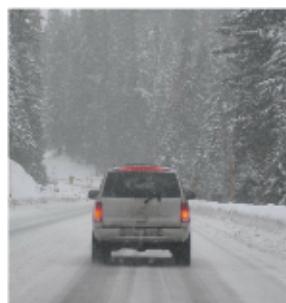


Object detection:-

Image localization vs detection:-

Image classification



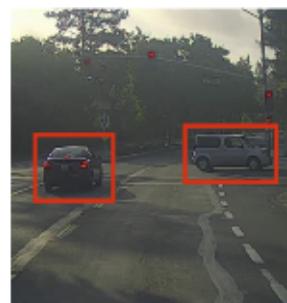
car

Classification with localization



car*

Detection

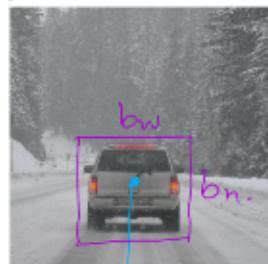


multiple object

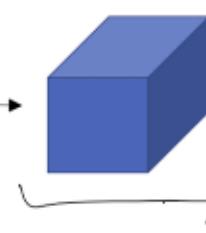
1 object

Classification with localization:-

(0,0)



(1,1)



conv net

softmax



...

b_x, b_y, b_w, b_h.

bounding box.

(b_x, b_y)

b_w → width of bounding box

coordinate · b_h → height of " "

Defining target label y:- (Supervised learning)

Classes:-

1 - pedestrian

need to output:-

b_x, b_y, b_w, b_h, Class label (1-4)

2 - Car

3 - motor bike

4 - background

$y =$

$\begin{bmatrix} P_e \\ bx \\ by \\ bw \\ bh \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$

if have object

Assuming having only one object per image:-

if a image have car then y will be no object present.

$$y = \begin{bmatrix} 1 \\ \square \\ \square \\ \square \\ 0 \\ 0 \end{bmatrix} \rightarrow \text{some value}$$

$$y = \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix} \rightarrow \text{don't care.}$$

loss :-

If there is an object in the image ($P_e = 1$) loss will be sum of square of all component.

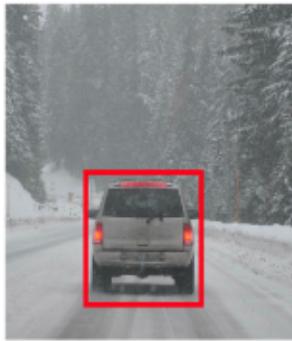
$$\mathcal{L}(\hat{y}, y) = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_8 - y_8)^2$$

if there is no object ($P_e = 0$)

$$\mathcal{L}(\hat{y}, y) = (\hat{y}_1 - y_1)^2$$

Instead of using the sum of square we can also use the log likelihood.

