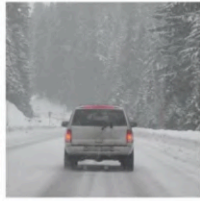


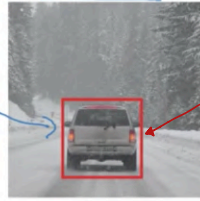
What are localization and detection?

Image classification



"Car"

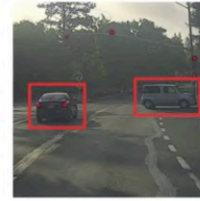
Classification with localization



"Car"

bounding box

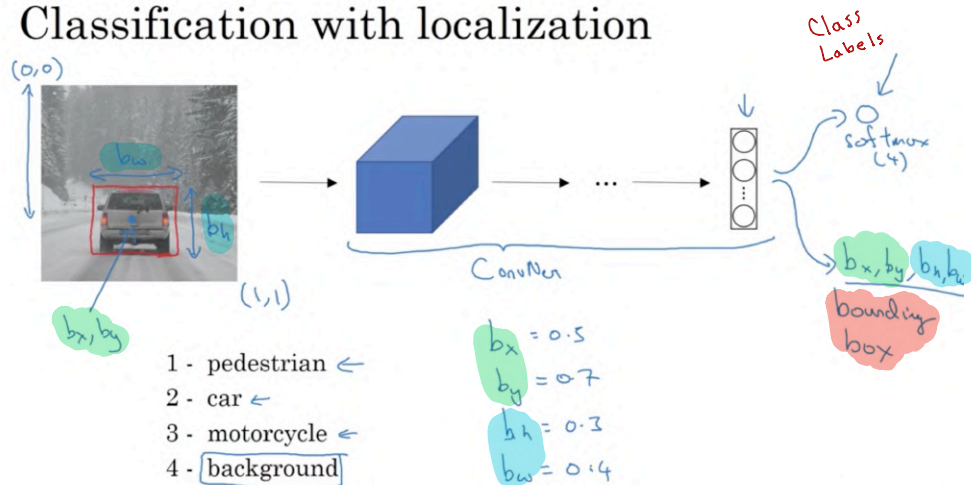
Detection



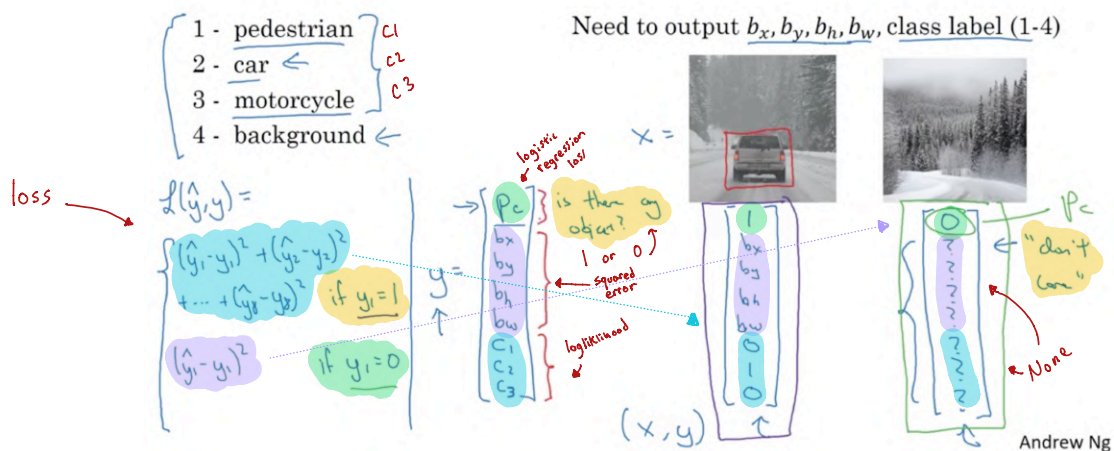
multiple objects

Andrew Ng

Classification with localization



Defining the target label y



Landmark detection

Diagram illustrating landmark detection using a ConvNet. The input image is processed by a ConvNet to produce a feature map. The feature map is then used to detect landmarks, such as the eyes, nose, and mouth, which are represented by red dots. The output is a bounding box around the face, labeled with b_x, b_y, b_h, b_w .

