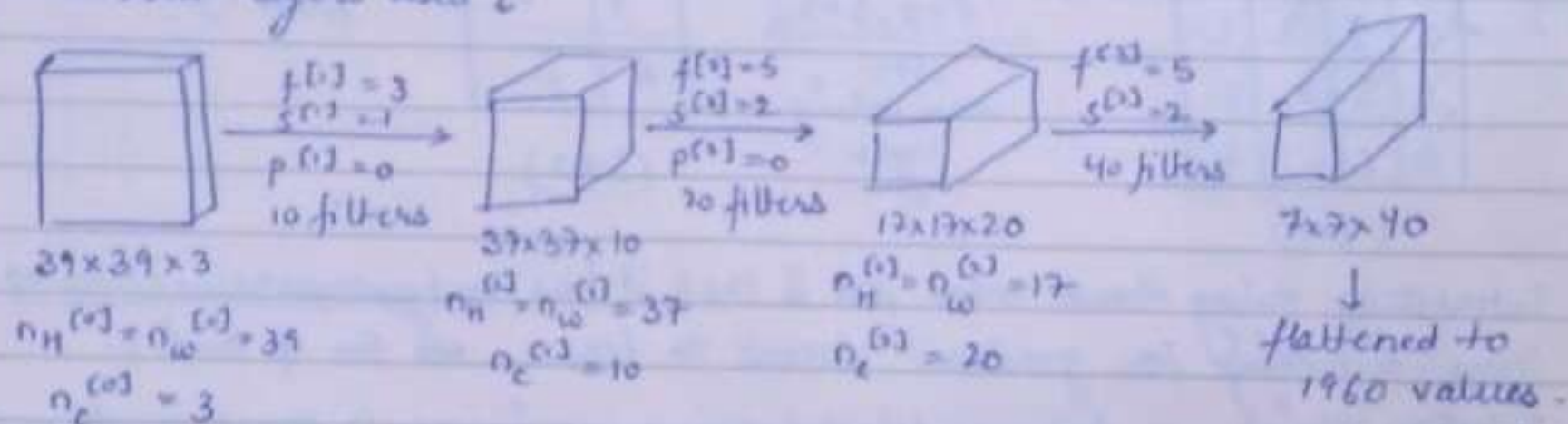$$a^{[l]} = n_H^{[l]} \times n_w^{[l]} \times n_c^{[l]}$$

$$A^{[l]} = m \times a^{[l]} \quad \{ \text{for } m \text{ inputs} \}.$$

weights : $\underbrace{f^{[l]} \times f^{[l]} \times n_c^{[l-1]}}_{\text{dimension}} \times \underbrace{n_c^{[l]}}_{\text{# filter}}$

of filter.

bias : $n_c^{[l]} \to (1,1,1, n_c^{[l]})$

→ Simple CNN example : (Conv Net)

3 convol⁴ layers used 2



$f^{[1]} = 3$
$s^{[1]} = 1$
$p^{[1]} = 0$
10 filters

$f^{[2]} = 5$
$s^{[2]} = 2$
$p^{[2]} = 0$
20 filters

$f^{[3]} = 5$
$s^{[3]} = 2$
40 filters

$39 \times 39 \times 3$

$37 \times 37 \times 10$

$17 \times 17 \times 20$

$7 \times 7 \times 40$

$n_H^{[0]} = n_w^{[0]} = 39$

$n_c^{[0]} = 3$

$n_H^{[1]} \times n_w^{[1]} = 37$

$n_c^{[1]} = 10$

$n_H^{[2]} = n_w^{[2]} = 17$

$n_c^{[2]} = 20$

↓
flattened to
1960 values.

Hence, $(39 \times 39 \times 3)$ i/p image converted to $(7 \times 7 \times 40)$ features for this image.



→ $\hat{y}$ (final o/p)

logistic / softmax

1960 values

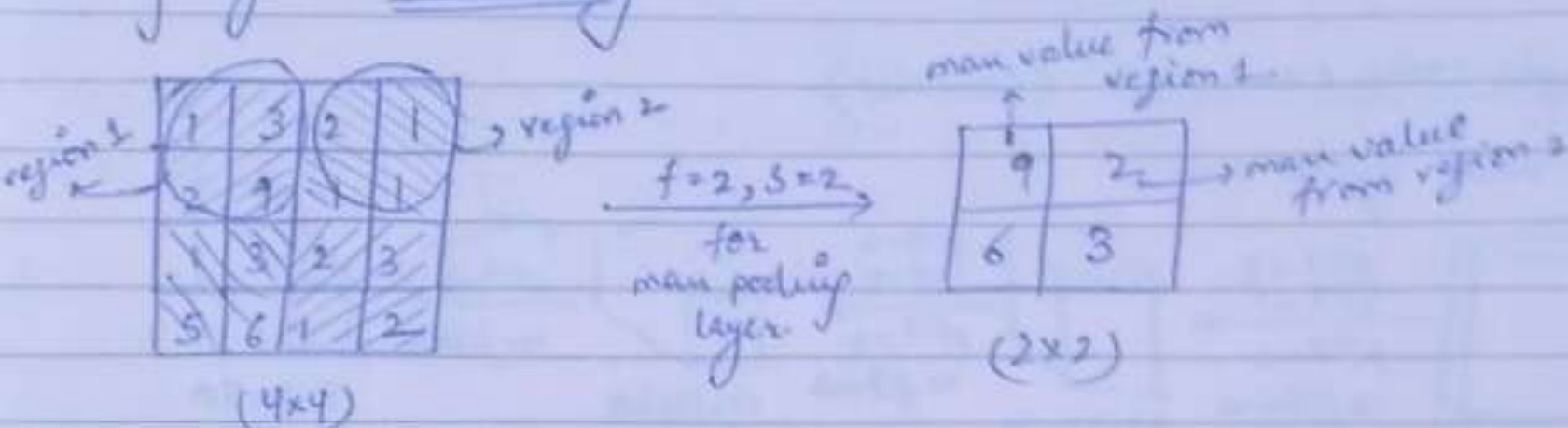Simple regression applied to 1960 values giving features of the i/p image.

Note As we go deeper in CNN, $n_H$ and $n_w$ generally ↓ while $n_c$ tes. Main work in CNN is to decide the hyperparameters i.e $f^{[l]}, s^{[l]}, p^{[l]}$, no. of filters.

Note Although its possible to design a good CNN using just convolution layers but most NN architectures also have a few pooling layers and a few fully connected layers.

So, Types of layers in CNN — Convolution
                                    Pooling
                            Fully connected.

① Pooling layer: to reduce the size of the representation, to speed up computation

Pooling layer : Max Pooling



$(4 \times 4)$ → region 2

$\xrightarrow{f=2, s=2,}$ for max pooling layer

$(2 \times 2)$

max value from region 1

| 9 | 2 |
| 6 | 3 |

→ max value from region 2

Interesting thing about max pool is that it has no hyperparameters to learn i.e. nothing for gradient decent to learn. we fix $f$ and $s$.

__Intuition behind using maxpool__ :  features detected anywhere in one of these quadrants remain preserved in the o/p of max pooling. If feature is not detected, the maximum value in quadrant itself is small.

for max pooling also, $n \cdot \dfrac{s + n + 2p - f}{s}$ is the o/p size where $n \times n$ is the i/p size.
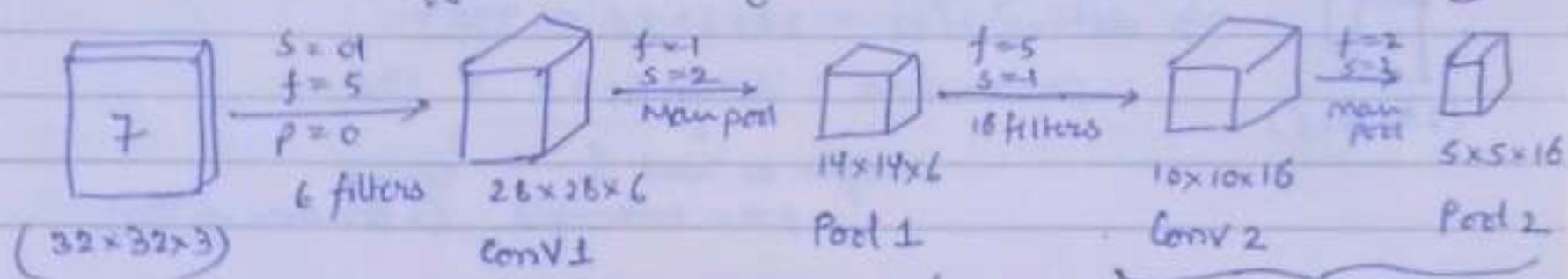
__Max pooling in 3D__: If $n \times n \times n_c$ is i/p then $n' \times n' \times n_c$ is o/p. The filter is applied independently to all the channels, unlike in case of 3D convol" where o/p was 2D even for 3D output and 3D filter. Here, in pooling we have 3D output for 3D input & 2D filter applied to each channel independently.

__Note__ Avg. pooling is also used but not much where avg. of quadrant is taken instead of max value.

$f=2, s=2$ are generally used to shrink the height x width to approximately $\frac{1}{2}$ in the o/p. $p=0$ generally in max pooling.
i·e (of half size)
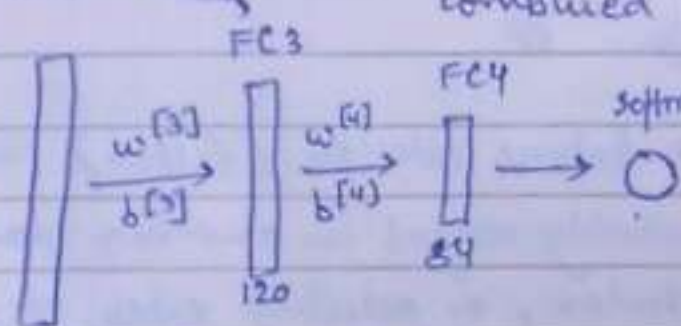
**Hyperparameters:** $f, s$, max/avg pooling
(binary 0 or 1)

→ Complex CNN example: (inspired by LeNet-5)
we need to identify which no. from 0 to 9 is +nt in the image.



