

## + Object Localization

→ Two sub problem, - (classification

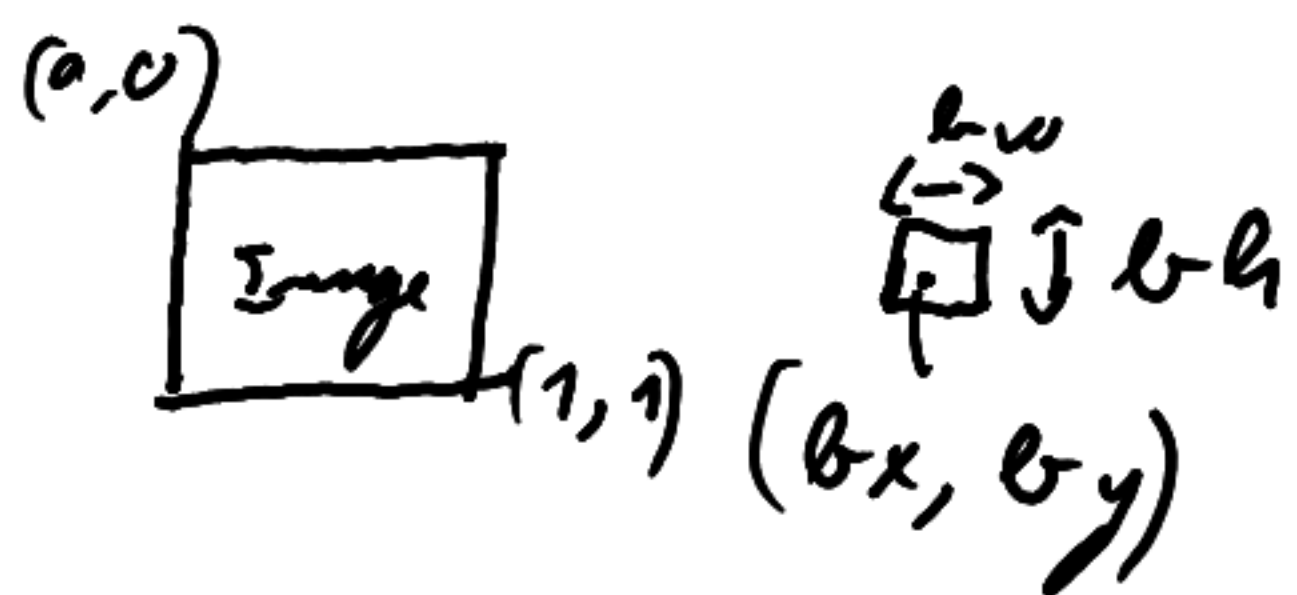
- Classification with localization
- Detection

Image  $\rightarrow$  Conv Net  $\rightarrow \begin{bmatrix} 0 \\ 0 \end{bmatrix} \rightarrow \text{Softmax}$ .

### • Classification with localization

Img  $\rightarrow$  Conv Net  $\rightarrow \begin{bmatrix} 0 \\ 0 \end{bmatrix}_{(1,4)} \rightarrow \text{Softmax}$

If we want to localize the detected img  
an network will output  $(b_x, b_y, b_h, b_w)$   
(bounding box)



So in our example  $\Rightarrow$  4 outputs  $\left\{ \begin{array}{l} 1. \text{ Pedestrian} \\ 2. \text{ Car} \\ 3. \text{ Motorcycle} \\ 4. \text{ Background.} \end{array} \right.$

Also output  $b_x, b_y, b_h, b_w$  class label (1-4)

Ex  $\Rightarrow y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$

$\rightarrow$  Probability of object.

$\rightarrow$  Bounding box.

$\rightarrow$  class

Example for car  $\Rightarrow$   $\begin{bmatrix} 1 \\ b_x \\ b_y \\ b_z \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}$  No object.  $\begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$   $\leftarrow$  "don't care"

Loss function:  $L(\hat{y}, y) = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots$

$\dots + (\hat{y}_8 - y_8)^2$  if  $y_1 = 1$

$(\hat{y}_1 - y_1)^2$  if  $y_1 = 0$

+ Landmark Detection:

If we didn't want a box, just a point, we could use  $l_x, l_y$ .