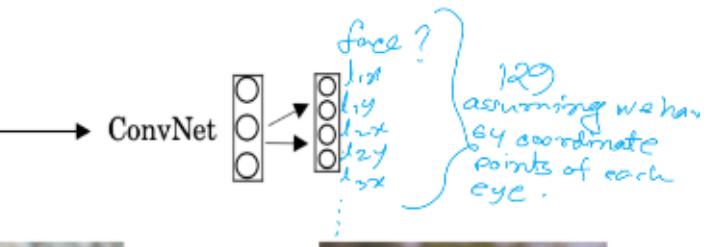
Landmark detection:-



b_x, b_y, b_h, b_w



$l_{1x}, l_{1y}, l_{2x}, l_{2y}$
 $l_{3x}, l_{3y}, l_{4x}, l_{4y}$



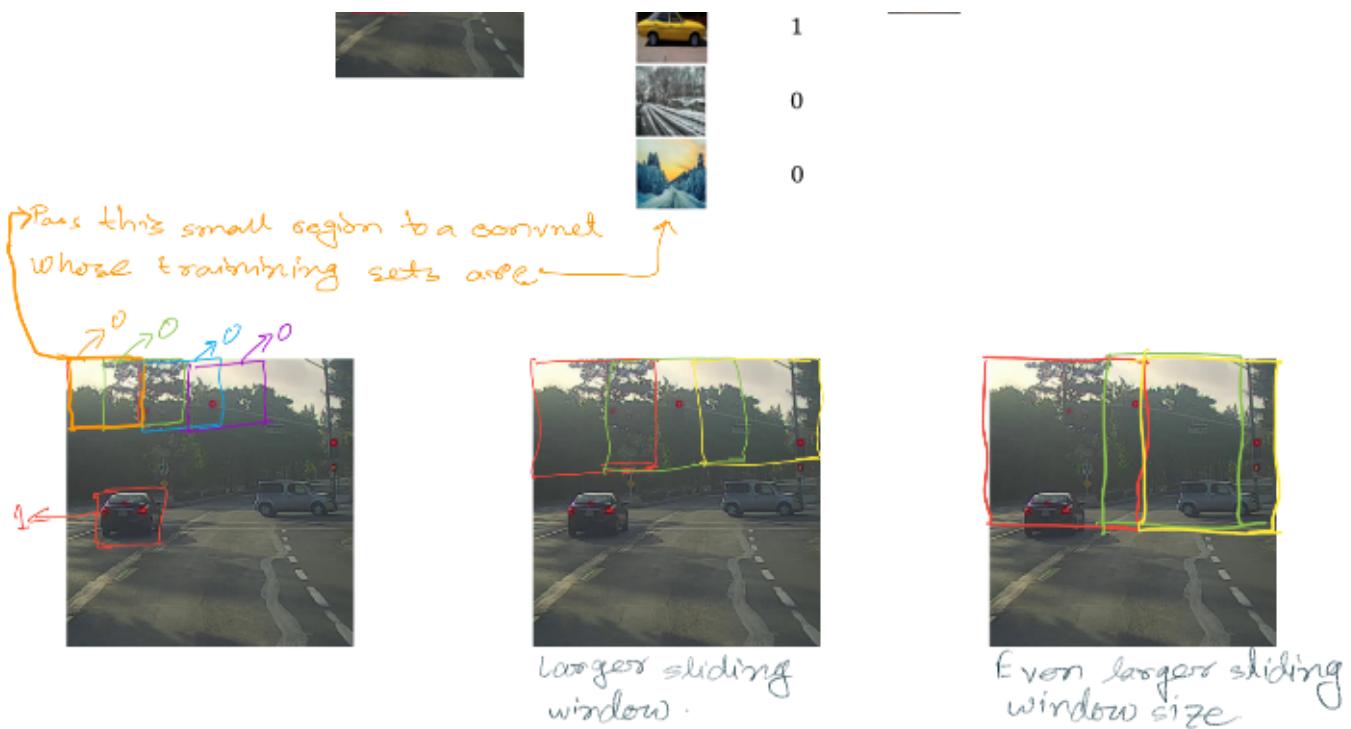
to find the shape of a body movement.

In order to find the eye of a person we need to mark the eye of person face image for training set. Let's we want to find the coordinate of two eyes of a person.