Handwritten notes: $l_{1x}, l_{1y}, l_{2x}, l_{2y}, l_{3x}, l_{3y}, \dots, l_{64x}, l_{64y}$ are grouped under x, y . A note says "e.g. Find eyes → Points on face". Another note says "Face?".

Andrew Ng

Car detection example

Diagram illustrating car detection using a ConvNet. The training set consists of images x and labels y . The images are processed by a ConvNet to produce a feature map, which is then used to detect cars. The output is a bounding box around the car, labeled with b_x, b_y, b_h, b_w .

Training set:

x	y
	1
	1
	1
	0
	0

Diagram showing a car image being processed by a ConvNet to produce a feature map, which is then used to detect cars. The output is a bounding box around the car, labeled with b_x, b_y, b_h, b_w .

Andrew Ng

Sliding windows detection



Sliding windows detection



Sliding windows detection



Pass crops to have it classify 0 or 1

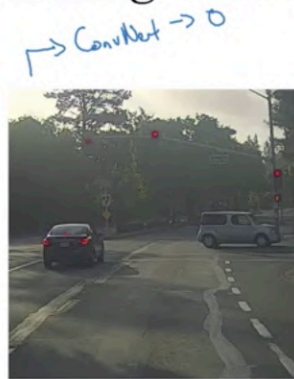
Sliding windows detection



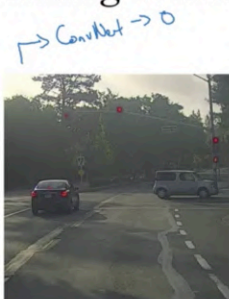
Sliding windows detection



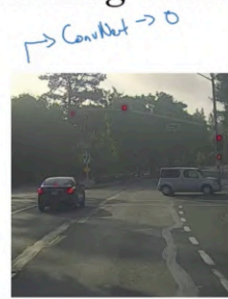
Sliding windows detection



Sliding windows detection



Sliding windows detection



Sliding windows detection

→ ConvNet → 0



→ ConvNet



Sliding windows detection

→ ConvNet → 0

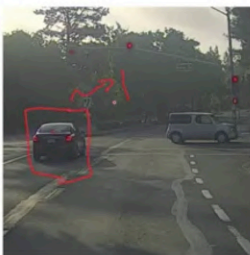


→ ConvNet



Sliding windows detection

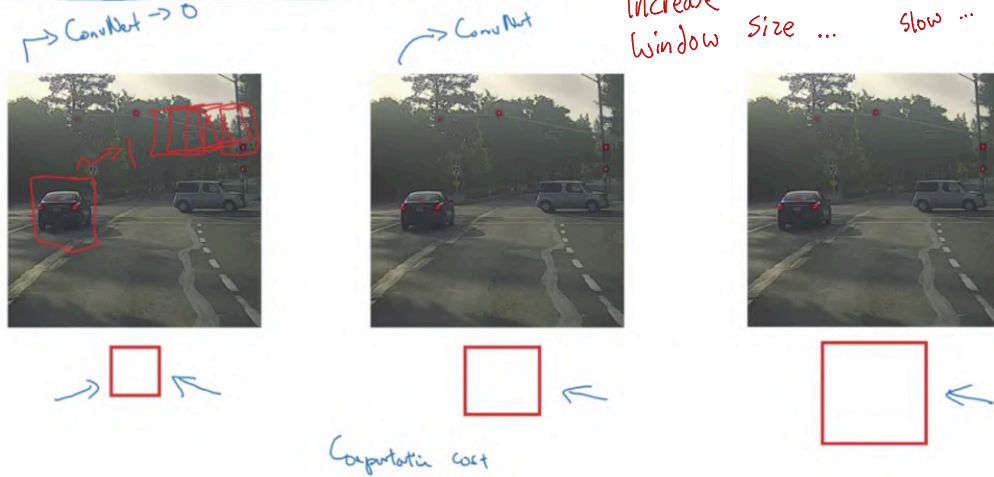
→ ConvNet → 0



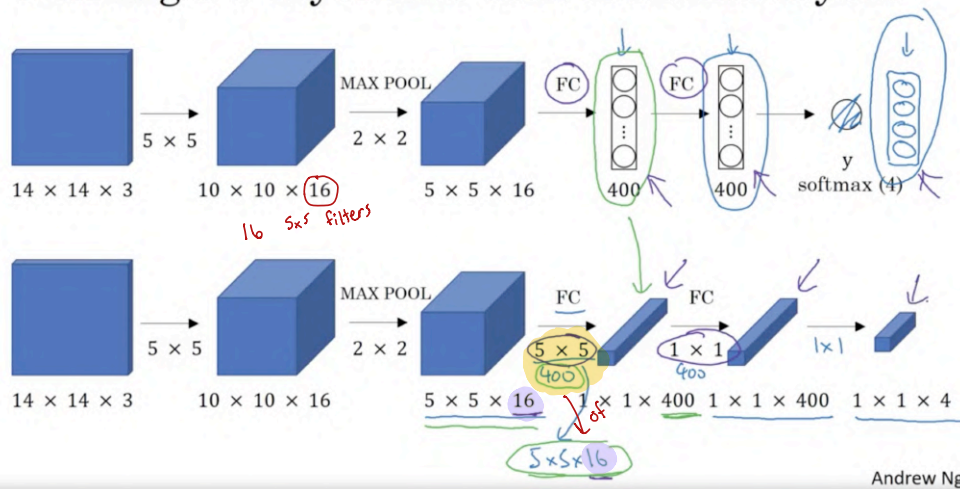
→ ConvNet



fixed
with
counets

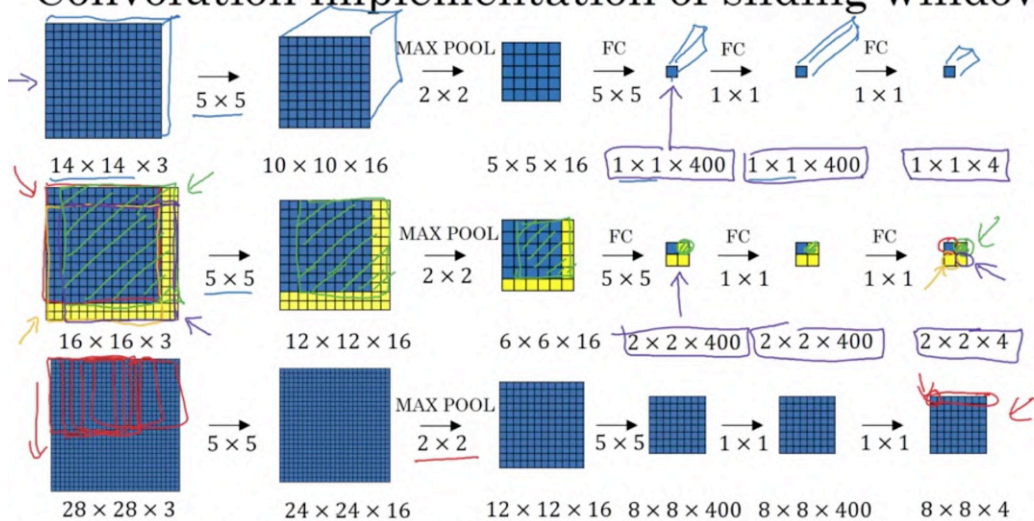


Turning FC layer into convolutional layers



Andrew Ng

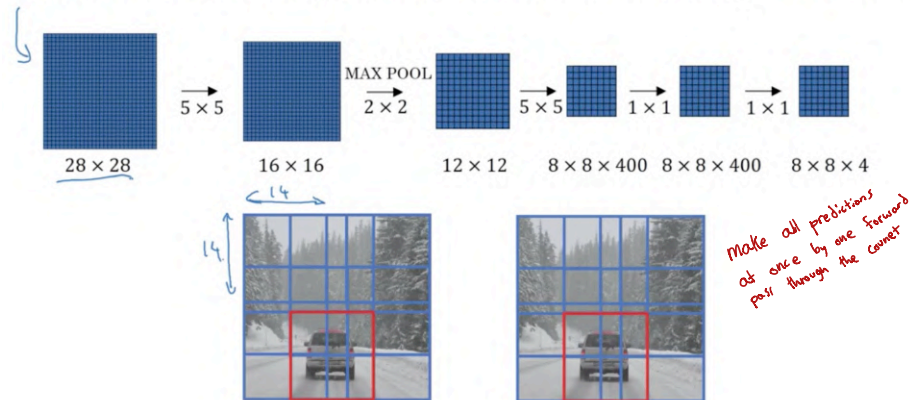
Convolution implementation of sliding windows