$32 \times 32 \times 3$

$\begin{array}{l} s=01 \\ f=5 \\ p=0 \end{array}$ 6 filters

$28 \times 28 \times 6$
Conv1

$\begin{array}{l} f=1 \\ s=2 \end{array}$ Max pool

$14 \times 14 \times 6$
Pool 1

$\begin{array}{l} f=5 \\ s=1 \end{array}$ 16 filters

$10 \times 10 \times 16$
Conv 2

$\begin{array}{l} f=2 \\ s=3 \end{array}$ max pool

$5 \times 5 \times 16$
Pool 2

called 1 layer of the NN, though parameters hyper are taken twice, but weights & biases are considered only once during the convolu operⁿ hence it is combined called 1 layer.

layer 2 of the NN

fully connected (layer 5)

FC3



$\xrightarrow[b^{[3]}]{w^{[3]}}$

120

400 i·e (5x5x16 flattened as 400 pixels)

FC4

$\xrightarrow[b^{[4]}]{w^{[4]}}$

84

$\longrightarrow$ O  softmax (० 10 outputs in this case i·e 0 to 9)

$w^{[3]} = (120, 400)$ dimensional
$b^{[3]} = (120)$
$w^{[4]} = (84, 120)$
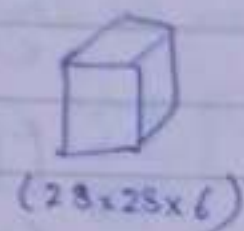$b^{[4]} = (84)$

{just like the original NN that we saw before CNN}
Fully connected as every i/p neuron is connected to every o/p neuron.

Common CNN pattern. (no. of layers may vary).

Conv layers → Pooling layer → Conv layers → Pooling layer → Fully Connected layer → Fully connected layer → Softmax.

Note Activⁿ size drops as we go deep in NN. No. of parameters is
0 for Pooling layer, too many for FC layer and relatively small
for CNN as we discussed earlier. The drop in activⁿ size should
not be too quick.

$$\Rightarrow \text{activⁿ size} = 28 \times 28 \times 6 = 4704$$

$(28 \times 28 \times 6)$

$$\text{parameters} = (5 \times 5 \times 6) + 6 = 156$$

weights as used in prev example    biases (1 for each filter)

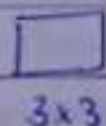→ Why convolⁿ over fully connected?

i) Less no. of parameters needed, hence less chance of overfitting.

In above e.g. conv needed 156 parameters. But if FC was used
instead of conv here then $4704 \times (32 \times 32 \times 3)$ parameters i.e 14 million
parameters would be needed. ↳ o/p.   ↓ i/p

Conv uses parameter sharing bcoz a feature detector (filter) that's
useful in one part of the image is probably useful in another part of
the image, like the vertical edge detector, or detecting eyes, etc as well.
Also, in conv, in each layer each o/p value depends only on small
no. of i/p i.e sparse connections unlike in FC, i.e.



6×6   *   3×3   =   4×4   → this o/p depends only on
the shaded region of the
i/p. No other pixels affect
this o/p.

ii) It captures the desirable property of translation invariance.
i.e an image of a cat slightly shifted is labelled same in o/p
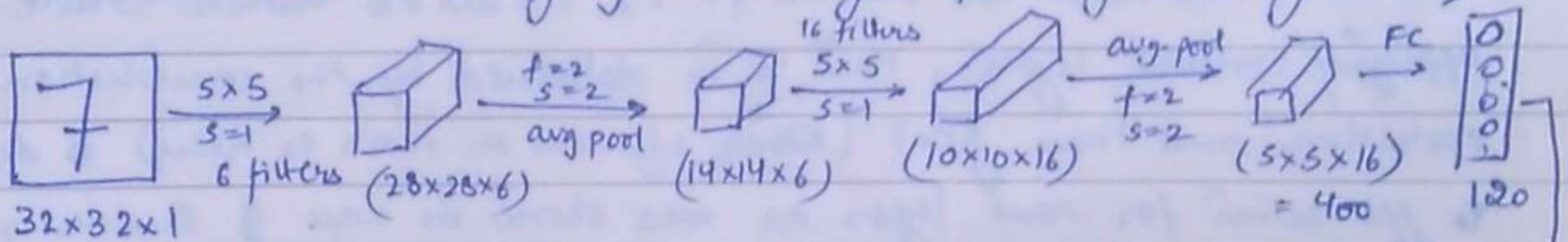as the image ī cat in middle.

$$\text{Cost } J = \frac{1}{m} \sum_{i=1}^{m} L(\hat{y}^{(i)}, y^{(i)}) \qquad \{m = no. \text{ of training examples}\}$$

Now, use gradient decent to optimize parameters to reduce J.

⇒ Some classic NN:
- Le Net-5
- Alen Net  } → lay the foundā for modern computer vision.
- VGG Net
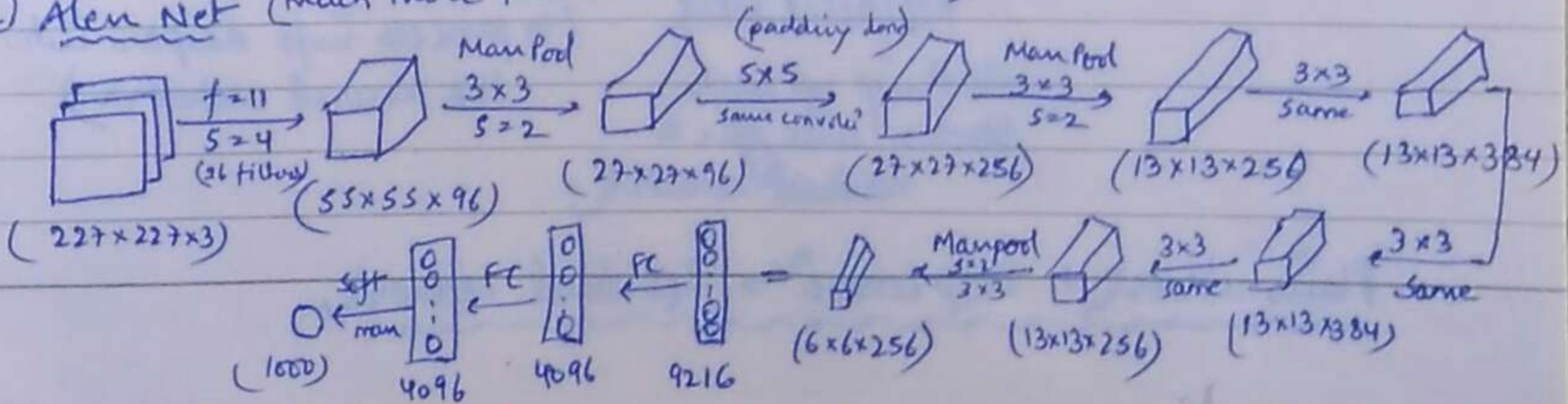- Resnet (152 layers)
- Inception

① Le Net-5 (trained on grey scale images for digit recognition).



Note A modern version of this NN
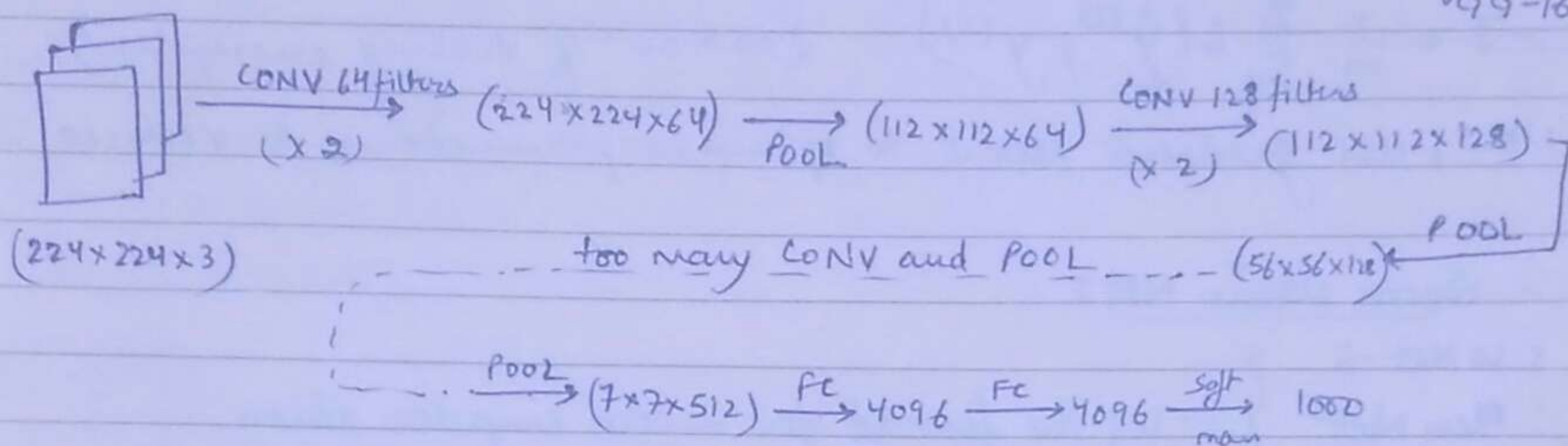– will use a softman layer ī a
10 way classificā output.

(can take
10 values from
0 to 9)

② Alen Net (much more parameters than le Net-5)

VGG-16 net (simplified architecture) has around 138 million parameters and thus a pretty large network.