For example we could have a variety of landmarks  $l_{64x}, l_{64y}$

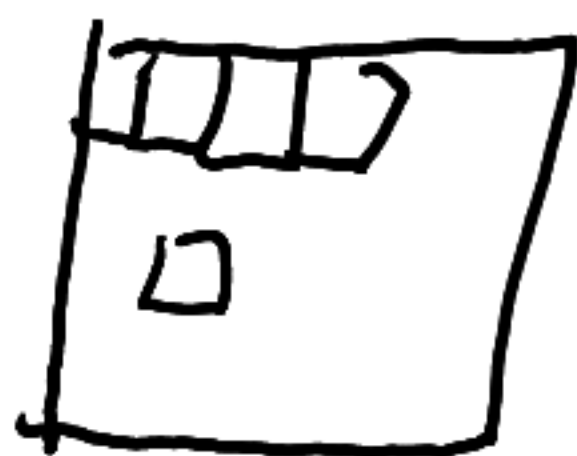
Image  $\Rightarrow$  conv Net  $\Rightarrow \begin{bmatrix} 0 \\ 0 \\ \vdots \end{bmatrix} \rightarrow \begin{bmatrix} \text{face} \\ \rightarrow l_x, \dots, l_n \\ l_y, \dots, l_n \end{bmatrix}$  129 outputs.

## + Object Detection:

Using a conv net with sliding window algorithms

Training set  $\Rightarrow$  closely cropped.

Once we have trained it, The sliding window algorithm, The image is segmented.



$\Rightarrow$  If there is a 1 or a 1 is output

The computational cost is very high.

## + Convolutional implementation of sliding windows

Turning fully connected layers into conv layers.

A FC layer  $\text{vec} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \Rightarrow 1 \times 1 \times 400$   
400 on a conv.

Even in the final layer  $1 \times 1 \times 4$

So To replace.

We crop  $\rightarrow$  We run through sequence not in sequence but in one go using a Conv net.

Main weakness  $\Rightarrow$  Bounding box will be inaccurate.



## + Bounding Box Predictions

YOLO  $\Rightarrow$  you only look once.



A label for each grid cell:

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_x \\ b_y \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \quad (3 \times 3 \times 8)$$

$X$    $100 \times 100 \times 3 \rightarrow$  CNN  $\rightarrow$  Max Pool  $\dots \rightarrow$    $y$ .  $3 \times 3 \times 8$

## + Intersection Over Union

$\frac{\text{Size of Intersection}}{\text{Size of Union}} \Rightarrow$  correct if  $IOU \geq 0.5$ .

## + Non-max Suppression

Some boxes will overlap.

Step 1. Discard boxes with  $p_c \leq 0.6$ .

2. The remaining boxes:

- Pick the box with highest  $p_c \Rightarrow$  prediction.
- Discard any remaining box with  $IOU \geq 0.5$ .

## + Anchor Boxes

Predetermined box shapes. (for different anchors)

So instead of  $y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$   $3 \times 3 \times 8$

$\Rightarrow$  Now  $\Rightarrow y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ \vdots \end{bmatrix}$   $\begin{matrix} = \text{Anchor 1} \\ \vdots \\ = \text{Anchor 2} \end{matrix}$

grid cell, anchor box  
( $3 \times 3 \times 16$ )

## Summary

Previously: Each object in Training image is assigned to grid cell that contains that object's midpoint

With ~~two~~ anchor boxes: Each object in Training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IOU



## + Semantic Segmentation with U-Net.

Rather than object detection, each pixel is labeled,

Used in self driving. or X-ray readings.

### • U-Net (segmentation)

If looking for a car in an image, 0 will be not a car and 1 car.

☺  
segmentation map is generated.

One key step is to make the image bigger  
The height and width will increase the  
deeper we go into the U-Net.

### + Transpose Convolution:

Normal Convolution turns the image into  
a smaller output

A Transpose can do the opposite.

$(2 \times 2) \Rightarrow \xrightarrow[\text{(3x3)}]{\text{Filter}} \text{Filter is applied to the output.} \quad (4 \times 4)$

filter  $f \times f$  (3x3)

Stride  $S=2$

Padding  $P=1$

## + U-Net Architecture

$h \times w \times 3 \Rightarrow \text{U-Net} \Rightarrow h \times w \times \text{ndim}.$