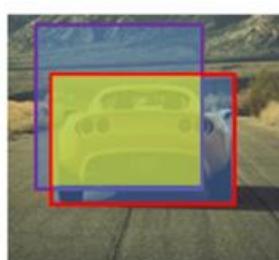


Today lecture:

1. Intersect over union
2. Non maxima suppression
3. Anchor box
4. Complete execution of Yolo algorithm
5. Mean average precision -mAP
6. Transfer learning

Evaluation

- Intersection over union



Intersection over union (IoU)

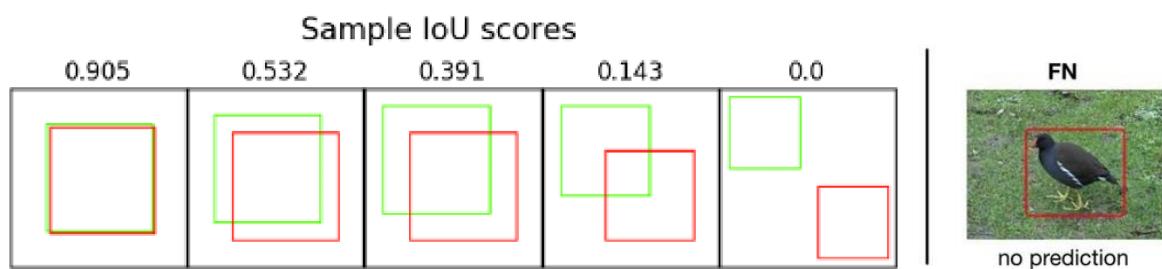
$$= \frac{\text{size of } \square}{\text{size of } \blacksquare}$$

General Criteria:

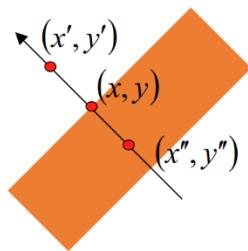
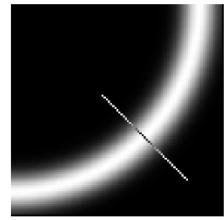
“Correct” if IoU ≥ 0.5

Evaluation

Sample IoU



Non-maximum suppression

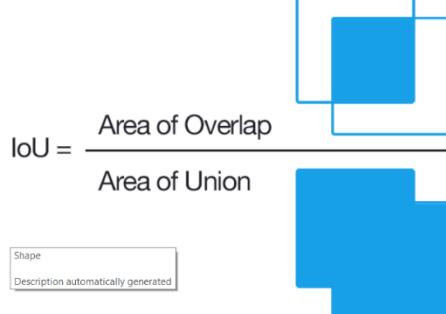


$$M(x, y) = \begin{cases} |\nabla S|(x, y) & \text{if } |\nabla S|(x, y) > |\Delta S|(x', y') \\ & \quad \& |\Delta S|(x, y) > |\Delta S|(x'', y'') \\ 0 & \text{otherwise} \end{cases}$$

x' and x'' are the neighbors of x along normal direction to an edge

Non-maxima suppression (NMS)

- Iterate over all detections
 - Pick the highest scoring box
- Find overlap
 - with all other boxes
- Remove boxes with high overlap
 - Threshold, usually 0.5



Comparing Boxes: Intersection over Union (IoU)



How can we compare our prediction to the ground-truth box?

Intersection over Union (IoU)
(Also called “Jaccard similarity” or
“Jaccard index”):

$$\frac{\text{Area of Intersection}}{\text{Area of Union}}$$



Justin Johnson

Lecture 15 - 36

November 6, 2019

Comparing Boxes: Intersection over Union (IoU)

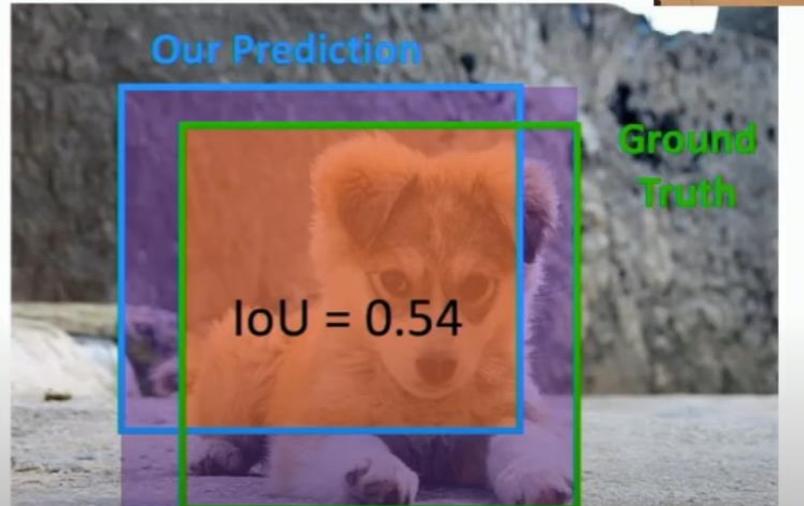


How can we compare our prediction to the ground-truth box?

Intersection over Union (IoU)
(Also called “Jaccard similarity” or
“Jaccard index”):

$$\frac{\text{Area of Intersection}}{\text{Area of Union}}$$

IoU > 0.5 is “decent”



Comparing Boxes: Intersection over Union (IoU)



How can we compare our prediction to the ground-truth box?

Intersection over Union (IoU)
(Also called “Jaccard similarity” or
“Jaccard index”):

$$\frac{\text{Area of Intersection}}{\text{Area of Union}}$$

IoU > 0.5 is “decent”,
IoU > 0.7 is “pretty good”,



Justin Johnson

Lecture 15 - 38

November 6, 2019

Comparing Boxes: Intersection over Union (IoU)

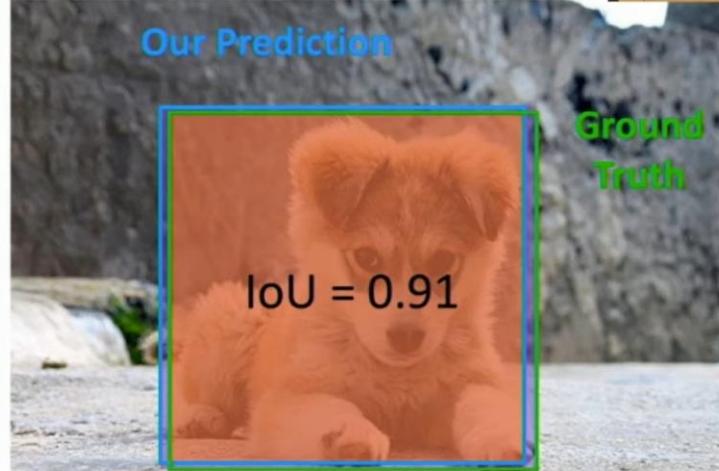


How can we compare our prediction to the ground-truth box?

Intersection over Union (IoU)
(Also called “Jaccard similarity” or
“Jaccard index”):

$$\frac{\text{Area of Intersection}}{\text{Area of Union}}$$

IoU > 0.5 is “decent”,
IoU > 0.7 is “pretty good”,
IoU > 0.9 is “almost perfect”



Justin Johnson

Lecture 15 - 39

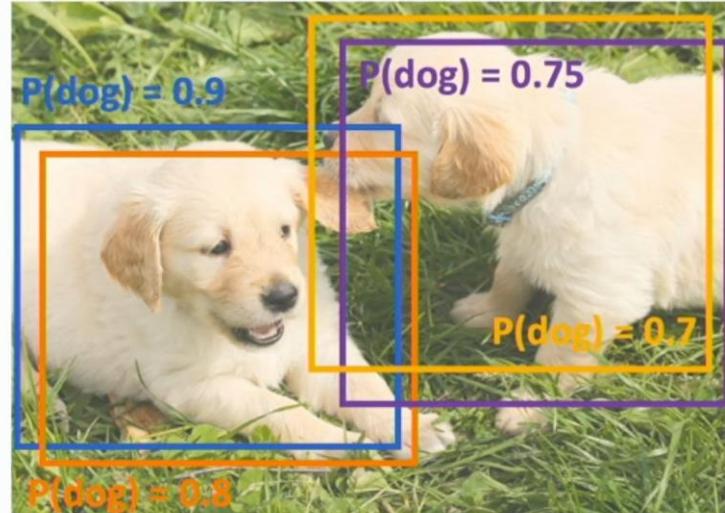
November 6, 2019

Overlapping Boxes

Problem: Object detectors often output many overlapping detections:



