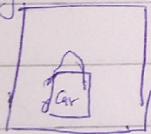


## Week3 Object detection Algorithms

Object localization

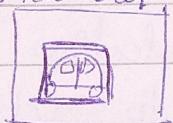
Img. classificat<sup>r</sup>



"car"

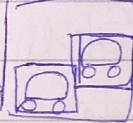
(To find "local" of object in the image)

Classificat<sup>r</sup> w/ localizat<sup>r</sup>



"car"

Detection



multiple objects

Here, the o/p of the softmax layer is a matrix:

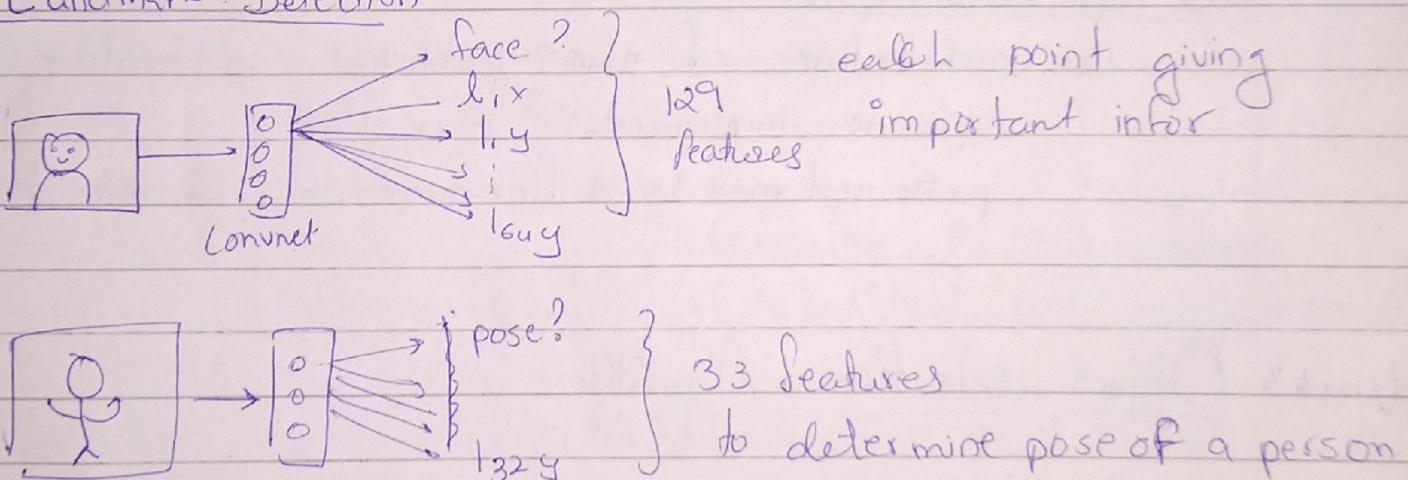
|       |                                |
|-------|--------------------------------|
| $P_c$ | → if object present or not     |
| $b_x$ | → midpoint of object           |
| $b_y$ |                                |
| $b_h$ | height width of the            |
| $b_w$ | object box                     |
| $C_1$ | → To the class it              |
| $C_2$ | belongs to.                    |
| $C_3$ | 1 → pedestrian, 2 → car, 3 → - |

### Loss function

$$\text{if } P_c = 1 \rightarrow \left\{ \begin{array}{l} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots \\ \dots + (\hat{y}_8 - y_8)^2 \end{array} \right.$$

$$\text{if } P_c = 0 \rightarrow ((\hat{y}_1 - y_1)^2) \quad \text{don't care of rest indices if obj. not present}$$

### Landmark Detection



### Sliding Window detection

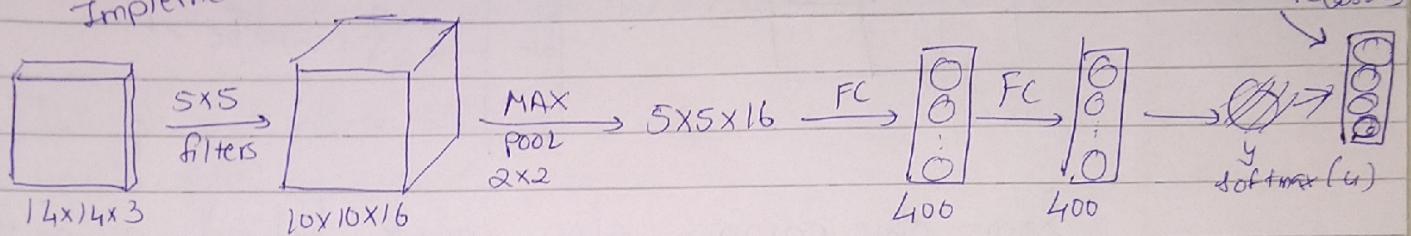
- 1) Train the model to determine if there is a car or not in the image.
- 2) Next, take a box, you slide it throughout the image, and show each image to model to find, if there is a car in the box or not?