(16 layers having weights) { Conv = 3×3 filter, S=1, same convolu^ (i.e padding done) } Always in this VGG-16.
{ maxpool = 2×2, S=2

(224×224×3)

$\xrightarrow[\text{(×2)}]{\text{CONV 64 filters}}$ (224×224×64) $\xrightarrow[\text{POOL}]{}$ (112×112×64) $\xrightarrow[\text{(×2)}]{\text{CONV 128 filters}}$ (112×112×128)

--- too many CONV and POOL --- (56×56×128) $\xleftarrow{\text{POOL}}$

--- $\xrightarrow{\text{POOL}}$ (7×7×512) $\xrightarrow{\text{FC}}$ 4096 $\xrightarrow{\text{FC}}$ 4096 $\xrightarrow[\text{max}]{\text{Soft}}$ 1000
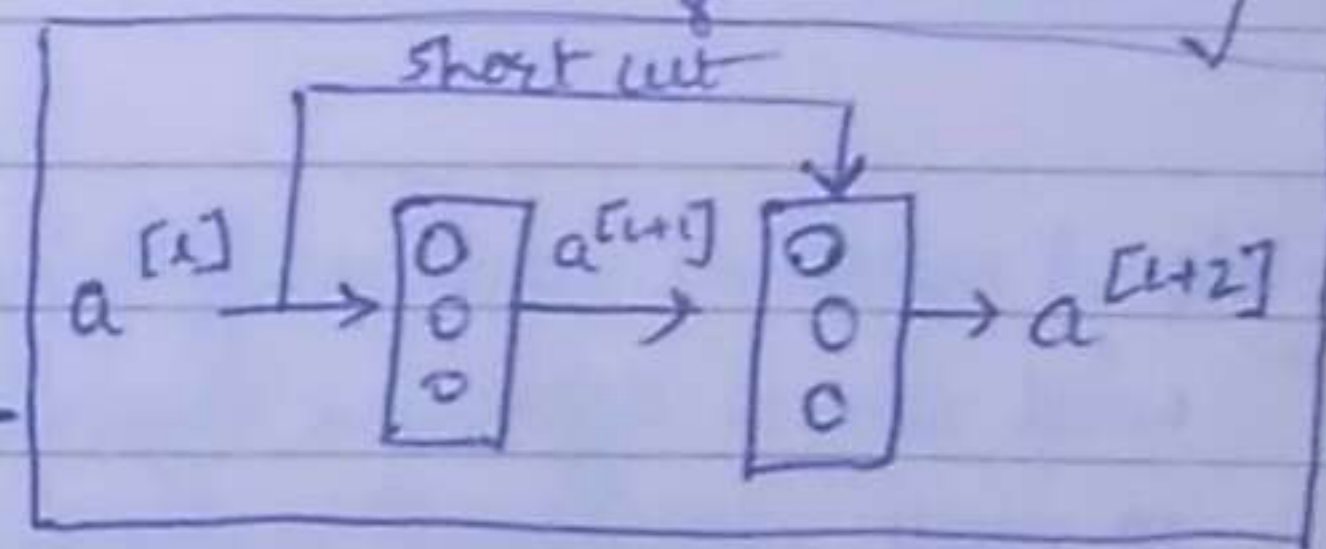
① ResNet (enables Us to train very deep NN, sometimes upto 100 layers)

Note ResNets are built out of Residual blocks.

Note In CNN, we don't use linear fN i.e z = wx + b rather while applying convolu^ layer, this z is obtained by the convolution operation and then g(z) {either sigmoid or tanh or Relu) is done to get activa^ for next layer as was done in case of linearity.

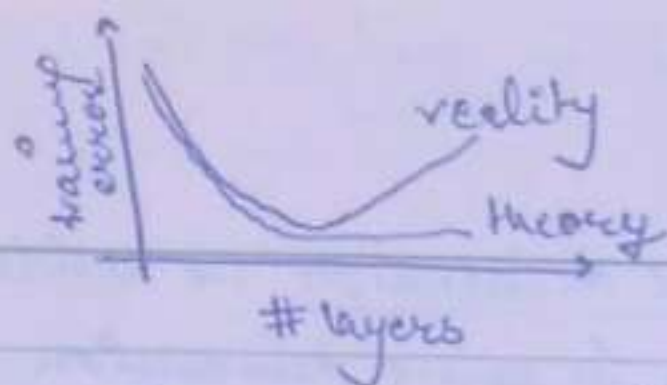In residual, $a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$



short cut = skip connection
(to process info deeper into the Neural Network)

Residual block
Stacking up such
blocks will give a
residual network.

Plain network + skip conne^ = Residual network

**reality**

**theory**

training error ↑

# layers →

(for Plain Network)

---

training error ↑

# layers →

**But!**

reality as well as theory

(Residual Network)

* { helps with the vanishing and exploding gradient problems }.

---

→ **why resnets do well?**

Going much deeper in the NN might hurt the network's ability to train the network to do well on the training set but this is much less true while training a ResNet.

$X \longrightarrow \boxed{Big\ NN} \longrightarrow a^{[l]}$

$X \longrightarrow \boxed{Big\ NN} \xrightarrow{a^{[l]}} \bigcirc \longrightarrow \bigcirc \longrightarrow a^{[l+2]}$

{ Supposing we use ReLu so, all activa$^n$ will be $\geq 0$ }

So, $a^{[l+2]} = g\left(z^{[l+2]} + a^{[l]}\right)$ —— ①

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

$L_2$ regularizi tends to shrink $w$ and $b$ with tiny layers and so at some point they might become 0

giving:

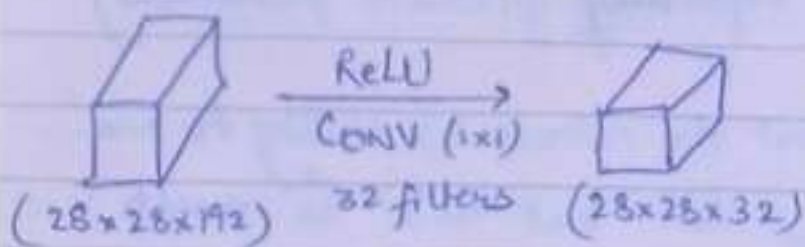$a^{[l+2]} = g\left(a^{[l]}\right) = a^{[l]}$  { as Relu of +ve is no. itself }.

**Note** Hence, it is easy for residual block to learn identity f^n, giving $a^{[l+2]} = a^{[l]}$. Thus adding the residual block in middle or in the end doesn't hurt the NN performance. These skip connec$^n$ bcoz of being not in plain network, the result get worse ↑ deepening of layer

The residual network doesn't hurt performance, and can sometimes actually help performance.

Adding $z^{[l+2]}$ and $a^{[l]}$ $\Rightarrow$ they have same dimension i.e Residual network uses same convolution resulting in same dimension matrices.
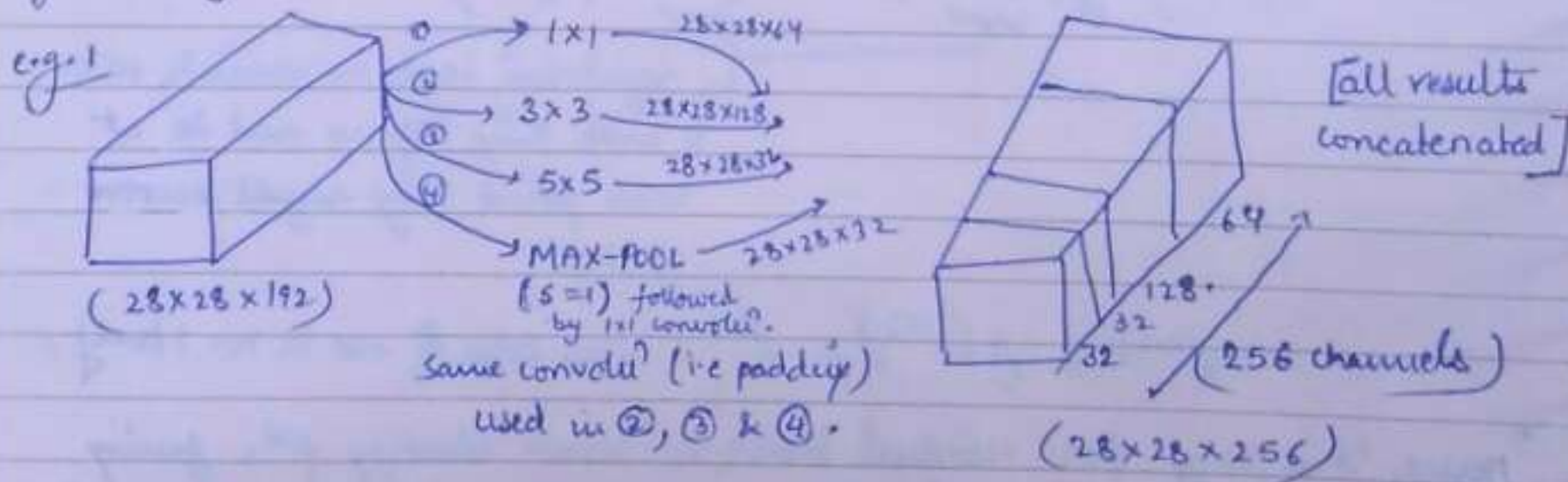
$\Rightarrow$ **1×1 Convolutions:** (also called Network in Network)



$(28 \times 28 \times 192)$   $\xrightarrow[\text{CONV} (1 \times 1)]{\text{ReLU}}$   32 filters   $(28 \times 28 \times 32)$

{ helps reduce no. of channels, from 192 to 32 like in this case unlike maxpool which helps reduce only $n_H$ and $n_w$. }

1×1 convolu' helps by adding non-linearity i.e taking Relu of the convoluted result, it helps ↑, ↓ and even keep same the $n_c$. This is useful in building the inception network.
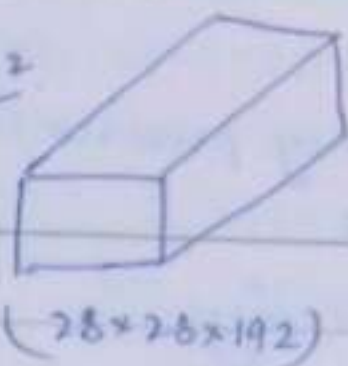
⑤ **Inception network :** It says why to chose if we want to use a conv layer or a pool layer and ē what parameters. It says lets try them all and let the network decide.
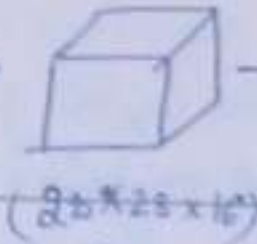
e.g. 1



$(28 \times 28 \times 192)$

$\textcircled{0} \rightarrow 1 \times 1 \rightarrow 28 \times 28 \times 64$
$\textcircled{2} \rightarrow 3 \times 3 \rightarrow 28 \times 28 \times 128$
$\textcircled{3} \rightarrow 5 \times 5 \rightarrow 28 \times 28 \times 32$
$\textcircled{4} \rightarrow$ MAX-POOL $\rightarrow 28 \times 28 \times 32$
$(s=1)$ followed by 1×1 convolu'.