Justin Johnson



Funny image is CC0 Public Domain

Lecture 15 - 40

November 6, 2019

Overlapping Boxes: Non-Max Suppression (NMS)

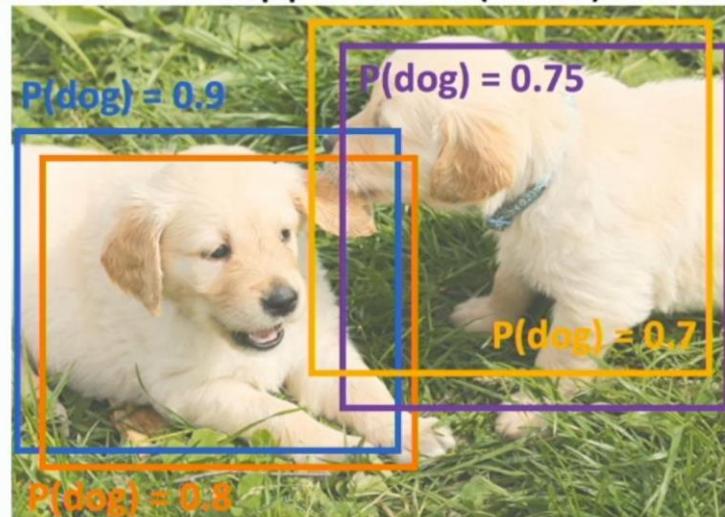
Problem: Object detectors often output many overlapping detections:

Solution: Post-process raw detections using **Non-Max Suppression (NMS)**

1. Select next highest-scoring box
2. Eliminate lower-scoring boxes with $\text{IoU} > \text{threshold}$ (e.g. 0.7)
3. If any boxes remain, GOTO 1



Justin Johnson



Funny image is CC0 Public Domain

Lecture 15 - 41

November 6, 2019

Overlapping Boxes: Non-Max Suppression (NMS)

Problem: Object detectors often output many overlapping detections:

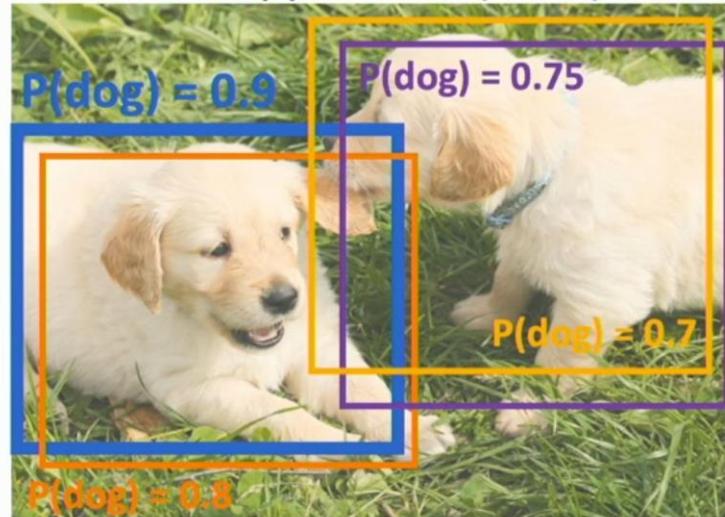
Solution: Post-process raw detections using **Non-Max Suppression (NMS)**

1. Select next highest-scoring box
2. Eliminate lower-scoring boxes with $\text{IoU} > \text{threshold}$ (e.g. 0.7)
3. If any boxes remain, GOTO 1



$$\begin{aligned}\text{IoU}(\text{blue}, \text{orange}) &= 0.78 \\ \text{IoU}(\text{blue}, \text{purple}) &= 0.05 \\ \text{IoU}(\text{blue}, \text{yellow}) &= 0.07\end{aligned}$$

Justin Johnson



Puppy image is CC0 Public Domain

Lecture 15 - 42

November 6, 2019

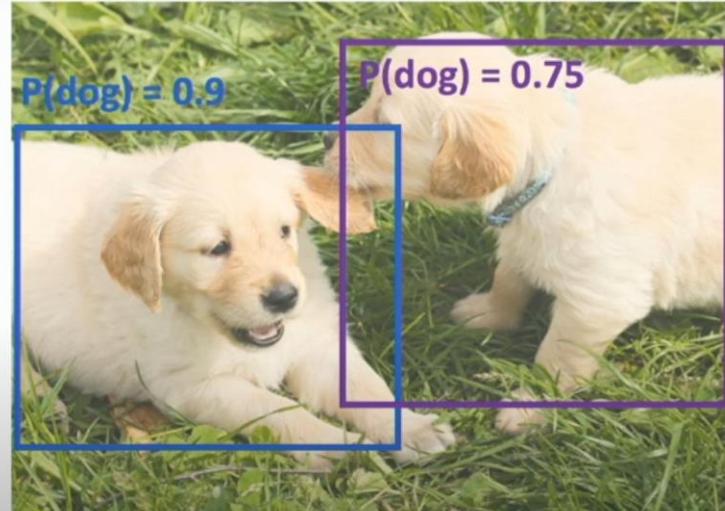
Lecture 15: Object Detection

Overlapping Boxes: Non-Max Suppression (NMS)

Problem: Object detectors often output many overlapping detections:

Solution: Post-process raw detections using **Non-Max Suppression (NMS)**

1. Select next highest-scoring box
2. Eliminate lower-scoring boxes with $\text{IoU} > \text{threshold}$ (e.g. 0.7)
3. If any boxes remain, GOTO 1



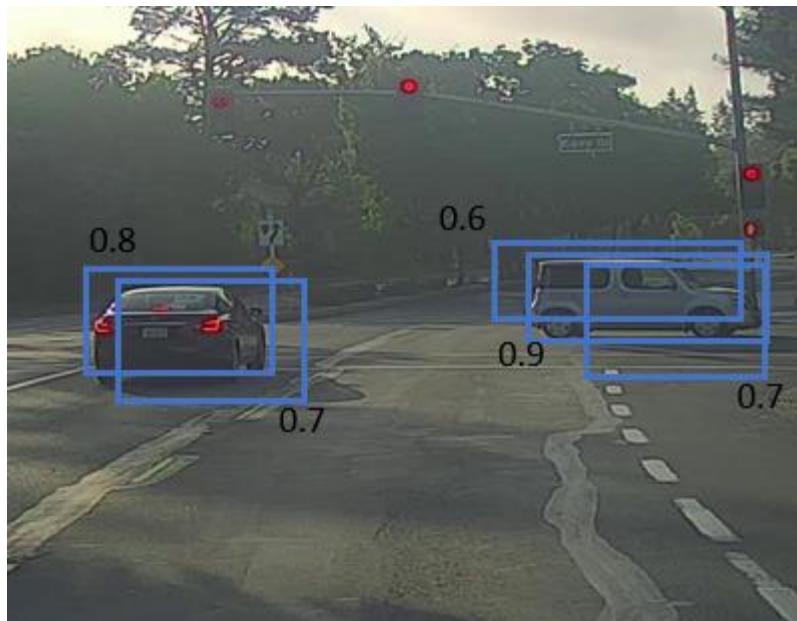
Puppy image is CC0 Public Domain

◀ ▶ Justin Johnson 36:22 / 1:12:31 • Overlapping Boxes: Non-Max Suppression (NMS) 15 - 44