3) Then you resize the window and slide it again

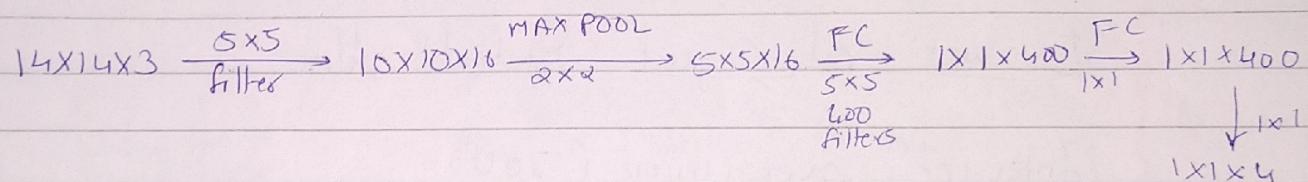
→ High Computat<sup>n</sup> cost

→ But if implemented Convolutionally, it can be faster.

Actual Implementation



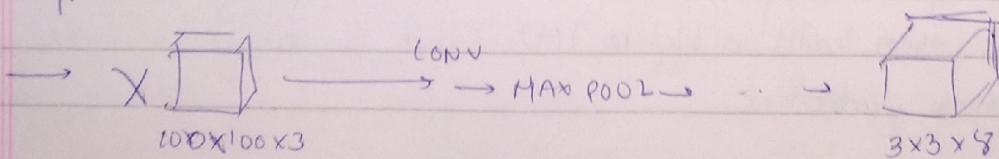
Convolutional Implementation



- Here, the o/p of all the ~~all~~ sliding windows is expressed at the best region. The reason this is faster, is because each sliding window shares some computat<sup>n</sup>. Also, instead of ~~giving~~ giving one window everytime to the model, it gives the whole image to the model.
- The problem is, the bounding box size is inaccurate.

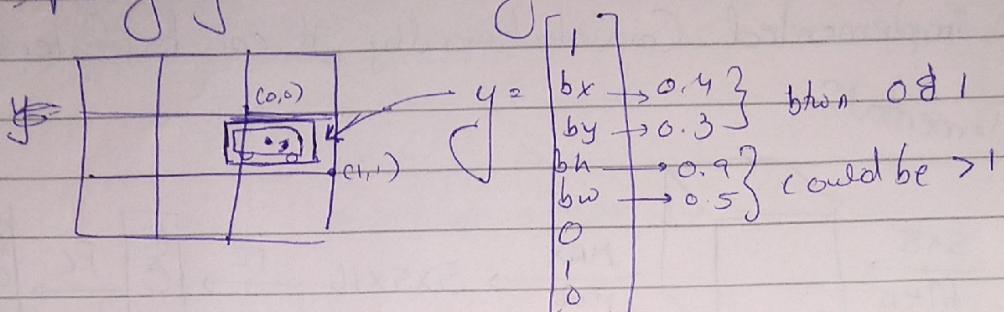
This can be solved by YOLO (You only Look Once)

The image is divided into  $3 \times 3$  grid. Each grid has a 8 element vector. So, when we put the image in above implementation, we get o/p as  $3 \times 3 \times 8$



here, ~~any~~ In reality, we use a finer grid  $19 \times 19$  instead of  $3 \times 3$ .

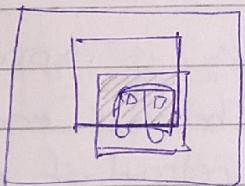
Specifying the bounding boxes in YOLO



Some other params also work behind giving the perfect box to the object

Evaluat<sup>n</sup> of Object Detect<sup>n</sup> Algos.

Intersect<sup>n</sup> over union (IOU)



$$\text{Size of Intersection} / \text{Size of Union}$$

"correct" if  $\text{IOU} \geq 0.5$

$\uparrow$  IOU  $\uparrow$  accuracy

To avoid your model to detect the same object again, you use "Non-Max Suppression"

(When you detect an object multiple times, it has a  $P_c$  (probability of detect<sup>n</sup>) associated with it.)

The one with highest  $P_c$  are chosen, getting rid of other estimat<sup>n</sup>s even with a high IOU. That is, suppress the non-max predictions.