Same convolu' (i.e padding) used in ②, ③ & ④.

[all results concatenated]

64
128
32
32
(256 channels)

$(28 \times 28 \times 256)$

Incep' network problem is computation cost.

**Note** This computa' cost can be reduced by using 1×1 convolu' layer.

e.g. 2



CONV (1×1) 16 filters

CONV (5×5) 32 filters

(28×28×192)

(28×28×16)

(28×28×32)

↑ Shrunken intermediate form & less ne wry 1×1 convolu".

↓ Same o/p obtained as in previous case but with less no. of computa"

Computa" cost for ③ in e.g. 1 = $(5 \times 5 \times 192) \times (28 \times 28 \times 32) = 120$ million.

" " for e.g. 2 = $[(1 \times 1 \times 192) \times (28 \times 28 \times 16)] + [(5 \times 5 \times 16) \times (28 \times 28 \times 32)]$

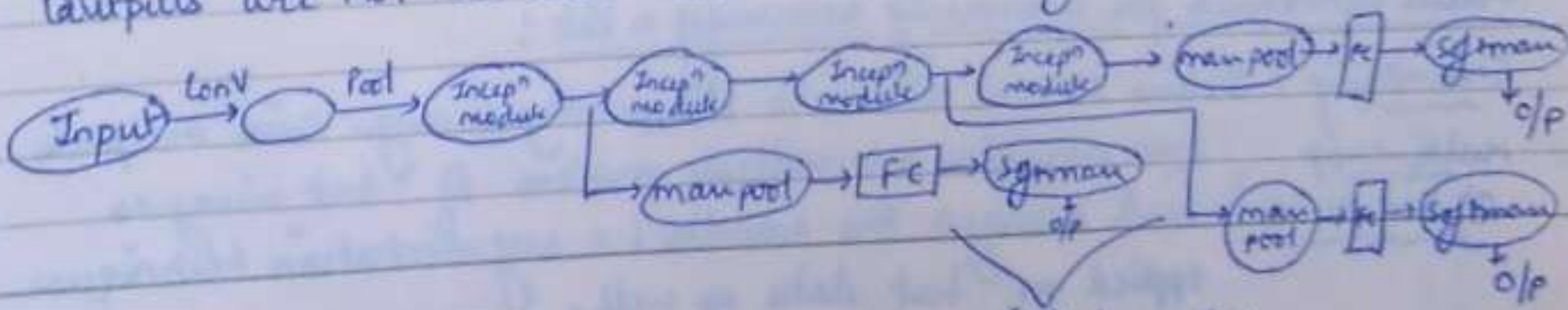= 12.4 million = $\frac{1}{10}$ th of above cost.



↳ 'bottleneck layer' (intermediate layer obtained after 1×1 convolu")

↓

This shrinking down doesn't hurt the performance much but reduces computa" cost significantly.

e.g. 1 is an inception module. Too many inception modules used together in a network give rise to inception network.

Note Side branches are added in between & the inception modules to show and confirm that if output is calculated from the intermediate modules using hidden layers + not upto it only, outputs are not too bad. This helps in regularize" as well.
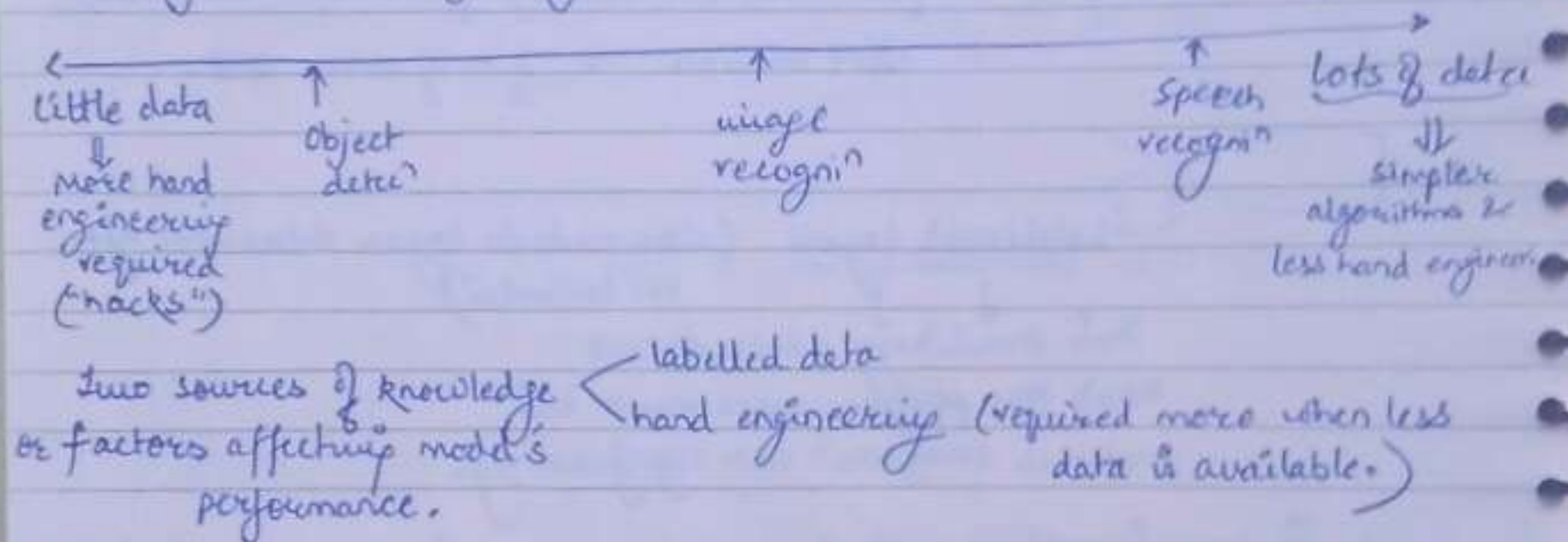


Side branches & o/p similar to the final o/p

⇒ **Data Augmentation:** to ↑ the training dataset present ∈ us. Various augmentaⁿ techniques given that they preserve the o/p as in the original image are:

make the learning also more robust to changes in the original images.

i) Mirroring      ii) Random cropping
iii) Rotaⁿ      iv) Shearing
v) Local warping      vi) Color shifting

* Loading data and getting trained can happen simultaneously.

←——————————————————————————————————→

| little data | ↑ | ↑ | ↑ | speech recognⁿ | lots of data |

little data ↑ Object detecⁿ     image recognⁿ     speech recognⁿ    lots of data
↓                                                  ↕
more hand engineering required ("hacks")                                       simpler algorithms & less hand engineering

Two sources of knowledge or factors affecting model's performance.
⟨ labelled data
⟨ hand engineering (required more when less data is available.)

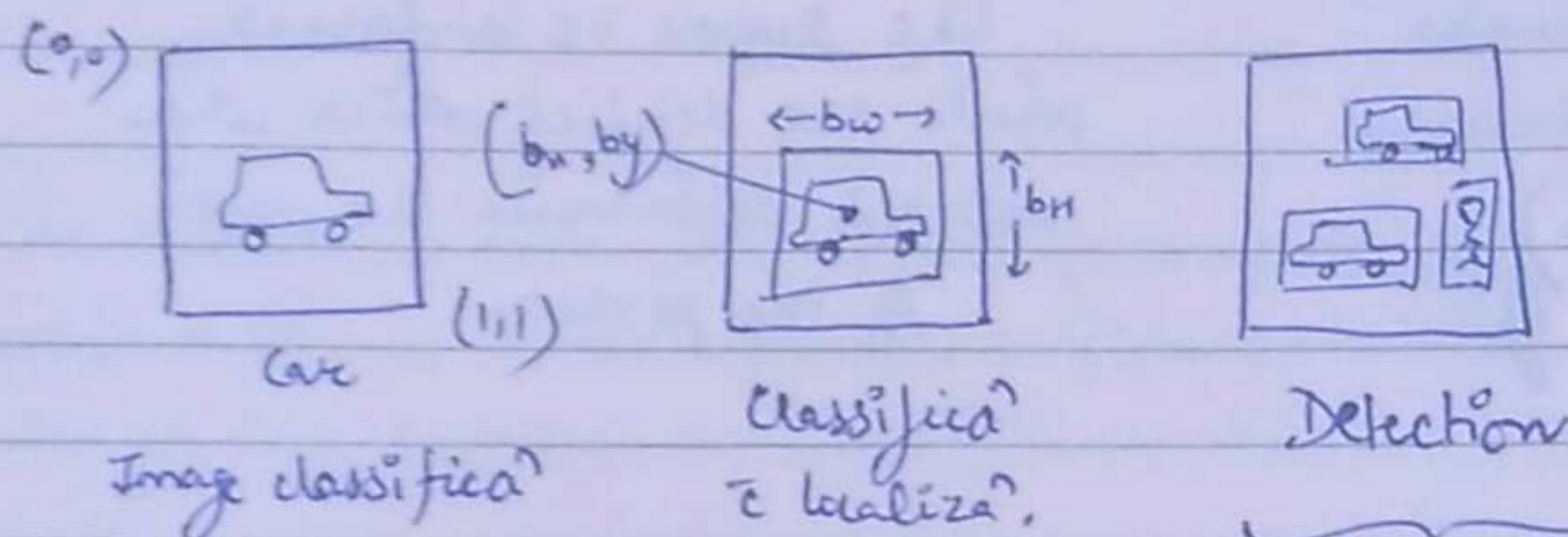**Note** Transfer learning helps in case of less availability of data.
↳ using other people's trained models and training just the last few layers of their network with our data along with changing the o/p of the softmax layer.