Object detection:-

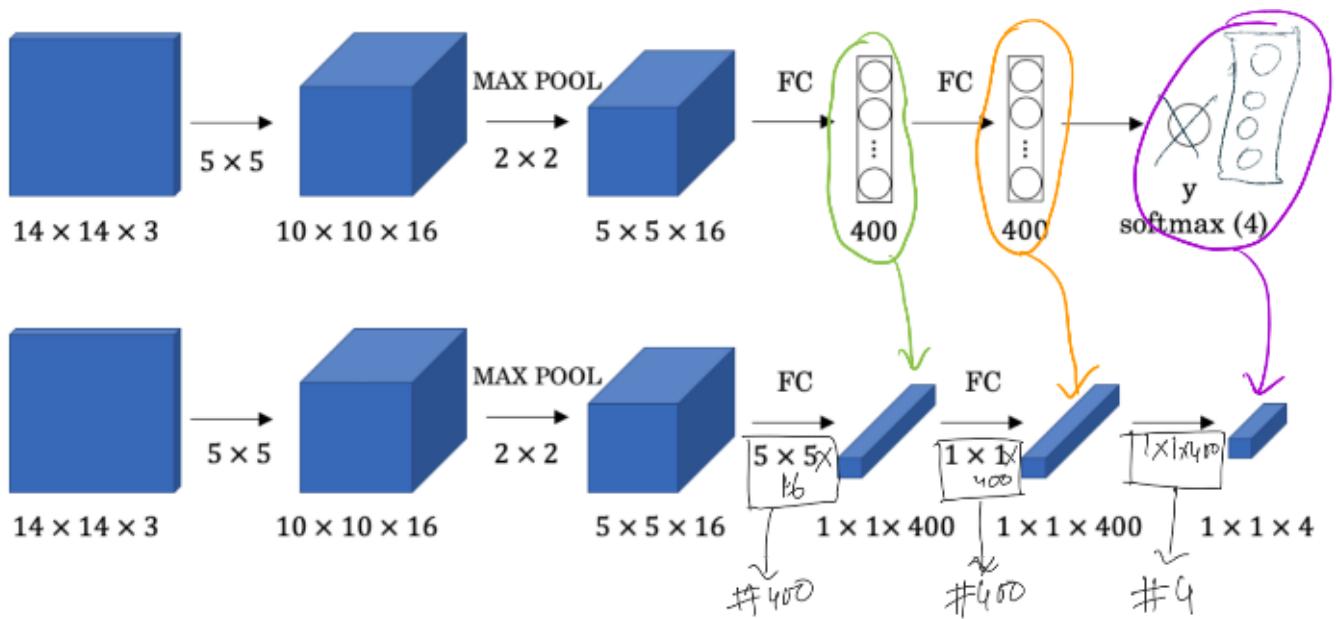
Let we want to localize cars in a image. One way to do that is sliding window detection. where we take a small portion of a image and pass it in a convnet. if there are any car in that small window then convnet will return yes. Training set for this model is the cropped version of car image.

Training set:		
x	y	
	1	
	1	



The computational cost of this model is much higher. We are using different sliding window size that makes it more costly. However we can reduce the cost by increasing the stride but the model will be inefficient in that case as we might miss some portion of the image.

Turning FC layers into convolutional layers:-

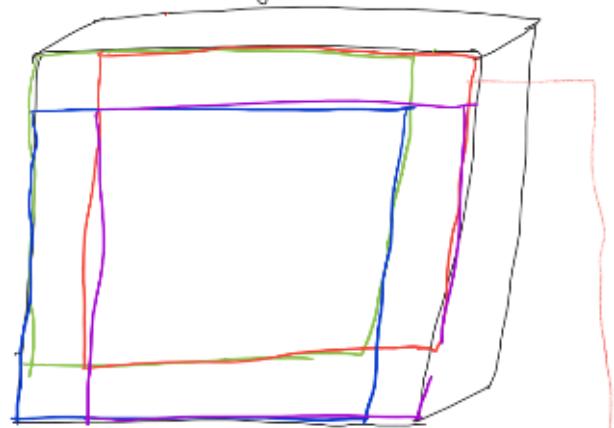
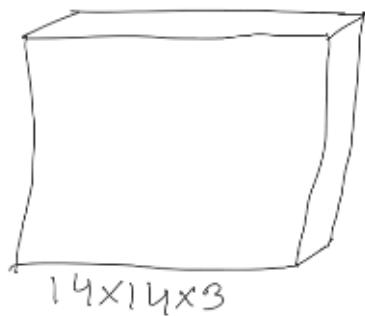