[Sermanet et al., 2014, OverFeat: Integrated recognition, localization and detection using convolutional networks]

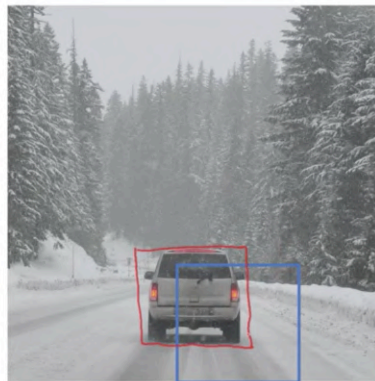
Andrew Ng

Convolution implementation of sliding windows



Andrew Ng

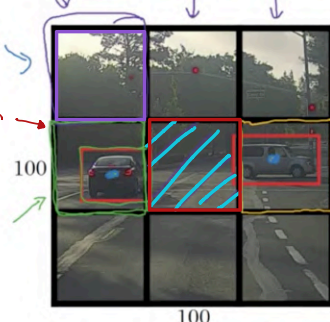
Output accurate bounding boxes



Andrew Ng

YOLO algorithm

Apply Image
Classification and
localization to each
of the 9 grids



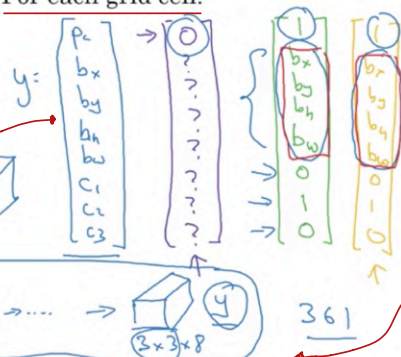
This image has two objects. And what the YOLO algorithm does is it takes the midpoint of each of the two objects and then assigns the object to the grid cell containing the midpoint. So the left car is assigned to this grid cell, and the car on the right, which is this midpoint, is assigned to this grid cell. And so even though the central grid cell has some parts of both cars, we'll pretend the central grid cell has no interesting object so that the central grid cell the class label Y also looks like this vector with no object, and so the first component P_C , and then the rest are don't care. Whereas for this cell, this cell that I have circled in green on the left, the target label Y would be as follows.

Target output:
 $3 \times 3 \times 8$

Grid

Labels for training

For each grid cell:

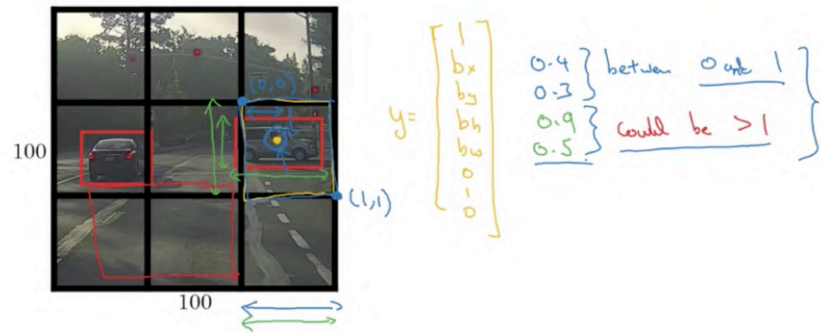


[Redmon et al., 2015, You Only Look Once: Unified real-time object detection]

Andrew Ng

Each object only assigned to 1

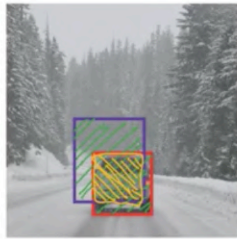
Specify the bounding boxes



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

Andrew Ng

Evaluating object localization



Intersection over Union (IoU)

$$= \frac{\text{size of } \text{yellow box}}{\text{size of } \text{green box}}$$

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

0.6 ← rarely below 0.5

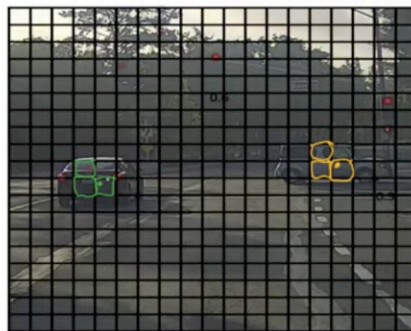
More generally, IoU is a measure of the overlap between two bounding boxes.