**Certain techniques for improving accuracy a bit:**

① Ensembling : Train several networks independly & avg. their outputs.
② Multi-crop at test time : Run classifier on multiple versions of test images and average the results i.e augmentation techniques applied on test data as well.
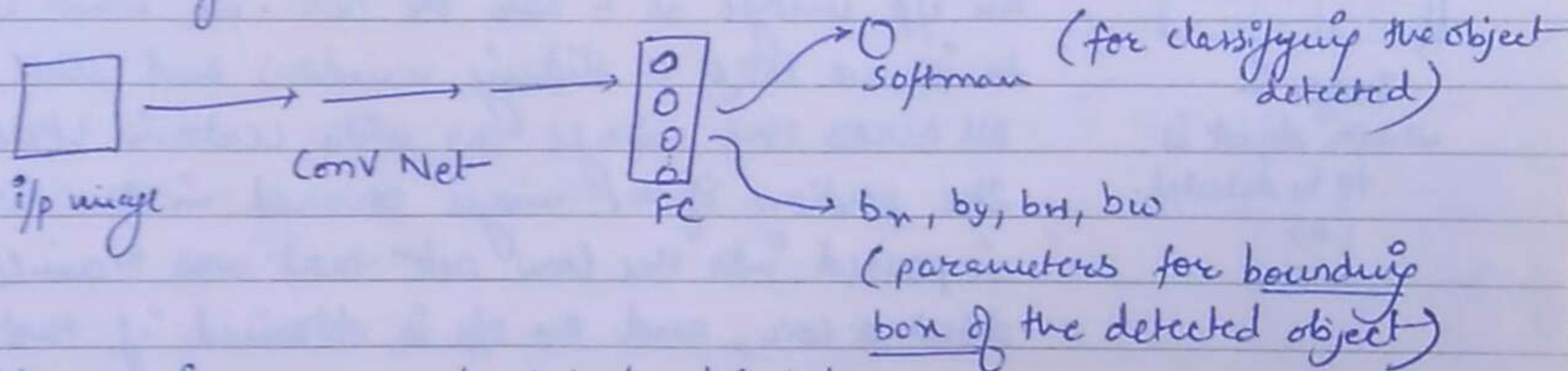
these techniques improve accuracy a little bit but use a lot of memory especially ensembling as it uses several networks. Hence, these techniques are used for winning competitions but not for producⁿ (deployment) purposes.

⇒ <u>Localization and Detection :</u>



Image classifica⁰          Classifica⁰          Detection
                          & localiza⁰.

usually have 1 big object       there can be
in the middle that needs to     multiple objects
be recognized & localized.



i/p image        Conv Net        FC

→ O Softman (for classifying the object detected)

→ $b_x, b_y, b_H, b_w$
(parameters for bounding box of the detected object)

eg. for an image $\tau$ one object to be detected :

$$y \text{ or output label} = \begin{bmatrix} P_c \\ b_x \\ b_y \\ b_w \\ b_H \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

$P_c$ → Probability that there is an object. (0 or 1)

$b_x, b_y, b_w, b_H$ → bounding box parameters if $P_c = 1$

$c_1, c_2, c_3$ → probability of belonging to a class of object if $P_c = 1$, provided object can be either a car, bike or cycle.

if $P_c = 0$, then rest i.e $b_x, b_y, b_w, b_H, c_1, c_2, c_3$ are dont cares (?).

**Landmark detecⁿ** {helpful in detecting face emotions}.

(i.e detecⁿ of important points in an image, now the o/p layer will have the coordinates of these important landmarks in the image as well) along c the o/p unit telling if its a face or not. ✓

or in snapchat, etc, or in pose detection as well.
↓
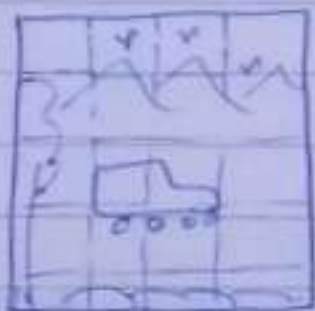like suppose 32 landmark points are defined which when detected determine the pose of the person.



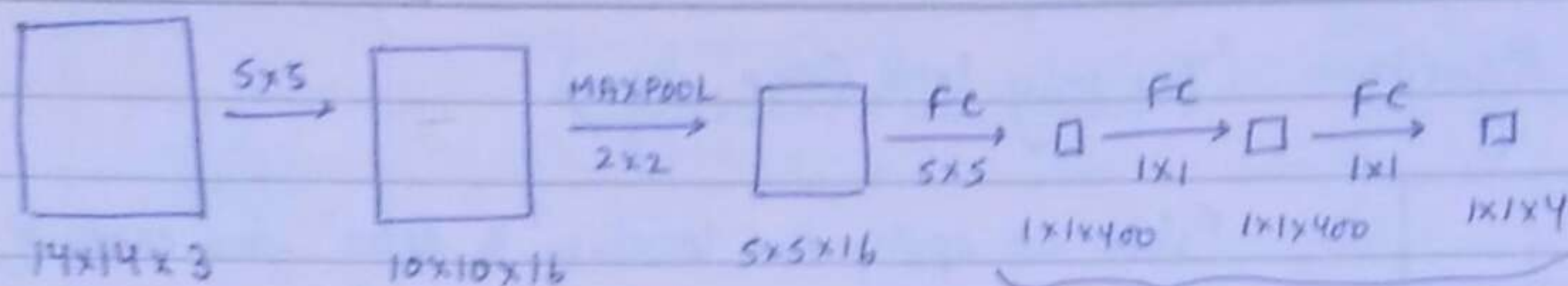Image in which object is to be detected.
(a)

**Object detecⁿ** using sliding window detection.

A conv net is trained suppose for detecting if an i/p image is a car or not. So, what we do is we take a sliding window and slide it all across the image (a) with certain stride. The portion of the image covered in the window is passed into the conv net that was trained to detect a car, and an o/p is obtained if that portion of the image contained a car or not. This process is repeated with diff window size in a hope that there will be some window that will bound the car tnt in (a) and output 1 i.e car tnt will be obtained in the conv net on passing that window as i/p to it.
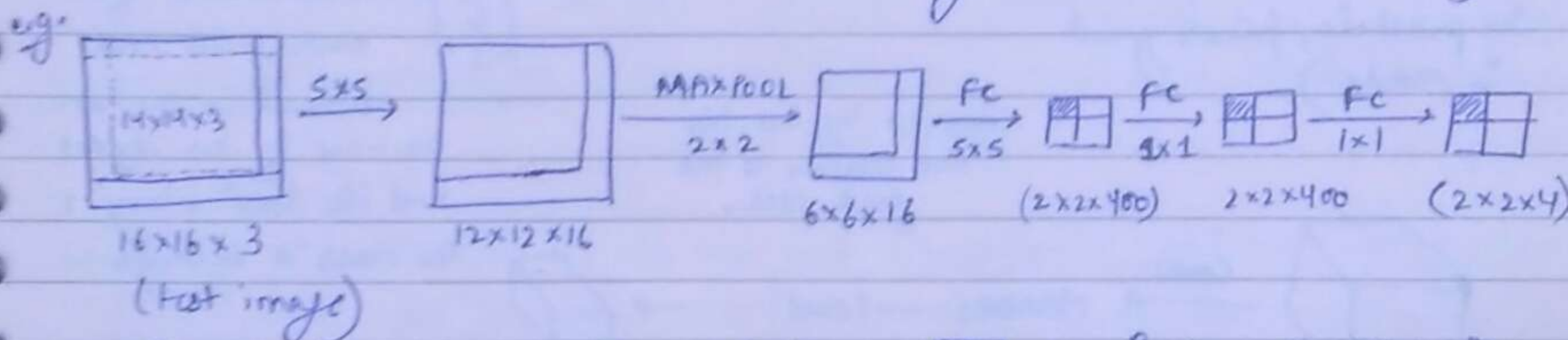
**Drawback** of sliding window detecⁿ : computation cost bcoz if fine strides are not used, we end up not being able to localize the objects that accurately within the image.

→ To reduce the computaⁿ cost of sliding window detecⁿ, **convolutional implementation** of sliding windows is done.
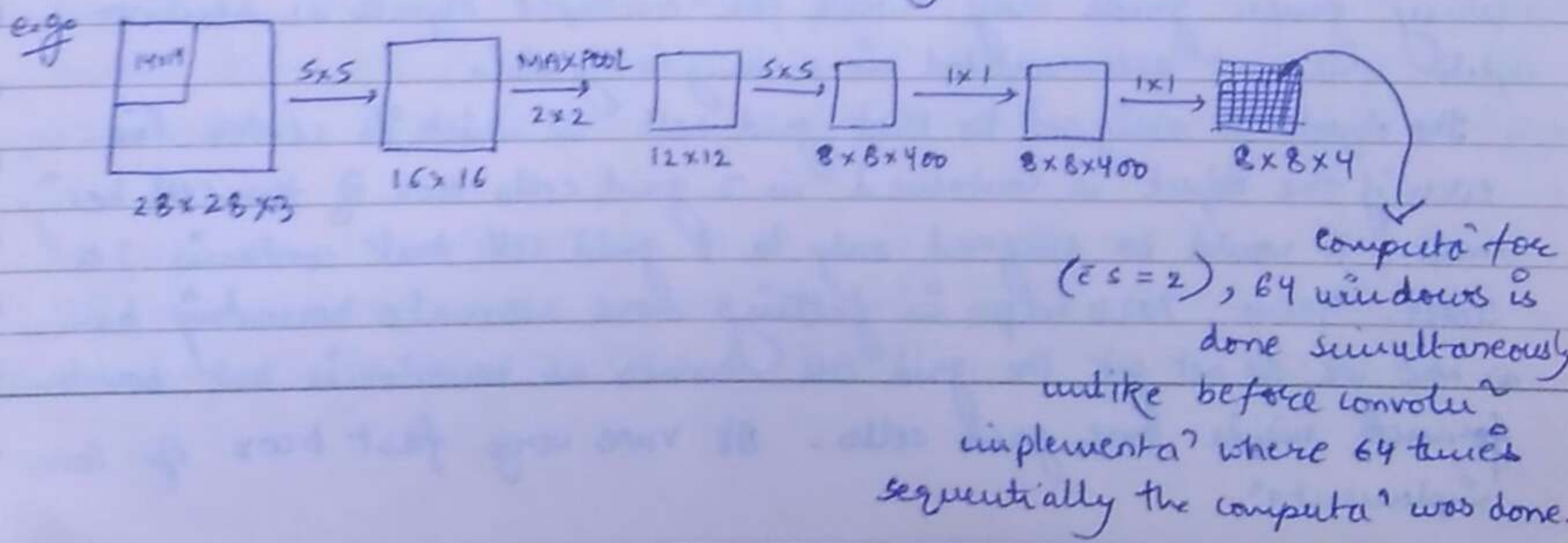
$14\times14\times3$ $\xrightarrow{5\times5}$ $10\times10\times16$ $\xrightarrow[2\times2]{MAXPOOL}$ $5\times5\times16$ $\xrightarrow[5\times5]{FC}$ $1\times1\times400$ $\xrightarrow[1\times1]{FC}$ $1\times1\times400$ $\xrightarrow[1\times1]{FC}$ $1\times1\times4$