November 6, 2019

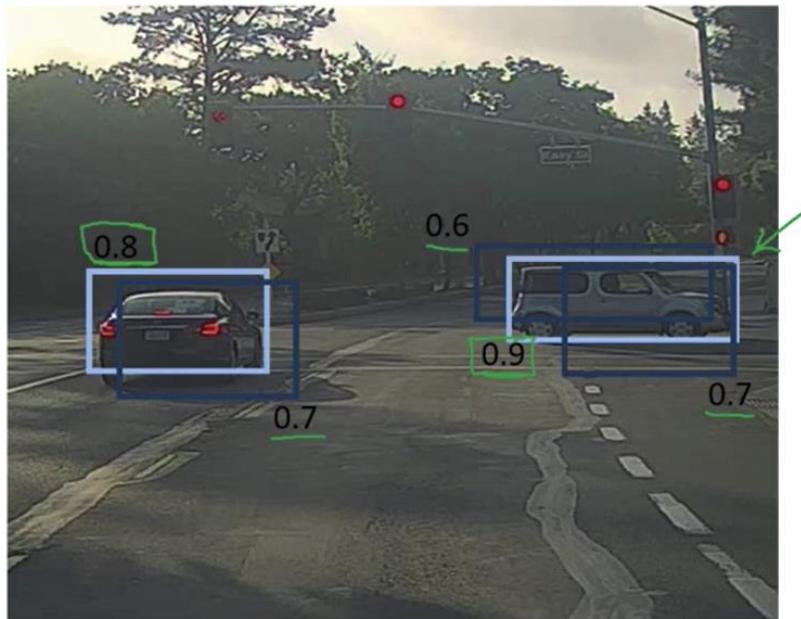
- One of the problems we have addressed in YOLO is that **it can detect an object multiple times**

- Non-max Suppression is a way to make sure that YOLO detects the object just once.
- For example:



i)
ii)

- Discard Box with low P_c
For a specific object class
- Output 0.9 box,
 - discard overlapping boxes 0.6 and 0.7, as have $IOU \geq 0.5$
 - Output 0.8 box
 - discard overlapping box 0.7, as have $IOU >= 0$.
 - We left with 2 boxes

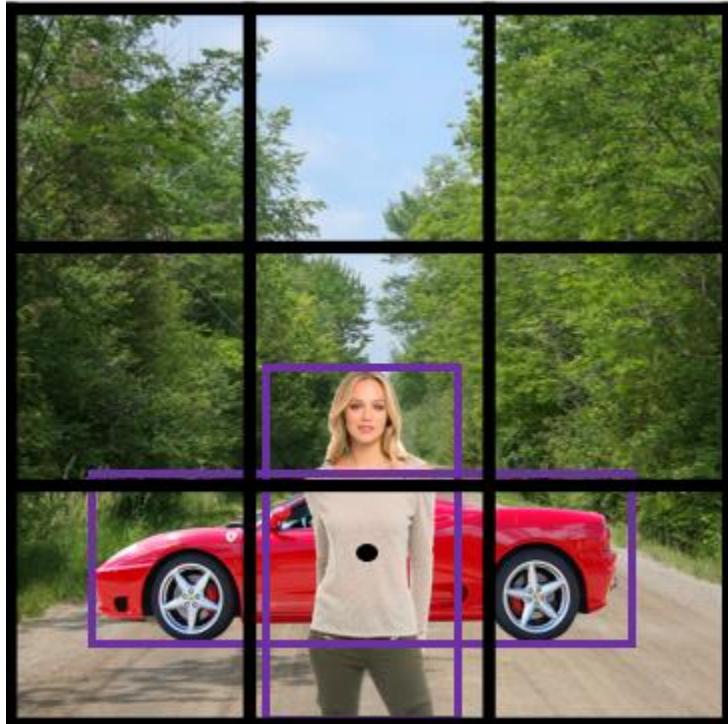


- Non-max suppression algorithm:
 - Lets assume that we are targeting **one class as an output class**.
 - Y shape should be $[P_c, bx, by, bh, hw]$ Where P_c is the probability if that object occurs.
 - **Discard all boxes with $P_c \leq 0.6$ (from prediction volume)**
 - While there are any remaining boxes:

- a. Pick the box with the **largest P_c** Output that as a prediction.
- b. **Discard** any remaining box with $\text{IoU} > 0.5$ with that box output in the previous step i.e any box with high overlap (greater than overlap threshold of 0.5).
- If there are **multiple classes/object types** c you want to detect, you should run the Non-max suppression c times, once for every output class.

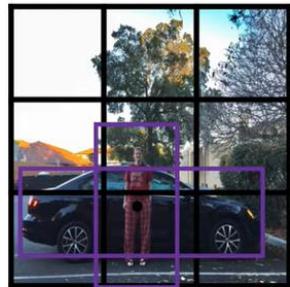
- **Anchor Boxes**

- In YOLO, a grid only detects one object. What if a grid cell wants to detect multiple objects?

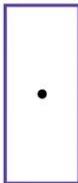


- Car and person rectangle center belongs to same region.
- In practice this happens rarely.
- Overlapping objects:

Overlapping objects:



Anchor box 1:



Anchor box 2:

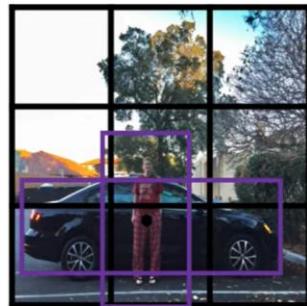


$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

[Redmon et al., 2015, You Only Look Once: Unified real-time object detection]

Andrew Ng

Overlapping objects:



Anchor box 1:



Anchor box 2:



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

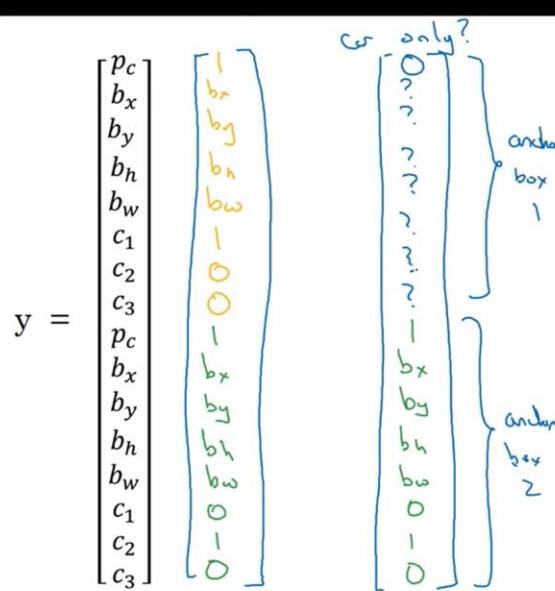
[Redmon et al., 2015, You Only Look Once: Unified real-time object detection]

Andrew Ng

Anchor box example



Anchor box 1: Anchor box 2:



Andrew Ng

- The idea of Anchor boxes helps us solving this issue.
- If $Y = [P_c, b_x, b_y, b_h, b_w, c_1, c_2, c_3]$ Then to use two anchor boxes like this:
 - $Y = [P_c, b_x, b_y, b_h, b_w, c_1, c_2, c_3, P_c, b_x, b_y, b_h, b_w, c_1, c_2, c_3]$ We simply have repeated the one anchor Y .
 - The two anchor boxes you choose should be known as a shape:

Anchor box 1: Anchor box 2:



- So 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. You have to check where your object should be based on its rectangle closest to which anchor box.
- Example of data:

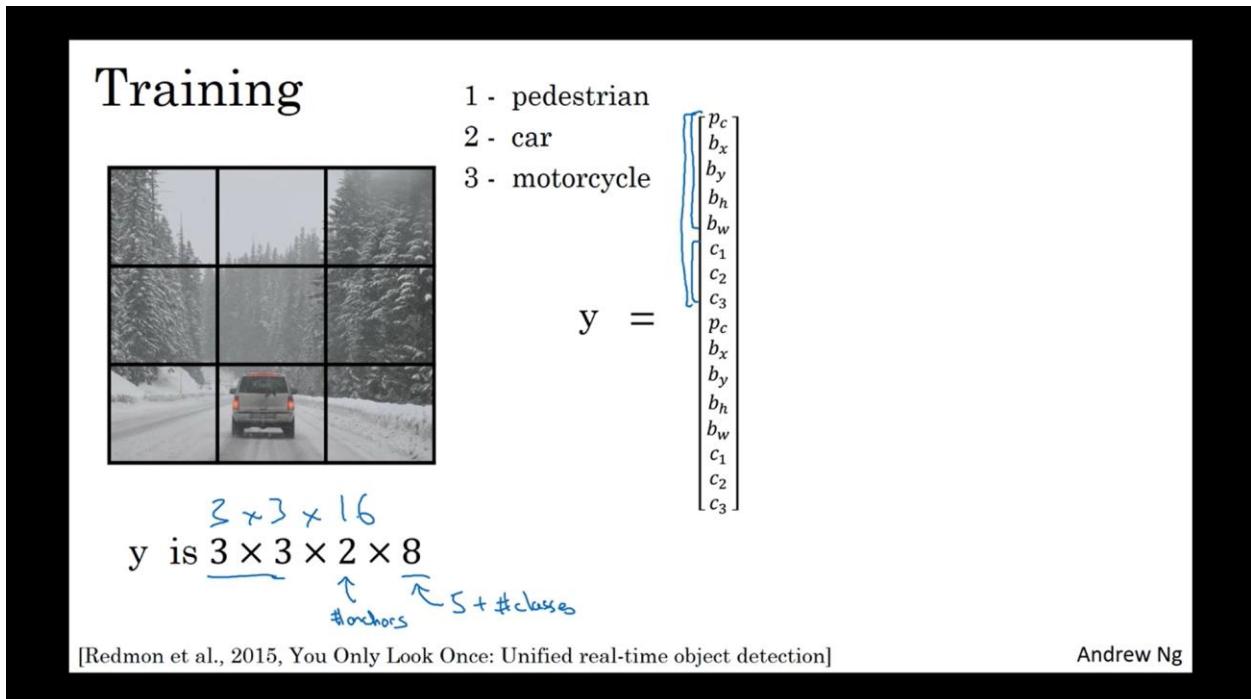
$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

or only?

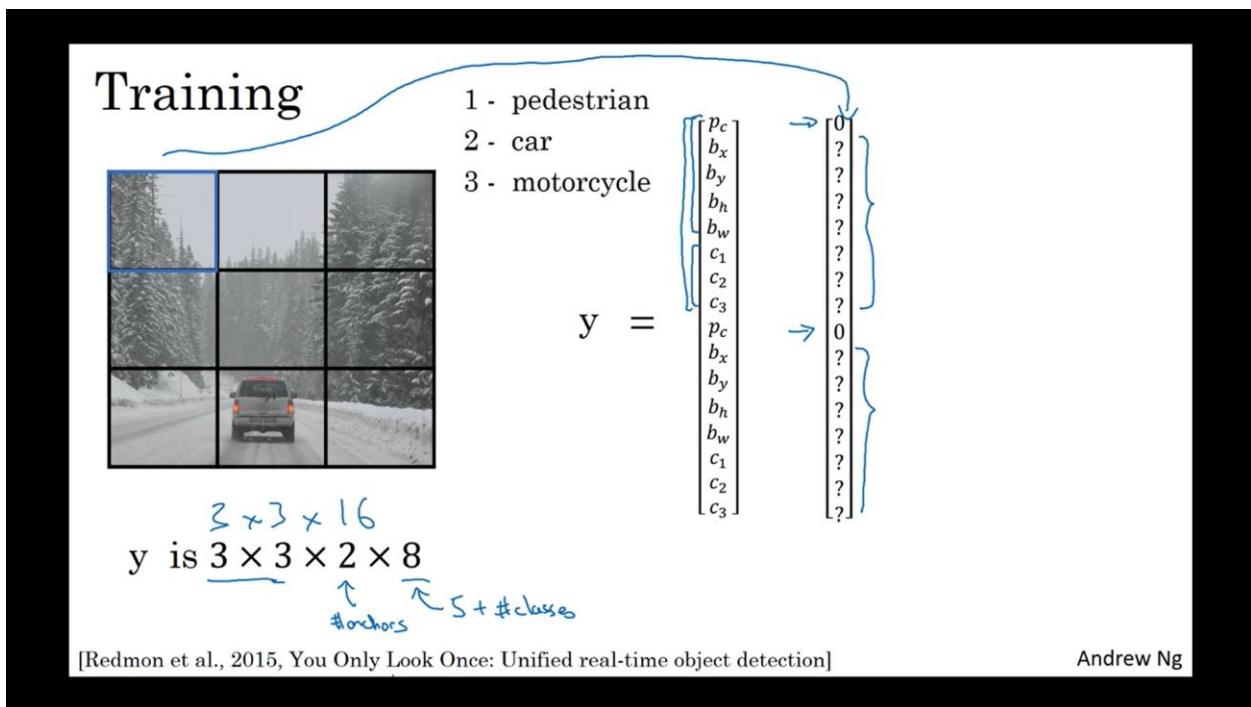
- Where the car was near the anchor 2 than anchor 1.
- You may have two or more anchor boxes but you should know their shapes.
 - how do you choose the anchor boxes and people used to just choose them by hand. Maybe five or ten anchor box shapes that spans a variety of shapes that cover the types of objects you seem to detect frequently.
 - You may also use a k-means algorithm on your dataset to specify that.
- Anchor boxes allows your algorithm to specialize, means in our case to easily detect wider images or taller ones.

- YOLO Algorithm

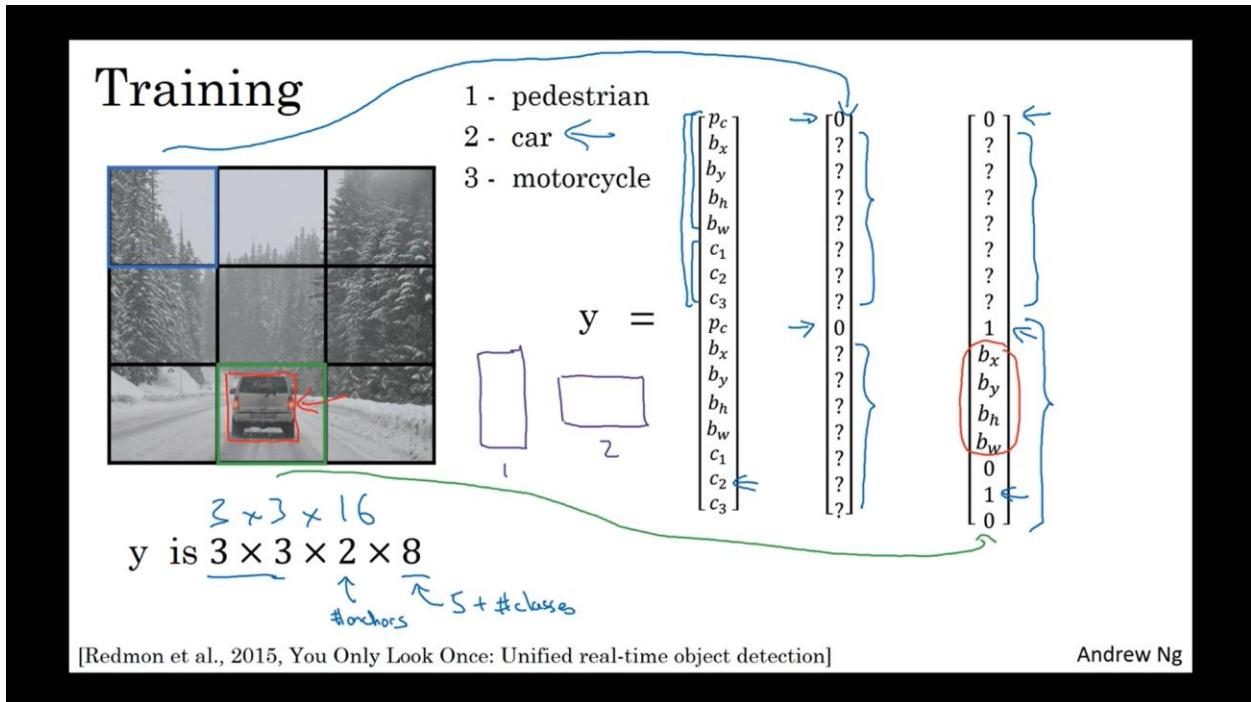
Training Dataset : x is input image of size $19 \times 19 \times 3$, y is output of size $3 \times 3 \times 16 = 3 \times 3 \times 2 \times 8$