Convolution implementation or sliding window -

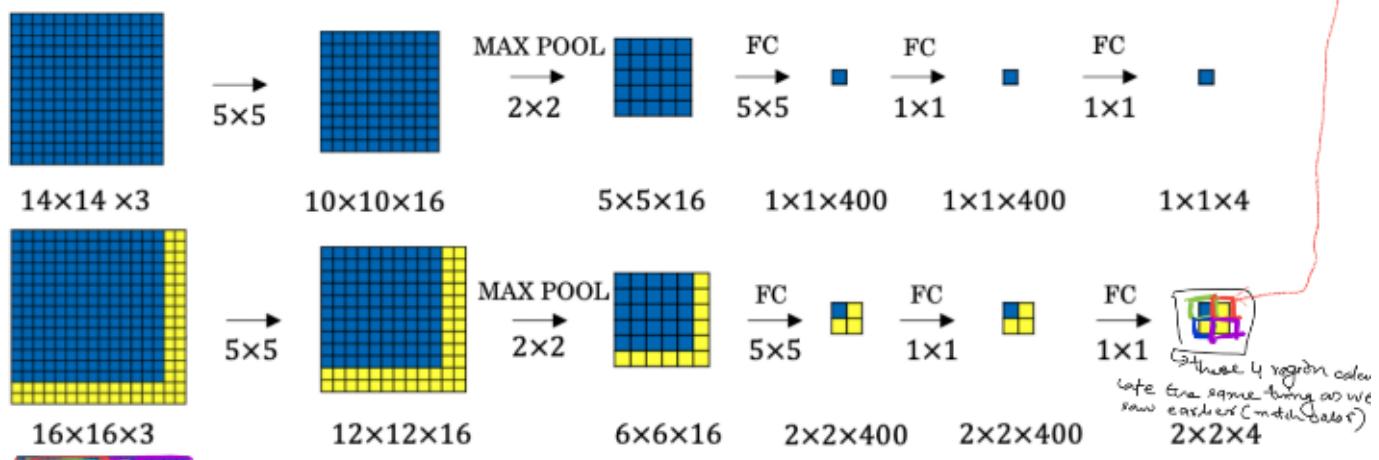
Sliding window of nearby region calculates duplicate information. We can implement convolutional network instead of sliding window to remove the duplicate calculation
For Example:-

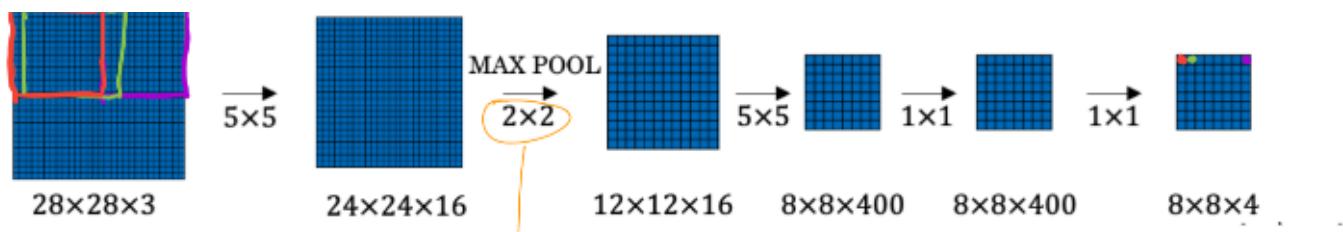
Let,

Our trainset is $14 \times 14 \times 3$ image. } test set is $16 \times 16 \times 3$ so we will use a sliding window of $14 \times 14 \times 3$.



As we can see we have 4 different region after implementing $14 \times 14 \times 3$ sliding window. And most of them share the same region so that leads duplicate calculations. It doesn't seem to be a problem if our sliding window & test set size doesn't differ much. For the above





If we use stride of two & sliding window of $14 \times 14 \times 3$ we get 8×8 result. The red section is the first sliding window the result of that is shown in first block of $8 \times 8 \times 4$ matrix. Which is calculated in a convolution way.

The problem with this approach is the position of the bounding boxes aren't accurate.