↳ Convolutional implement of fully connected layers

In sliding window detecⁿ, lot of computaⁿ is repetetive, so convolutional implementaⁿ allows us to share the common computaⁿ thus reducing the computaⁿ cost. e.g.
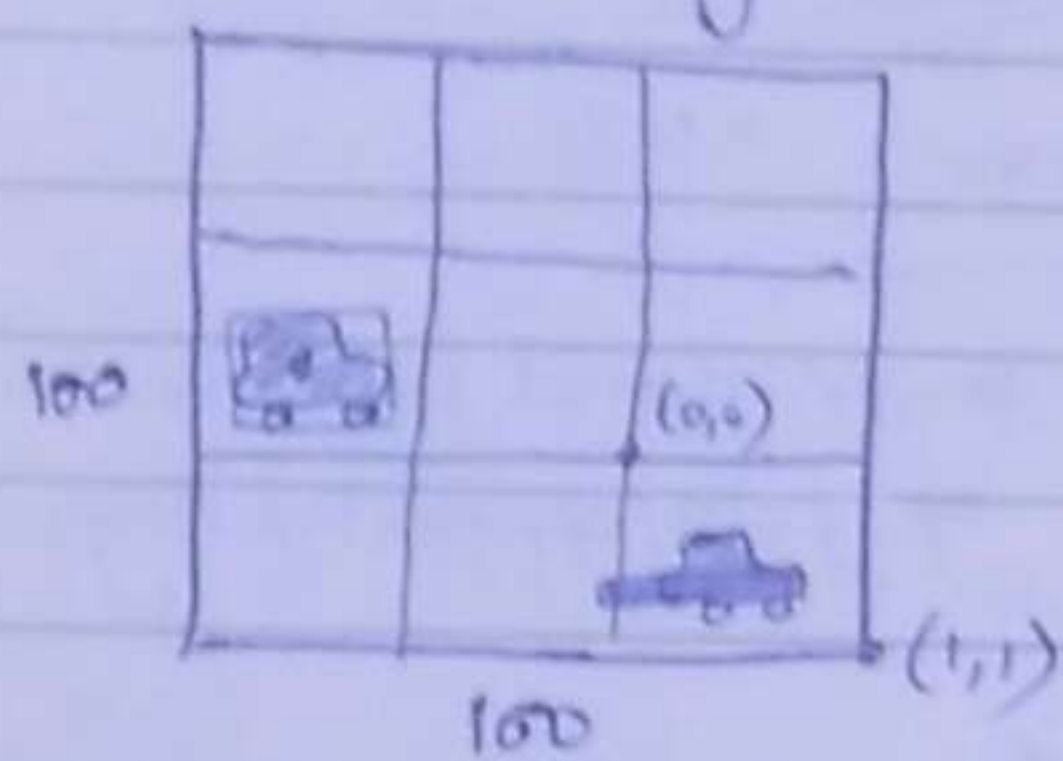
e.g.
$16\times16\times3$ $\xrightarrow{5\times5}$ $12\times12\times16$ $\xrightarrow[2\times2]{MAXPOOL}$ $6\times6\times16$ $\xrightarrow[5\times5]{FC}$ $(2\times2\times400)$ $\xrightarrow[1\times1]{FC}$ $2\times2\times400$ $\xrightarrow[1\times1]{FC}$ $(2\times2\times4)$

(Test image)

So, in the above process, the last layer gives the classificaⁿ result for 4 filters simultaneously which have a stride of 2. This save 4 times windows computaⁿ cost as the common computaⁿ of the 4 windows is done simultaneously.

e.g.
$28\times28\times3$ $\xrightarrow{5\times5}$ $16\times16$ $\xrightarrow[2\times2]{MAXPOOL}$ $12\times12$ $\xrightarrow{5\times5}$ $8\times8\times400$ $\xrightarrow{1\times1}$ $8\times8\times400$ $\xrightarrow{1\times1}$ $8\times8\times4$

↳ computaⁿ for ($\bar{c}s = 2$), 64 windows is done simultaneously unlike before convoluⁿ implementaⁿ where 64 times sequentially the computaⁿ was done

Drawback of this implementa? is that we don't get exact boundary boxes, approx. boundaries are obtained bcoz of using a fixed stride in the windows.
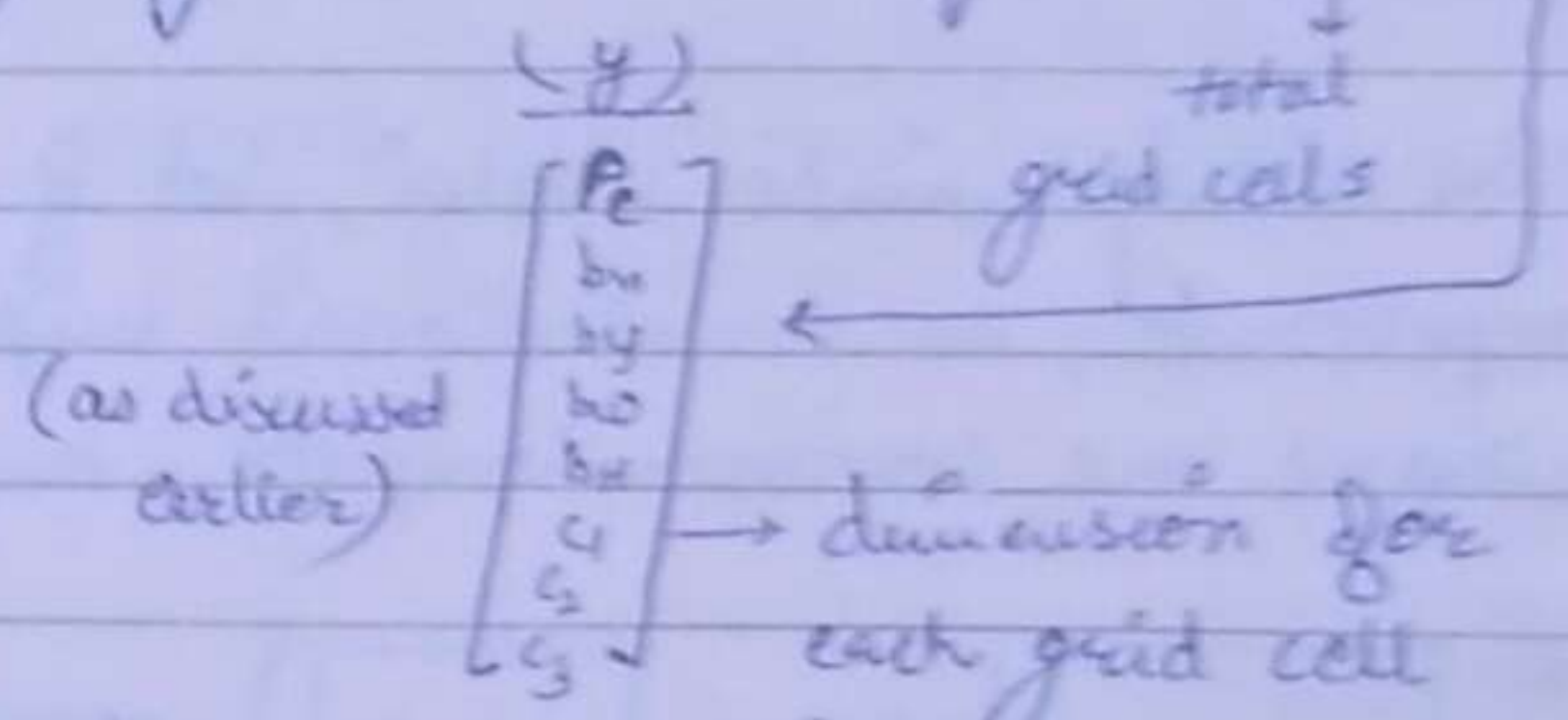
Solu?: YOLO Algorithm (You only look Once Algo)

First we apply the localiza? & detec? for each of the grid cell giving an output y i.e $3 \times 3 \times 8$.

(y) → total grid cells

$$\begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

(as discussed earlier)

→ dimension for each grid cell indicating the presence or-nce of an object and its loca? along & the class it belongs to.

(i/p image)

100

(0,0)

(1,1)

100

(In practice, finer grid is made.)

Note $b_x, b_y, b_w, b_h$ are relative to the grid cell.