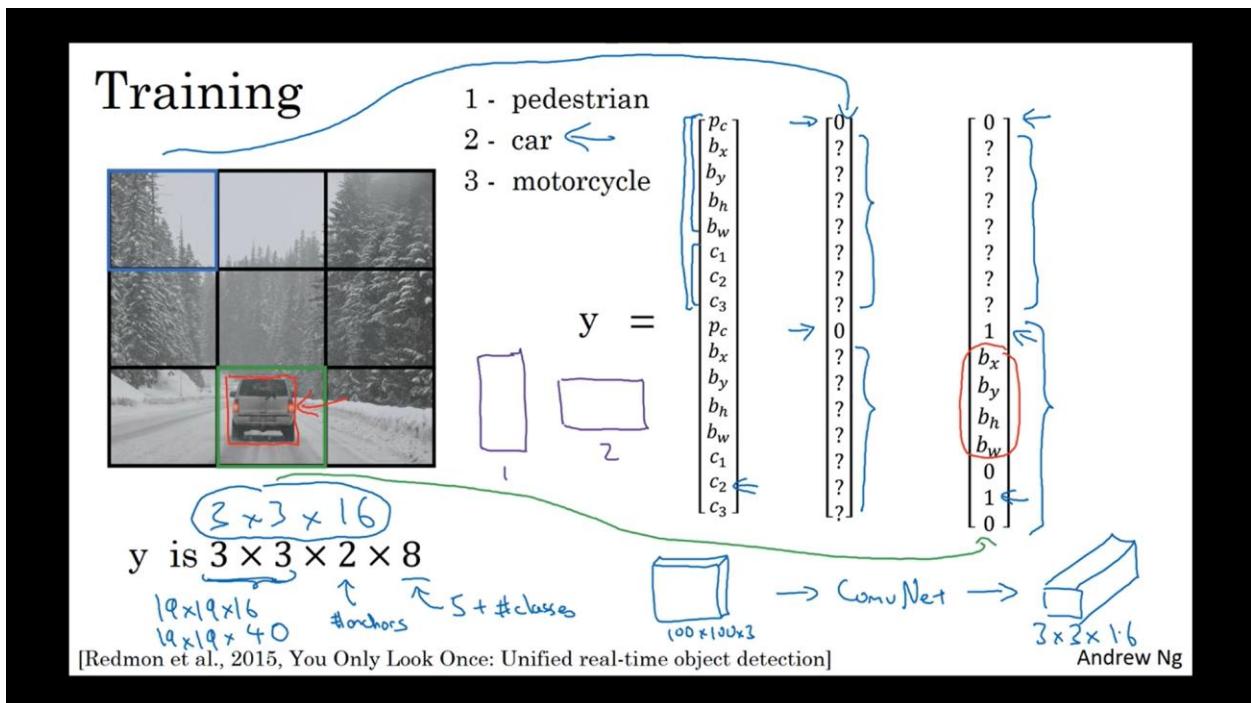
- If a $14 \times 14 \times 3$ grid does not contain any object, the y vector $[0, ?, ?, ?, ?, ?, ?, ?, 0, ?, ?, ?, ?, ?, ?, ?, ?]$



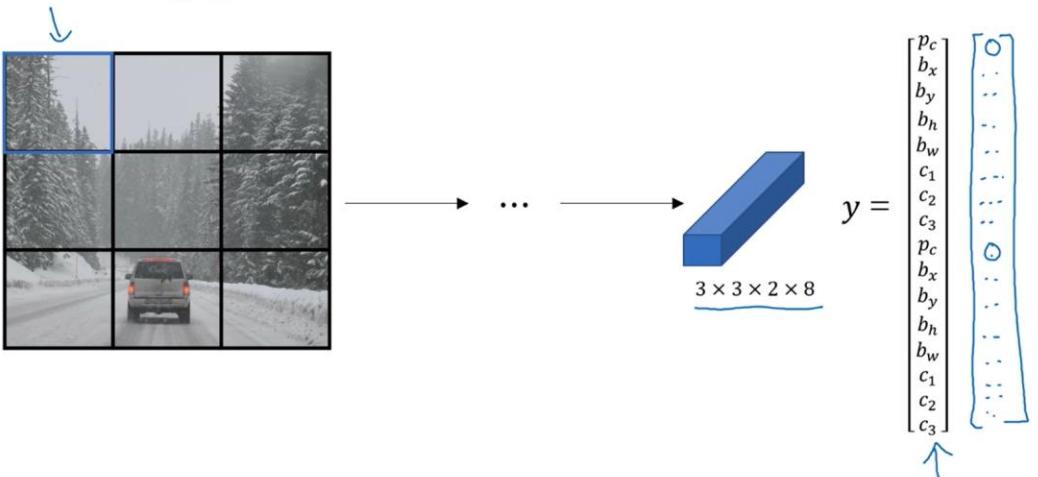
If a grid of 14x14 contains the object then $y = [0, ??????, 1, bx, by, bw, bh, 0, 1, 0]$



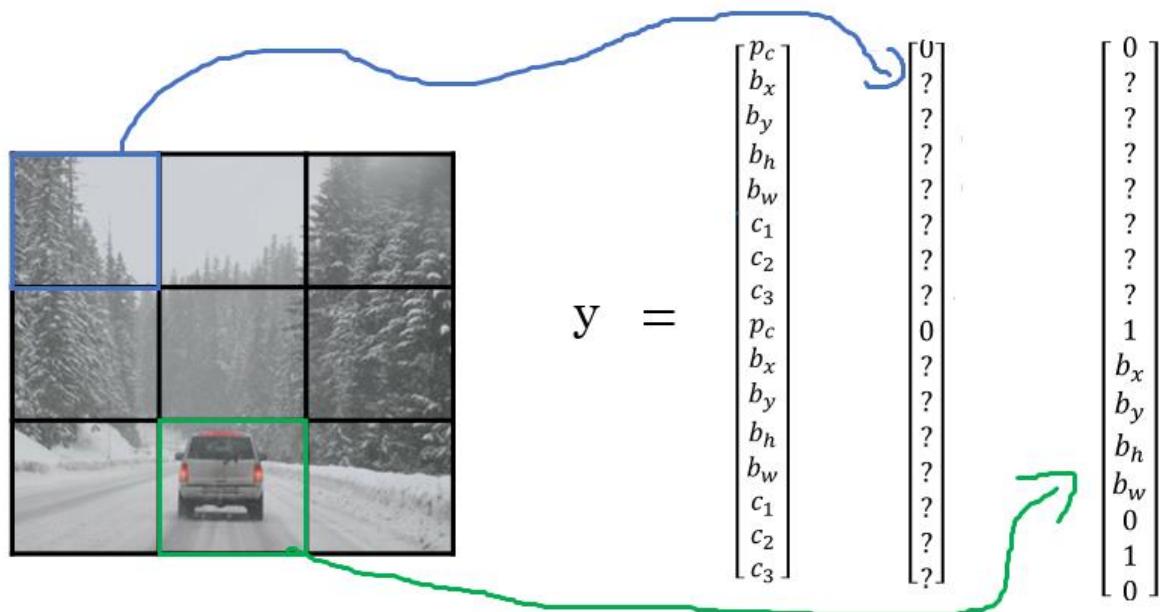
See the relationship between input image $X(100 \times 100 \times 3)$ and ouput $y(19 \times 19 \times ?)$ dimensions