Algorithm:

Input: Image of  $19 \times 19 \times 8$

1) Discard all boxes with  $P_c \leq 0.6$

2) ~~while~~ boxes left:

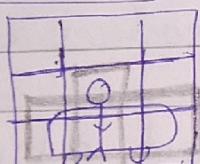
o/p the one with highest  $P_c$  as a predict<sup>n</sup>

Discard any remaining boxes w/  $J_{OV} \geq 0.5$  w/ the box o/p in previous step  
(discarding all other boxes of the same object)

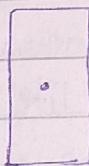
This process will be done for a single class, for multiple classes,  
you need to run it multiple times

When there are more than one object in a grid cell  $\rightarrow$  Anchor Boxes

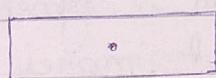
Anchor boxes



Anchor Box 1



Anchor Box 2



$y =$

$$\begin{bmatrix} P_c \\ b_x \\ b_y \\ b_w \\ b_h \\ c_1 \\ c_2 \\ c_3 \\ P_c \\ b_x \\ b_y \\ c_3 \end{bmatrix} \left\{ \begin{array}{l} \text{Anchor Box 1} \\ \text{Anchor Box 2} \end{array} \right.$$

We define two different sized anchor boxes to define each of the objects

Previously:

Each object in training image is assigned to grid cell that contains that object's midpoint

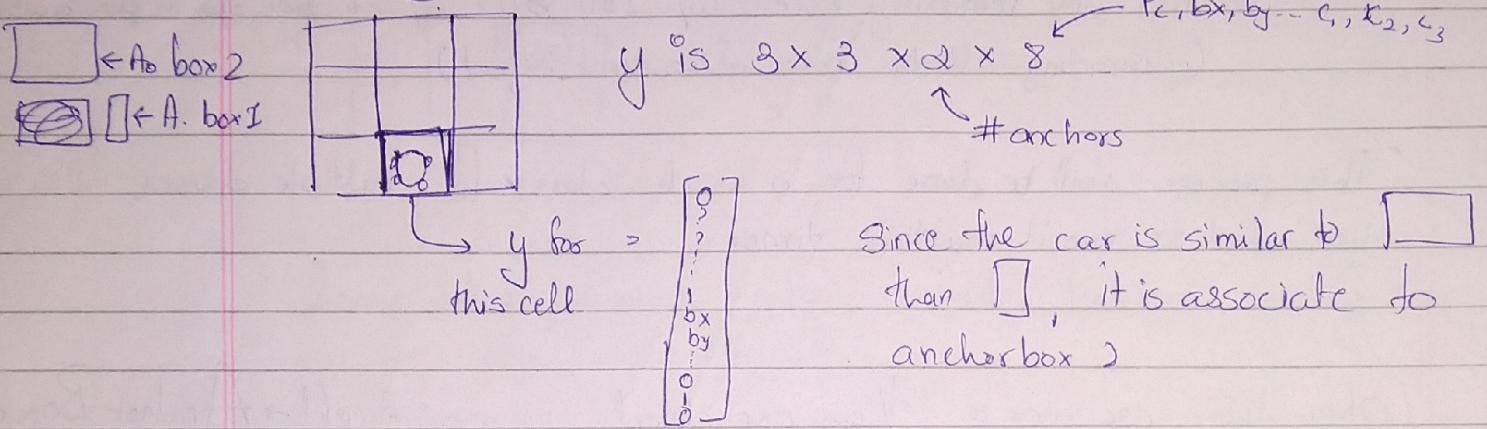
$y: 3 \times 3 \times 8$

With anchor boxes:

Each object in training image is assigned to grid cell that contains object's midpoint & anchor box for the grid cell with highest  $J_{OV}$   
o/p  $y: 3 \times 3 \times 16$

This overlapping generally doesn't happen where m.p. of two objects is on the same place. As, in real we generally use 19x19 grid & not 3x3

## Putting all back to YOLO



## Region Proposals - R-CNN

They run a segmentation algorithm. This produces something like Heat images (looks like it to me)

Here, the ones which have a different color or are highlighted are only fed for object detectn instead of the complete image.

Bounding boxes (square or rectangles) are drawn on those highlighted objects.

R-CNN outputs one region at a time and predicts the label. (and bounding box) This makes it slower

Fast R-CNN Uses convolution implementatn of sliding windows to classify all proposed regions. But, it was still slow

Faster R-CNN Use convolutional n/w to propose regions instead of segmentatn algo. Runs a bit faster than fast R-CNN algo.  
But this Faster R-CNN was a bit slower than YOLO Algorithm