CONV → MAXPOOL → CONV ····→

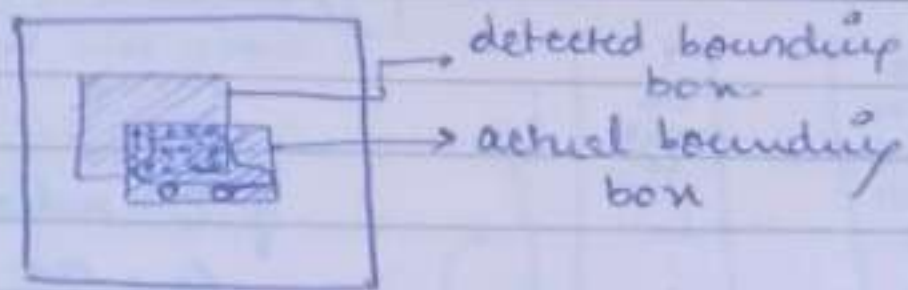$100 \times 100 \times 3$

(i/p image)

$(3 \times 3 \times 8)$

output (y)

This algo works fine as long as each grid has just one object in it. Making finer grids may work for multiple objects as now no 2 objects will get accomodated in a single grid.

Note: The object is assigned to that grid cell in which its center lies. So, even if the object is contained in 2 grid cells bcoz of the cell being small; it would be assigned only to 1 grid cell that contains its center. Hence, YOLO helps in getting more accurate boundary boxes as now we do not get the grid cell boxes as boundaries but boundaries formed inside these grid cells. It runs very fast bcoz of conv implementa?.

→ <u>Intersec over union</u> {to evaluate the object detecⁿ algorithm}
  (IoU function)



3/p image.

detected bounding box

→ actual bounding box

$$IoU = \frac{\text{Size of } \boxed{\cdot} \text{ (intersec area)}}{\text{Size of } \boxed{/\!/\!/} \text{ (union area)}}$$

3/ $IoU \geqslant 0.5$ (o/p is correct)

More the IoU, better is the bounding box predicⁿ, i·e the object is correctly localized if $IoU = 1$ and approx. correct if $IoU \geqslant 0.5$.

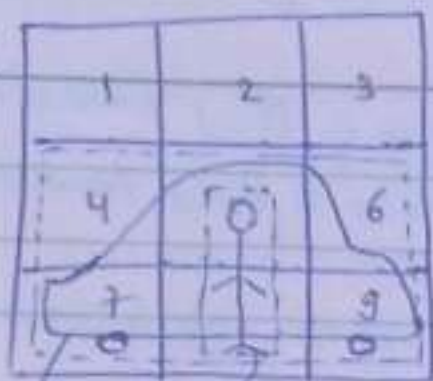→ <u>Non-man suppression</u> {to make sure that our algo detects each object only once}



Though the car has just 1 midpoint, but we run the algo for each grid cell and bcoz every cell will consider its object centre as the centre of the object, car will be detected multiple times.

Non man Suppression removes the extra detections of an object using the probability Pc. The output ŷ with highest Pc i·e the man probability that it contains an object is kept and the rest of the boxes with low Pc or higher IoU with the box having higher Pc are removed. Hence the name, i·e Suppressing the ones with non-man Pc.

To apply non-man suppression on an image with multiple class objects, the non-man suppression is run as many times as the no· of classes, i·e individually for each class of object.

→ Anchor boxes {to enable a grid cell to detect multiple objects}



Lets define 2 anchor boxes $\boxed{1}$, $\boxed{2}$. Even more can be defined.

$$ y = \begin{bmatrix} P_c \\ b_n \\ b_w \\ b_H \\ c_1 \\ c_2 \\ c_3 \\ P_c \\ b_n \\ b_y \\ b_H \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} 1 \\ b_n \\ b_y \\ b_H \\ b_w \\ 1 \\ 0 \\ 0 \\ 1 \\ b_n \\ b_y \\ b_H \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix} $$

→ corresponds to anchor box 1.

→ corresponds to anchor box 2

if $c_1$ = pedestrian
$c_2$ = car

↘ car  ↘ pedestrian

(mid point of both of them lie in the same grid cell)

we see that bounding box of pedestrian resembles anchor box 1 so we assign the pedestrian, $\boxed{1}$ and the bounding box of car is similarly assigned $\boxed{2}$.

(for grid cells)

So, the output comes out as $(3 \times 3 \times 16)$ for anchor box algorithm, as we have $3\times 3$ grid cells and 2 anchor boxes and 8 outputs corresponding to each anchor box

where $8 = 5 + (\# \text{ classes})$

$$ y = \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ 1 \\ b_n \\ b_y \\ b_H \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix} $$

if the grid cell has only a car and no pedestrian and the car has bounding box similar to anchor 2.

Drawbacks: i) three objects in a grid cell but only 2 anchor boxes available.
ii) both objects resemble the same anchor box.

the above 2 condiⁿ can't be handled by the above method of anchor boxes.

Note: Anchor box shapes are defined in such a way that a particular object if present in any shape can be detected by the anchor box. This helps algorithm specialize in classifying different object on basis of their shapes and sizes.

RCNN: Regions 2 CNN that is the regions that make sense to run CNN on. Only those regions in an image which have some interesting object to be detected are selected to run CNN on or window on.

Segmenta^n algorithm helps finding such regions (blobs) and run CNN on just those blobs, thus reducing the computa^n cost significantly.

⟹ Face Recognition:

Face verifica^n: Input image, name/ID
         Output whether the i/p image is that of the claimed person.

Face recognition: Input image and identify the name of the person if it exists in the database.

One shot learning: learning from just one example to recognize the person again, this learning is generally required in face recogni^n systems where we have just 1 picture of our employee available to train the network. This is in contrast c̄ our DNN which we have studied so far and which needs lot of training data to learn the features.

So, for one shot learning, we use a `similarity function`. Neural network learns the f^n ⟍

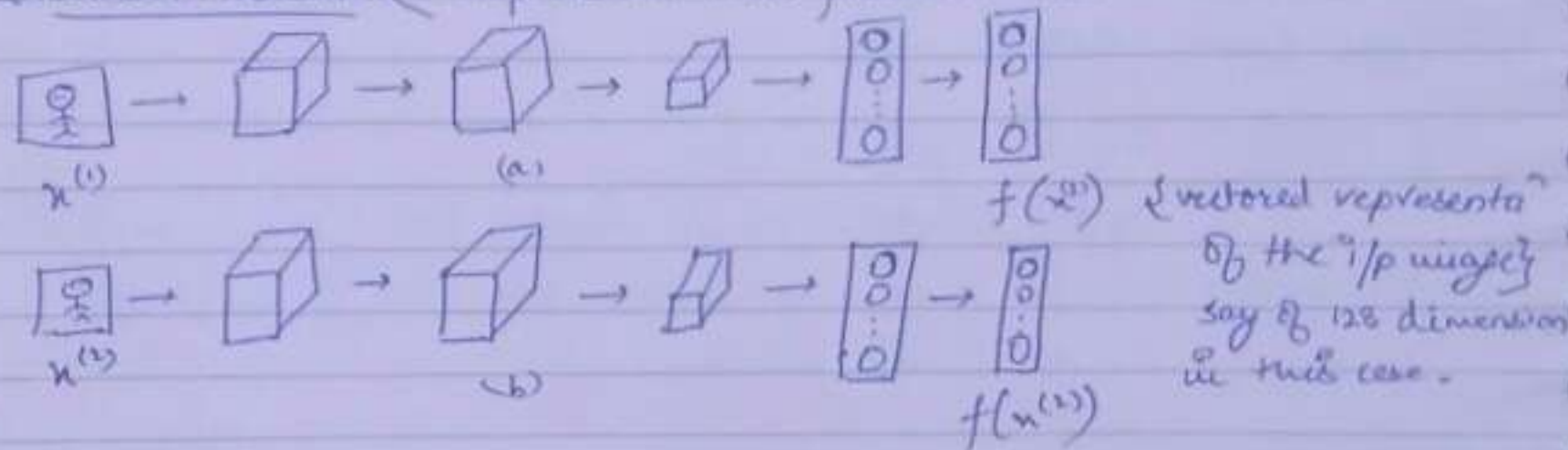$$d(img\ 1, img\ 2) = \text{degree of difference b/w images.}$$

If $d(img1, img2) \leq \tau$    "Same person"     } Image
$\phantom{If d(img1, img2)} > \tau$    "different person"    } verification.

$\tau$ is a hyperparameter.

for face recogn$^n$ the same formula is applied with every image present in the database. If the i/p image matches some image in the database $d(img1, img2)$ for it is quite low.

So, if any new person joins the team, simply a new image can be added to the database and 'd' remains the same.

Siamese Network help learn this function 'd'.



$x^{(1)}$        (a)        $f(x^{(1)})$ {vectored representa$^n$ of the "i/p image"}

$x^{(2)}$        (b)        $f(x^{(2)})$    say of 128 dimension in this case.

Now, $d(x^{(1)}, x^{(2)}) = ||f(x^{(1)}) - f(x^{(2)})||^2_2$

* Network (a) and (b) are the same with same parameters learnt such that if $x^{(i)}$ and $x^{(j)}$ are the same person, then $d(x^{(i)}, x^{(j)})$ is small, and " " " " " different " " " " " " large.

→ <u>Triplet loss function</u>: this aims at tuning the network parameters such that encoding of the images of same person is close enough and encoding of images of different persons is farther enough. As, in this we deal with 3 images, it is called so.

| A | | P | | | A | | N | |

Image of　Another image　Image of　Another
person A　of person A　person A　image but
　　　　　　　　　　　　　　　　not of person A.

Now, triple fN wants : $d(A,P) \leq d(A,N)$

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

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

$\downarrow$ margin to avoid
NN from giving trivial
solⁿ i.e encoding all
images as zero vector or
as same vectors.

Note this $\alpha$ is a hyperparameter
which actually pushes AP pair
and AN pair further away from
each other.

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

(loss on a single
triplet)

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

overall cost
fN for NN

{where m is the no. of triplets
for training purpose.}.