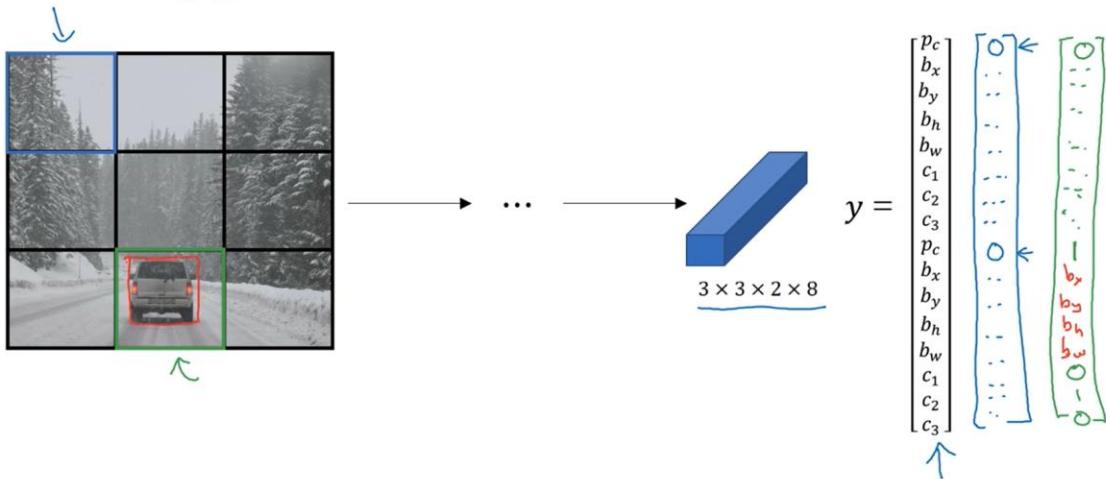
Making predictions



Andrew Ng



Making predictions



Andrew Ng

Outputting the non-max suppressed outputs



- For each grid cell, get 2 predicted bounding boxes.

Andrew Ng

Outputting the non-max suppressed outputs



- For each grid call, get 2 predicted bounding boxes.
- Get rid of low probability predictions.

Andrew Ng

Outputting the non-max suppressed outputs

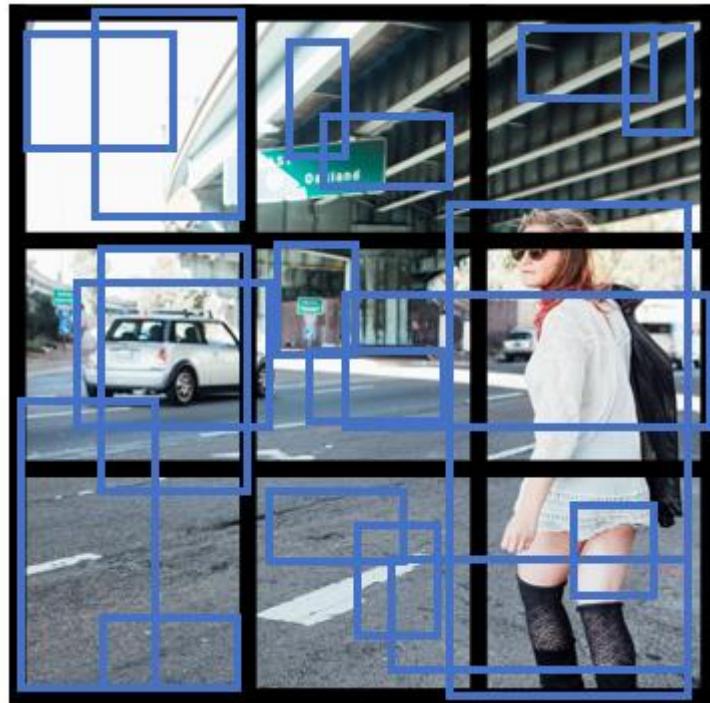


- For each grid call, get 2 predicted bounding boxes.
- Get rid of low probability predictions.
- For each class (pedestrian, car, motorcycle) use non-max suppression to generate final predictions.

Andrew Ng

.....

- YOLO is a state-of-the-art object detection model that is fast and accurate
- Lets sum up and introduce the whole YOLO algorithm given an example.
- Suppose we need to do object detection for our autonomous driver system. It needs to identify three classes:
 - i. Pedestrian (Walks on ground).
 - ii. Car.
 - iii. Motorcycle.
- We decided to choose two anchor boxes, a taller one and a wide one.
 - i. Like we said in practice they use five or more anchor boxes hand made or generated using k-means.
- Our labeled Y shape will be [Ny, HeightOfGrid, WidthOfGrid, 16], where Ny is number of instances and each row (of size 16) is as follows:
 - i. [Pc, bx, by, bh, bw, c1, c2, c3, Pc, bx, by, bh, bw, c1, c2, c3]
- Your dataset could be an image with a multiple labels and a rectangle for each label, we should go to your dataset and make the shape and values of Y like we agreed.
 - i. An example:
 - ii. We first initialize all of them to zeros and ?, then for each label and rectangle choose its closest grid point then the shape to fill it and then the best anchor point based on the IOU. so that the shape of Y for one image should be [HeightOfGrid, WidthOfGrid, 16]
- Train the labeled images on a Conv net. you should receive an output of [HeightOfGrid, WidthOfGrid, 16] for our case.
- To make predictions, run the Conv net on an image and run Non-max suppression algorithm for each class you have in our case there are 3 classes.
 - i. You could get something like that:



- Total number of generated boxes are $\text{grid_width} * \text{grid_height} * \text{no_of_anchors} = 3 \times 3 \times 2$
- ii. By removing the low probability predictions you should have:



- iii. Then get the best probability followed by the IOU filtering:

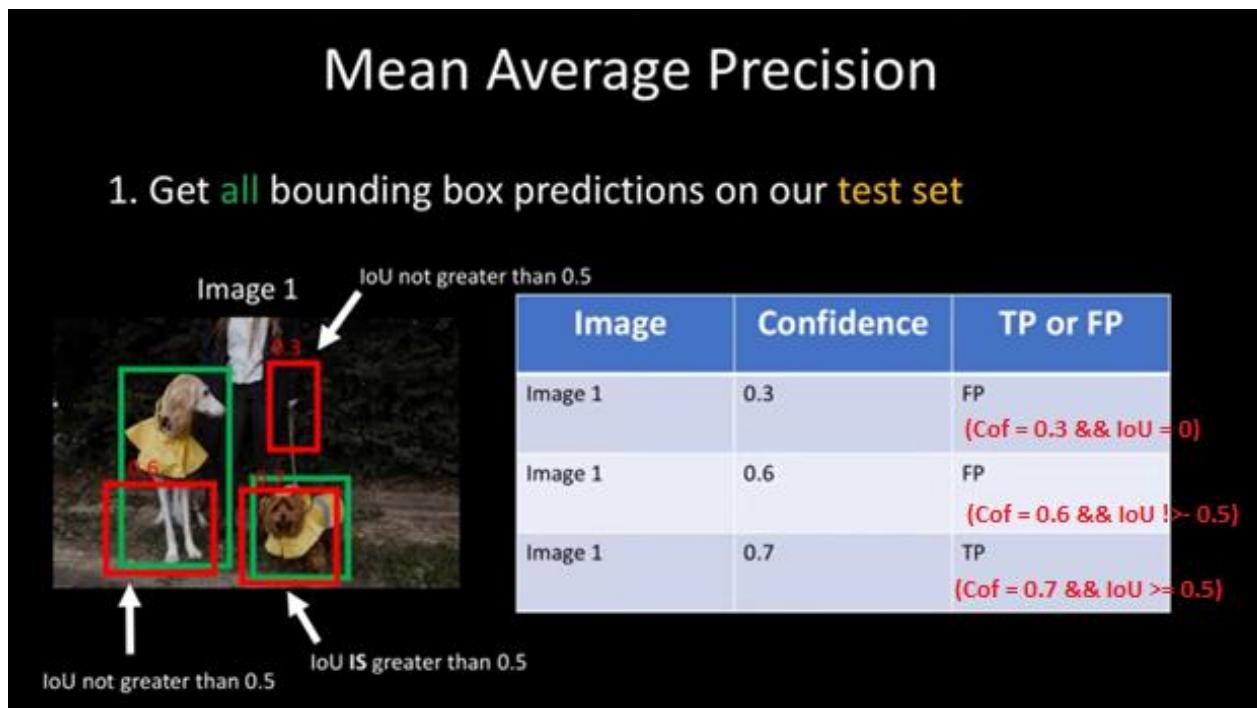
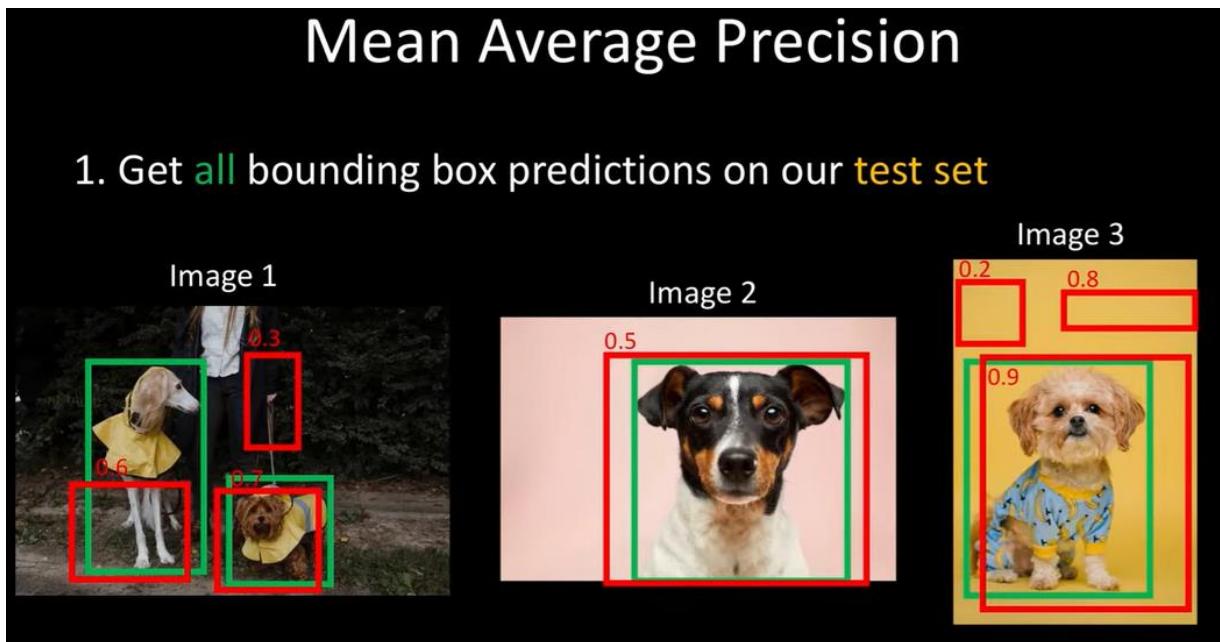


- YOLO are not good at detecting smaller object.
- [YOLO9000 Better, faster, stronger](#)

.....

Mean Average Precision: MAP is a key evaluation metric in object detection that assesses a model's accuracy across all classes by considering both classification and localization performance

<https://www.youtube.com/watch?v=FppOzcDvaDI>



Mean Average Precision

1. Get **all** bounding box predictions on our **test set**

Image 2

A photograph of a dog with a bounding box drawn around it. The confidence score '0.5' is displayed in the top-left corner of the red bounding box.

Image	Confidence	TP or FP
Image 2	0.5	TP

As (Confidence ≥ 0.5 && IoU ≥ 0.5) so it is TP

If (Confidence ≥ 0.5 && IoU ≥ 0.5) then it is TP, else it is FP

Mean Average Precision

1. Get **all** bounding box predictions on our **test set**

Image 3

A photograph of a dog wearing a patterned shirt, with three bounding boxes drawn around different parts of its body. The confidence scores are 0.2, 0.8, and 0.9, respectively.

Image	Confidence	TP or FP
Image 3	0.2	FP (As Conf = 0.2 && IOU= 0)
Image 3	0.8	FP (As Conf = 0.8 && IOU= 0)
Image 3	0.9	TP (As Conf = 0.9 && IOU ≥ 0.5)

1. Get **all** bounding box predictions on our **test set**

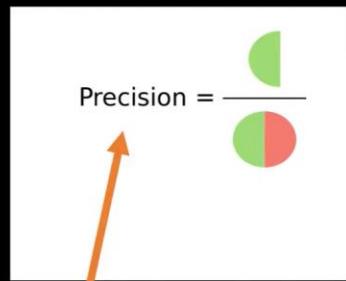
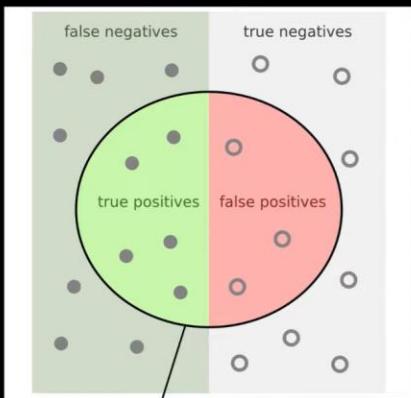
Image	Confidence	TP or FP
Image 1	0.3	FP
Image 1	0.6	FP
Image 1	0.7	TP
Image 2	0.5	TP
Image 3	0.2	FP
Image 3	0.8	FP
Image 3	0.9	TP

Mean Average Precision

2. Sort by **descending confidence** score

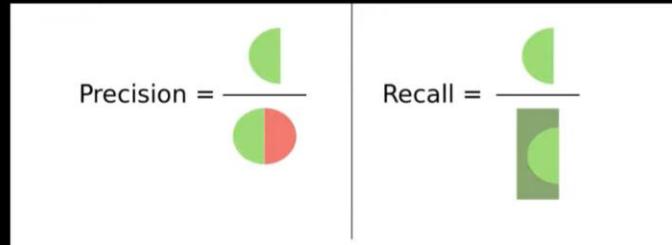
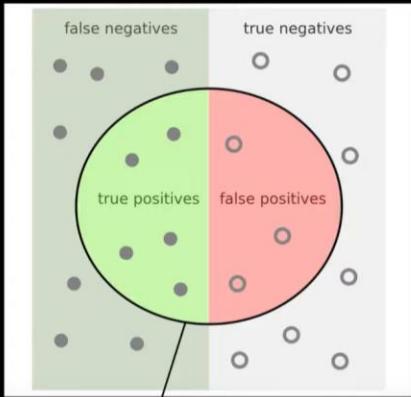
Image	Confidence	TP or FP
Image 3	0.9	TP
Image 3	0.8	FP
Image 1	0.7	TP
Image 1	0.6	FP
Image 2	0.5	TP
Image 1	0.3	FP
Image 3	0.2	FP

Precision and Recall

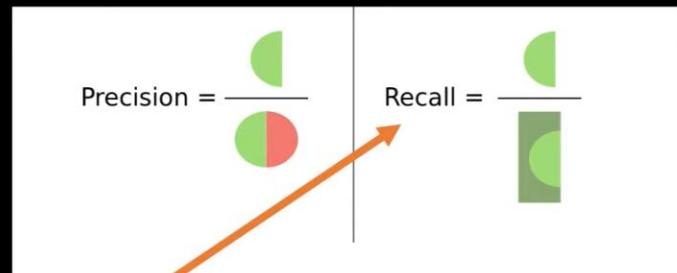
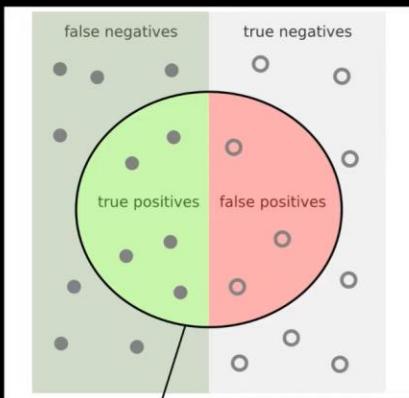


(Of **all** bounding box **predictions**, what fraction was actually correct?)