YOLO algorithm (You only look once) :-

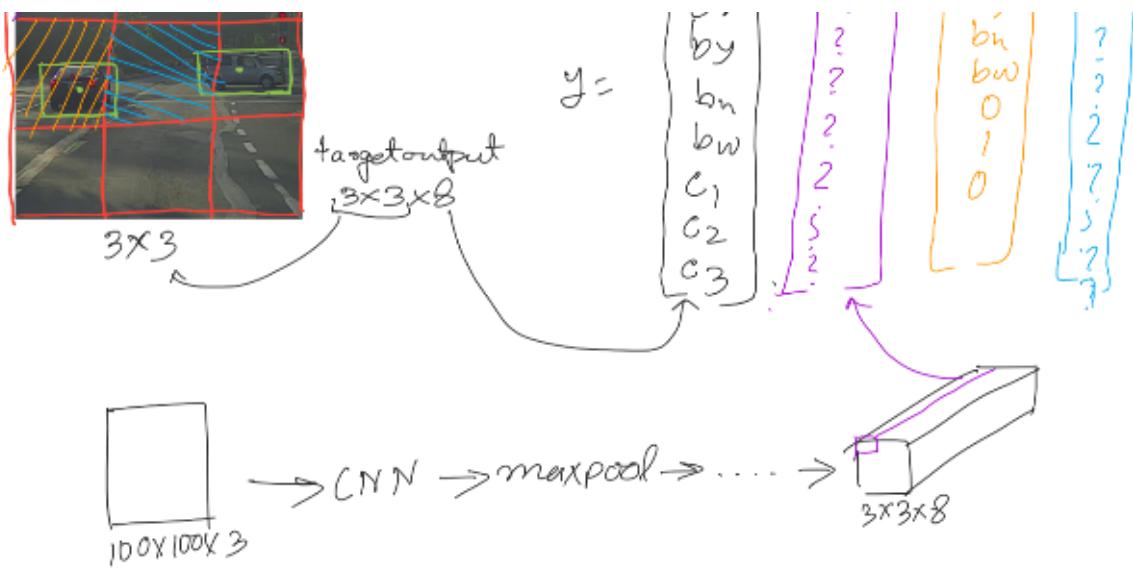
YOLO algorithm solves the previous approach problem. It correctly provide the bounding boxes.

YOLO first make a grid from the picture shown below we show 3×3 grid but in practice the grid size is more dense (19×19). Each of these grid will be evaluated by object detection algorithm as we show in the first algorithm. In the training set if a single object stays in two grid then the object is assigned to the grid where ever the middle point of bounding box reside. The main problem with this approach is if any grid contains more than a object then it cannot detect that



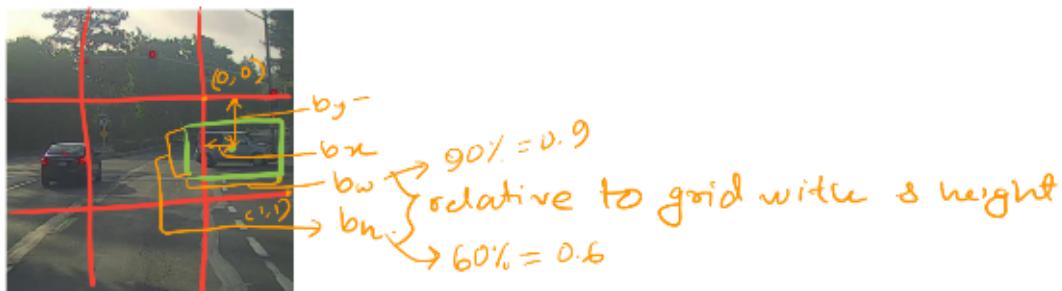
Label for training each grid cell:

$$\begin{bmatrix} P_c \\ b_x \\ b_y \end{bmatrix} \quad \begin{bmatrix} O \\ ? \end{bmatrix} \quad \begin{bmatrix} 1 \\ b_x \\ b_y \end{bmatrix} \quad \begin{bmatrix} O \\ ? \end{bmatrix}$$

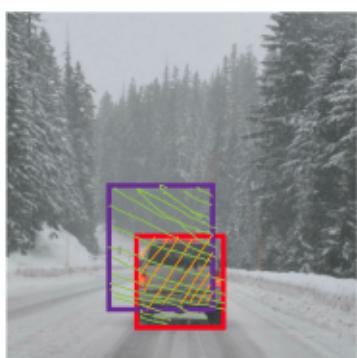


Specify the bounding box:

Bounding boxes is calculated relatively on the grids not on the whole picture.