If we have 10k images of
1k people, then triplets
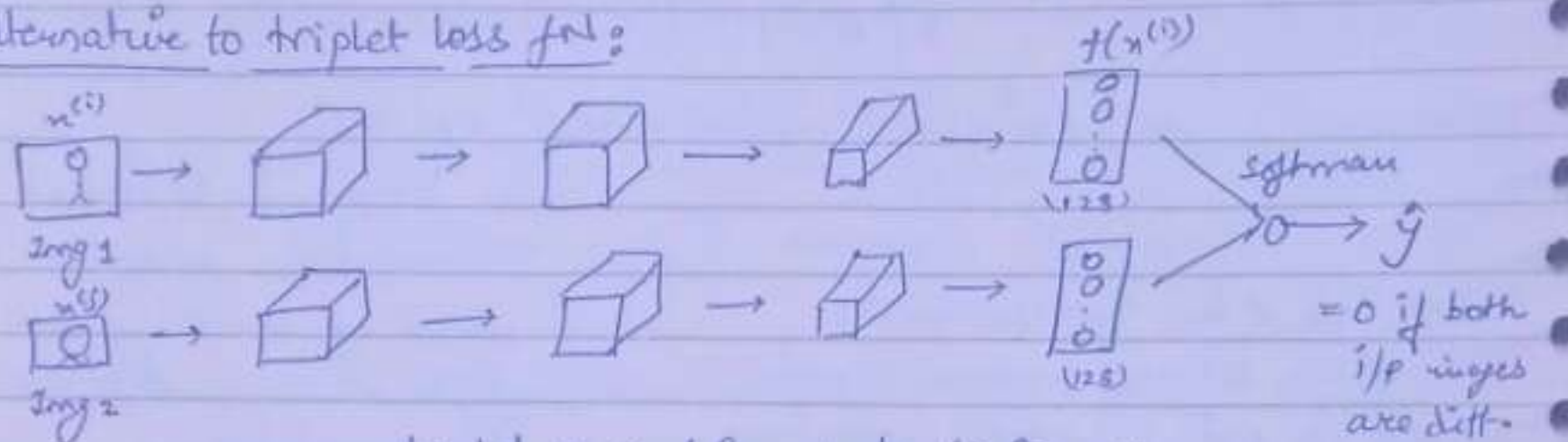APN need to be formed for train
dataset.

Note ↓
Hence, a good train data is needed to train the NN for one shot
learning. Once the hyperparameters are properly tuned to encode
an image, one shot learning is possible.

Triplets that are 'hard' to train on are chosen. Bcoz if triplets
are chosen randomly, there is higher probability of $L(A,P,N)$ to
be satisfied. Such triplets are chosen where $d(A,P) \approx d(A,N)$ so

that the network tries hard to learn the hyperparameters so as to meet the constraints i.e $d(A,P) + \alpha \leq d(A,N)$. In case of hard triplets only the gradient descent would learn something b'coz otherwise in case of randoms triplets, the network would get it right everytime and hence not learn anything.

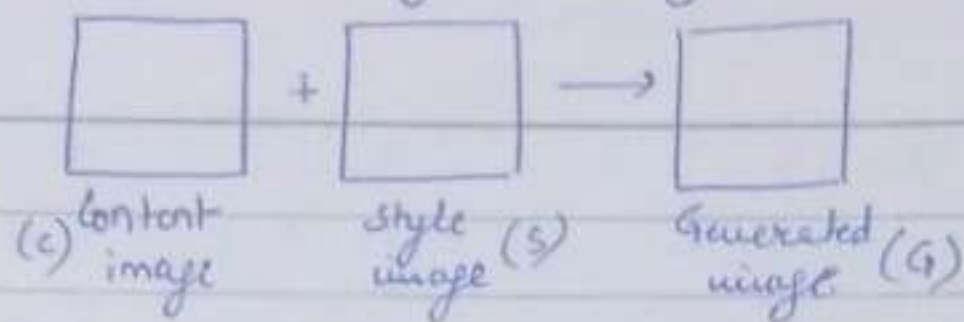## Alternative to triplet loss fn:



treated as a binary classifier problem

The 128 features are fed into the logistic regression unit i.e softmax layer

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

↓
there can be other variations to this to compute the difference b/w the two f's.

→ parameters as used in any logistic regression.

$= 0$ if both i/p images are diff.

$= 1$ if i/p images are same.

Note ↓

So, in this approach the training set is not a triple but a pair of images where the target label is 1 if the pair has images of same person and 0 if the pair has images of diff person.

$\Rightarrow$ Neural Style Transfer:

$$\square + \square \longrightarrow \square$$

(c) Content image    style image (s)    Generated image (G)

Going deeper in NN, more and more complex features are detected by the hidden units. Starting ε lines and edges, going to detecting shapes and so on.

- for Generating a Neural style transfer image, a cost fN is defined:

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

cost fN in generating G.     measures how similar is the content of generated image to the content of the content image C.     measures how similar is the style of the image G to the style of the image S.

steps involved:

i) Initialize G randomly.

ii) Use Gradient decent to minimize $J(G)$.

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

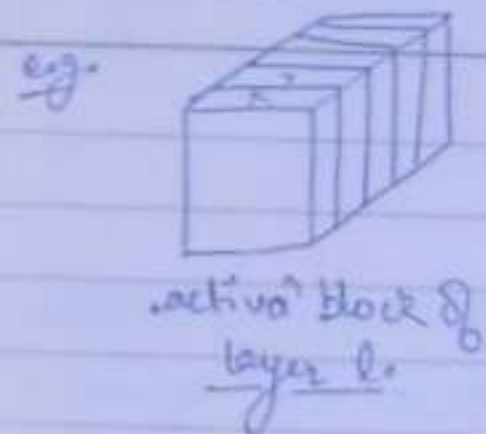Note $\boxed{J_{content}(C,G) = \| a^{[l](C)} - a^{[l](G)} \|^2}$

    activa of $l^{th}$ hidden layer for image C.    activa of $l^{th}$ hidden layer for image G.

{ $l$ is a middle layer not too shallow, not too deep in the conv Net. (which is pretrained to generate features of the i/p image) }

middle layer is taken so that neither too simple nor too complex features are extracted from the image.

If layer $l$'s activa? is used to measure "style", then style is defined as correlation between activations across channels.

e.g.



activa? block of layer $l$.

If correla? b/w R and Y channel is calculated, it means calculating the probability that whatever part of the image has features detected by Y also has features which are detected by channel R. So, it is basically finding out what features tend to occur together and what features are rarely seen together in different parts of an image.

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

(with labels H, W, C pointing to $i$, $j$, $k$)

↓ Style matrix of hidden layer $l$ having $n_c$ channels.

↓ to measure the correla? b/w each pair of channels.

$$G_{kk'}^{[l]} = \sum_{i=1}^{n_H^{[l]}} \sum_{j=1}^{n_w^{[l]}} a_{ijk}^{[l]} \, a_{ijk'}^{[l]}$$

where $k$ and $k' \in \{1, 2, \dots n_c\}$

calculate for every $k$ and $k'$ to find $G^{[l]}$ or the style matrix.

Note If the two channels $k$ and $k'$ are correlated $G_{kk'}$ will be large otherwise if " " " " uncorrelated " " " small.

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

↓ style matrix of image S.   ↓ style matrix of generated image G.

$$J_{style}^{[l]}(S,G) = \sum_{k} \sum_{k'} \left( G_{kk'}^{[l](S)} - G_{kk'}^{[l](G)} \right)^2$$
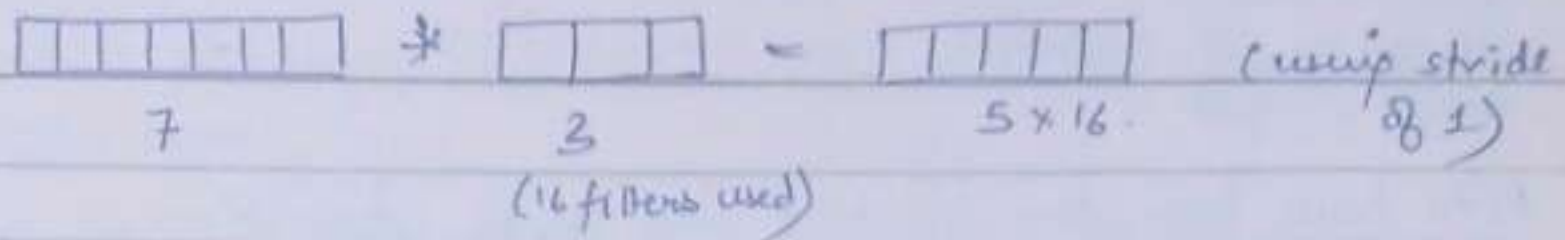
$$J_{style}(S,G) = \sum_{l} \lambda^{[l]} J_{style}^{[l]}(S,G)$$

↳ (hyperparameter)

{ taking into account both low level & high level feature while calculating correlations }

**Note** Convolu' opera' can be performed on 1D, 2D as well as 3D inputs.

e.g. of 1D.



7 $\quad$ 3 $\quad$ 5 × 16. $\quad$ (using stride of 1)

(16 filters used)

e.g. of 2D.

$$(10 \times 10) * (5 \times 10) = (6 \times 16)$$

$\downarrow$
16 filters used

e.g. of 3D. $\quad (10 \times 10 \times 16) * (5 \times 5 \times 16) = (6 \times 32)$

$\downarrow$
32 such filters used.

For **1D data**, recurrent NN are more used, though CNN can also be applied as seen above.

→ **3D convolu':**



3D volume
(e.g. CT scan)
say (14×14×14)

(e.g. movie data)

3D filter
(5×5×5)
(16 filters)

3D output
(10×10×10)×16

{ Data can also have different no. of channels as we saw in case of 2D. In this e.g. $n_c = 1$ }

**Note** Activation functions:

| sigmoid | tanh | Relu |
|---|---|---|
| $\frac{1}{1+e^{-x}}$ | $2\,\text{sigmoid}(2x) - 1$ | $\max(0, x)$ |
| range: $(-1, 1)$ | range: $(-1, 1)$ | |
| vanishing gradient decent | less vanishing GD | range: $(0, \infty)$ non-linear fn, hence no problem like linear fn. |

- No vanishing gradient problem
- leaky Relu over dying Relu
- Relu is sparse i-e not all activa' are processed to describe the o/p of network hence reducing cost, as 50% network yields 0 b'coz of property of Relu.

Disadv