Precision and Recall

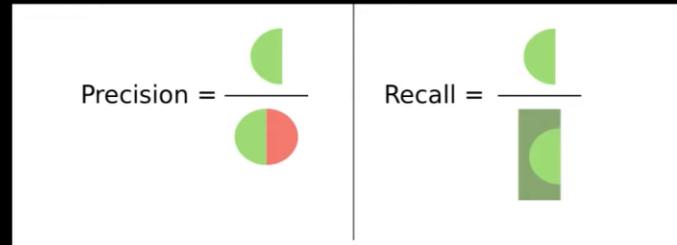
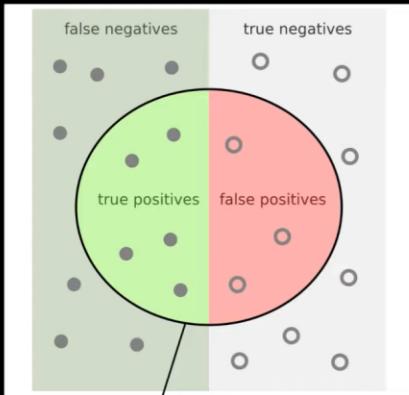


Precision and Recall



(Of all target bounding boxes, what fraction did we correctly detect?)

Precision and Recall



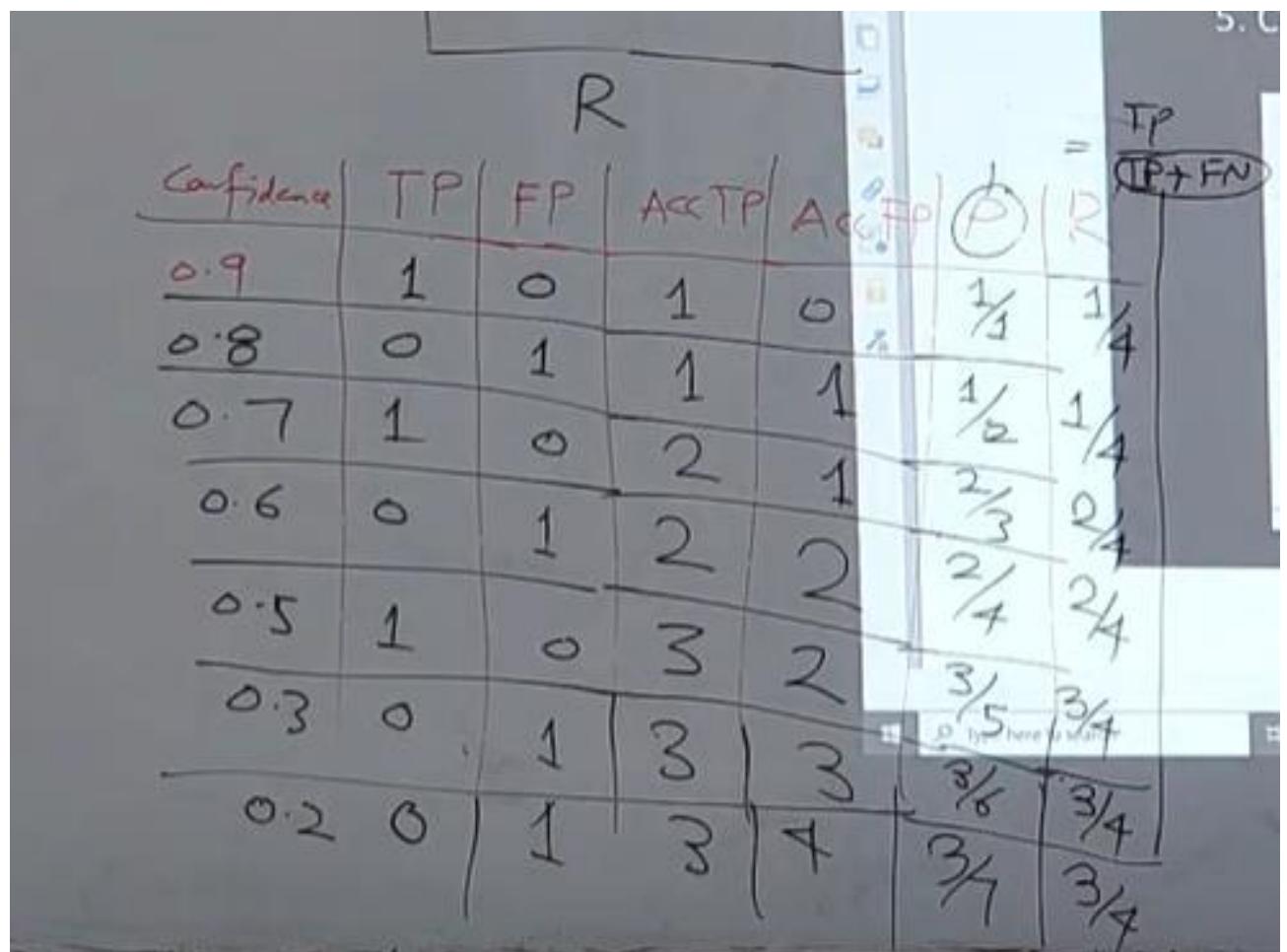
There's a **battle** between precision and recall! **Different applications** may prioritize recall and others precision

Not missing any person in prediction in self-driving car (Relax Pc and IOU threshold)

Mean Average Precision

3. Calculate the **Precision** and **Recall** as we go through all outputs

Image	Confidence	TP or FP	Precision	Recall
Image 3	0.9	TP		
Image 3	0.8	FP		
Image 1	0.7	TP		
Image 1	0.6	FP		
Image 2	0.5	TP		
Image 1	0.3	FP		
Image 3	0.2	FP		



Mean Average Precision

3. Calculate the Precision and Recall as we go through all outputs

Image	Confidence	TP or FP	Precision	Recall
Image 3	0.9	TP	1 / 1	1 / 4
Image 3	0.8	FP	1 / 2	1 / 4
Image 1	0.7	TP	2 / 3	2 / 4
Image 1	0.6	FP	2 / 4	2 / 4
Image 2	0.5	TP	3 / 5	3 / 4
Image 1	0.3	FP	3 / 6	3 / 4
Image 3	0.2	FP	3 / 7	3 / 4

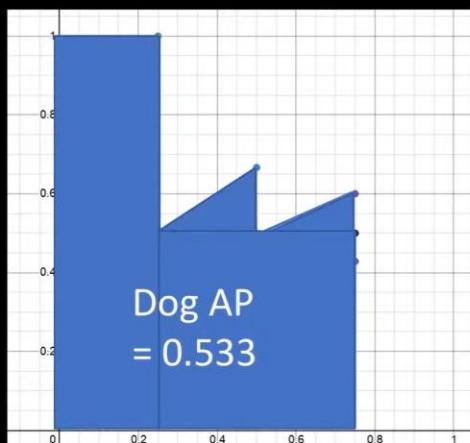
Mean Average Precision

4. Plot the Precision-Recall graph

Precision	Recall
1 / 1	1 / 4
1 / 2	1 / 4
2 / 3	2 / 4
2 / 4	2 / 4
3 / 5	3 / 4
3 / 6	3 / 4
3 / 7	3 / 4

Mean Average Precision

5. Calculate **Area** under PR curve



Precision	Recall
1 / 1	1 / 4
1 / 2	1 / 4
2 / 3	2 / 4
2 / 4	2 / 4
3 / 5	3 / 4
3 / 6	3 / 4
3 / 7	3 / 4

Mean Average Precision

6. This was **only** for dog class, we need to calculate for **all classes**. Let's say we do this for cats and dogs

- **Cat** AP = 0.74
- **Dog** AP = 0.533

$$mAP = (0.533 + 0.74) / 2 = 0.6365$$

Mean Average Precision

7. All this was calculated given a **specific IoU** threshold of 0.5, we need to redo all computations for many IoUs, example: **0.5, 0.55, 0.6, ..., 0.95**. Then **average this** and this will be our **final result**. This is what is meant by mAP@0.5:0.05:0.95

While F1 score can be considered for evaluating object detection, mAP generally offers more comprehensive and informative insights