Evaluating object localization:-



red is the ground truth mask
purple --- detected bounding box.

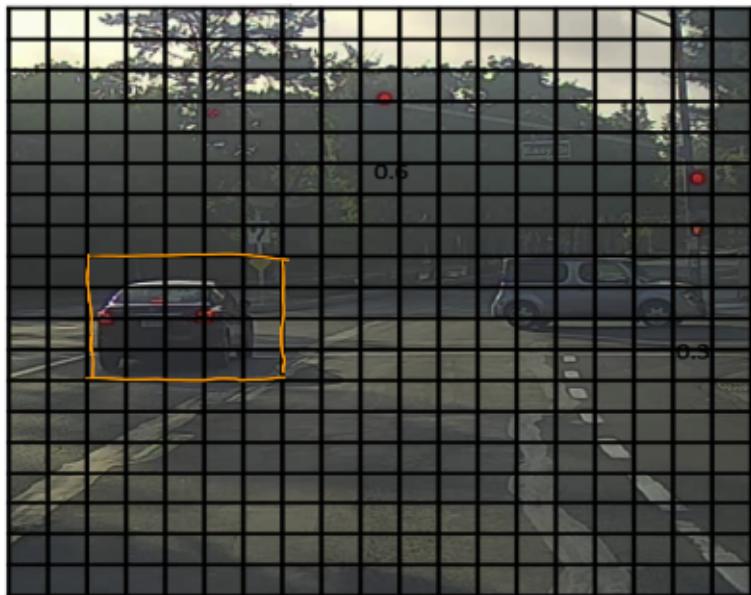
$$\text{Intersection over union (IoU)} = \frac{\text{size of intersection}}{\text{size of union}} = \frac{\text{orange hatched area}}{\text{green hatched area}}$$

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

Non max suppression:-

The problem with making this grid approach is any object

can be identified multiple times.



19×19

portion of this car is present in 20 grids so the car will be recognized 20 times.

Let multiple grid recognizer same car. Each of the grid the their confidence level (P_c) associated with it.



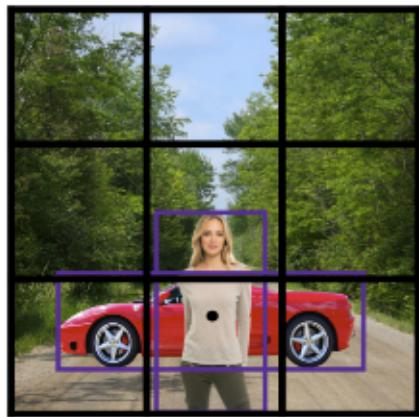
The grid with highest P_c will be considered of having object. The other overlapping grid having larger IoU will be suppressed.

Anchors boxes:-

if we have multiple object in a same grid then

its not possible to detect them with previous approach

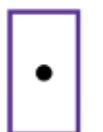
Anchor boxes can help to detect those



Car only

p_c	1	0
b_x	b_x	?
b_y	b_y	?
b_h	b_h	?
b_w	b_w	?
c_1	1	?
c_2	0	?
c_3	0	?
p_c	1	1
b_x	b_x	b_x
b_y	b_y	b_y
b_h	b_h	b_h
b_w	b_w	b_w
c_1	0	0
c_2	1	1
c_3	0	0

Anchor box 1: Anchor box 2:



→ If same anchorbox object (two human) present in the same grid box then this algorithm cannot detect them.

→ If there are three object in the same grid cell this algorithm won't handle that.

→ People choose the size & shape of the anchor boxes by hand

→ we can also use k-means algorithm to choose anchor box