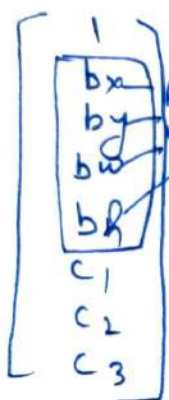


# Yolo A fresh approach to object detection

Yolo ①

- Not classification as detection
- New idea.
  - Evaluated rigorously by the authors
  - The fact that new idea is better is validated by rigorous evaluation.
- Specifying spatially separated bounding boxes and associated class

OD as a regression problem



- Single neural network

- It can be optimized end to end directly on detection problem

- Unified architecture is extremely fast, it outperforms

CVPR 2016

self driving car

- Generalizes better

$\sim 20ms$

5 feet

60km/h  
0.5 km

Yolo learns generalizable representation of objects

- Less likely to break down when applied to new domains

- Accuracy concern

- The first version of the yolo algorithm lacks behind in accuracy

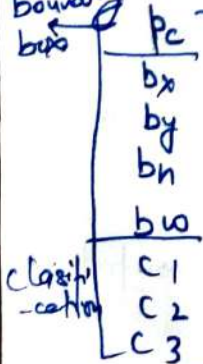
R-CNN, fast R-CNN :- 7-8%

## Preparing the training data

-  $S \times S$  size grid.

- Predict the mid-point to determine the bounding box.

bounding box



$\rightarrow$



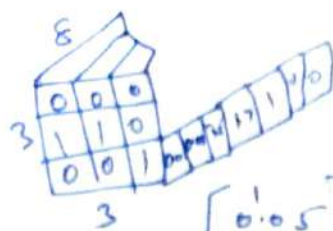
related to grid  
related to image



$> 1$



$8 \times 8$



$$\begin{bmatrix} 0.05 \\ 0.05 \\ 1.5 \\ 1.7 \\ 1 \\ 8 \end{bmatrix}$$

- Intuition for object detection.
- Yolo reasons globally about the image
  - Yolo sees the entire image during training and test time
  - Unlike sliding or R-CNN
  - It encodes contextual information about classes as well as their appearance

IoU Intersection over Union  $3 \times 3 \times 8$

- object detection is working well or not.

IoU: Intersection: Common part

Union:

$$IoU = \frac{\text{area of Intersection}}{\text{area of Union}}$$

measure to determine the bounding box similarity

$$IoU \geq 0.5$$



Understanding the output tensor

- Output tensor size

Why multiple bounding box?

$$s \times s \times (B * 5 + C)$$

size of grid

# bounding box

# of classes

1000 images

car

pedestrian

motorcycle

- Multiple grid cell may claim an object

Test output: - only one object

$$B=2, C=3 \begin{matrix} \rightarrow \text{bicycle} \\ \rightarrow \text{car} \\ \rightarrow \text{dog} \end{matrix}$$

$$13 \times 13 \times (2 * 5 + 3)$$

$$= 13 \times 13 \times 13$$

$$13 \times 13 \times 2 = 169 \times 2 = 338$$

Non max Suppression

Each grid cell predicts.

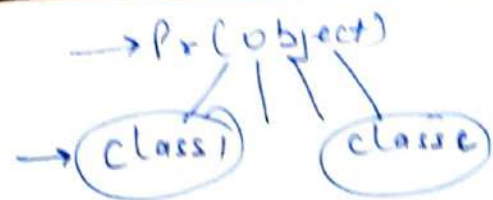
(Probability of confidence)

$$P_c = P_r(\text{object}) * IoU$$

- C conditional class probabilities

$$P_r(\text{Class} | \text{object}) * P_r(\text{object}) * IoU$$





→ Discard  $p_c \leq 0.6$

→ for remaining boxes

highest  $p_c$  → output as prediction

- Discard any remaining boxes  $IoU \geq 0.5$

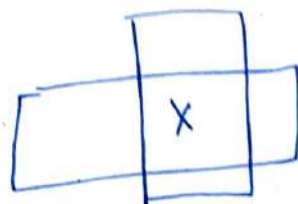
- Independently carry out non-max suppression for each class

Anchor boxes Limitation of Yolo! - one object

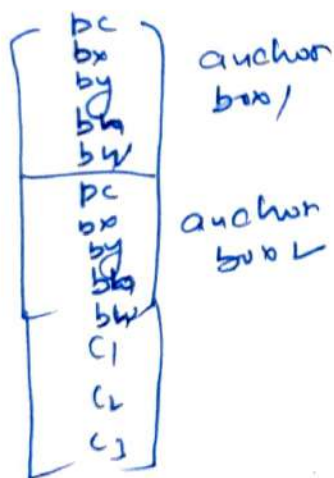
per grid cell - Anchor boxes solve it up to some extent

- can be also reduced by increasing  $s$

Overlapping objects



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



each object → 1 grid cell

mid point.

each object → mid point 1 grid cell + anchor box for the grid cell with highest IoU

Anchor boxes! - help in specialization

- Parts of the network specialize in detection.

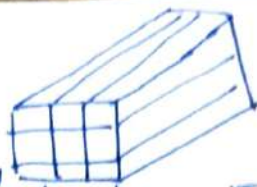
Yolo! - Training with anchor boxes but 4 no. of bounding boxes

$y$  is  $3 \times 3 \times 2 \times 8$

5+3 classes

Image

Trained  
Yolo architecture



Non max suppressed  
output

$P_c < 0.6$   
Discarded

Limitation of YOLO

Improving object detection



Anchor box 1



Anchor box 2

Classification

Classification  
+ Localization

Object  
Detection

You only look once

Instance segmentation  
-tion.

- Why should it work?

• Authors had an idea / intuition.  
That alone is not sufficient.

• Prepare for evaluation.  
This requires laborious task of data set  
preparation for training.  
so you better have confidence in your idea  
intuition.

- The fact that new idea is better is validated  
by rigorous evaluation.

- Train / debug / improve
- Performance on test set
- Generalize well.

- The above steps are true for any new deep  
learning algorithm.  
Repeat them for AlexNet (but of course  
for classification)