Andrew Ng

Non-max suppression example

Since you're running the image classification and localization algorithm on every grid cell, on 361 grid cells, it's possible that many of them will raise their hand and say, "My Pc, my chance of thinking I have an object in it is large." Rather than just having two of the grid cells out of the 19 squared or 361 think they have detected an object. So, when you run your algorithm, you might end up with multiple detections of each object. So, what non-max suppression does, is it cleans up these detections. So they end up with just one detection per car, rather than multiple detections per car. So concretely, what it does, is it first looks at the probabilities associated with each of these detections.

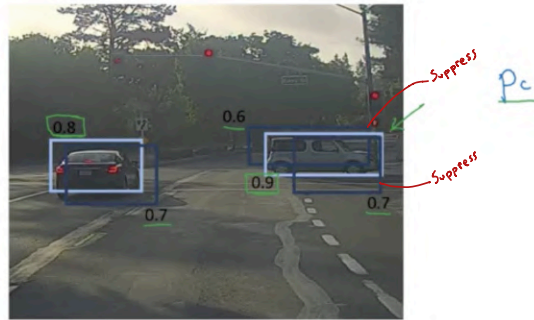
let's just say is Pc with the probability of a detection. And it first takes the largest one, which in this case is 0.9 and says, "That's my most confident detection, so let's highlight that and just say I found the car there." Having done that the non-max suppression part then looks at all of the remaining rectangles and all the ones with a high overlap, with a high IOU, with this one that you've just output will get suppressed.



19x19

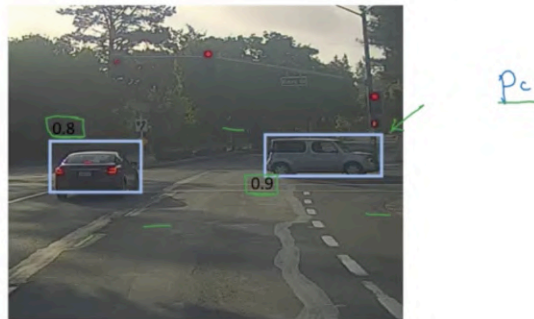
Andrew Ng

Non-max suppression example



Andrew Ng

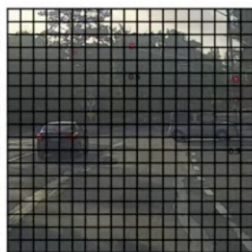
Non-max suppression example



Andrew Ng

Non-max suppression algorithm

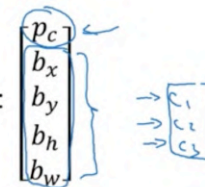
If you actually tried to detect three objects say pedestrians, cars, and motorcycles, then the output vector will have three additional components. And it turns out, the right thing to do is to independently carry out non-max suppression three times, one on each of the outputs classes.



19x19

Assume only car detection for this example

Each output prediction is:



Discard all boxes with $p_c \leq 0.6$

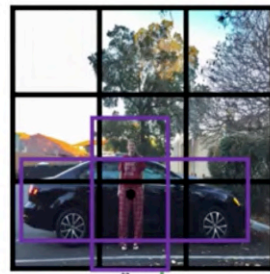
→ While there are any remaining boxes:

- Pick the box with the largest p_c . Output that as a prediction.
- Discard any remaining box with $\text{IoU} \geq 0.5$ with the box output in the previous step

Andrew Ng

Overlapping objects:

With the idea of anchor boxes, what you are going to do, is pre-define two different shapes called, anchor boxes or anchor box shapes. And what you are going to do is now, be able to associate two predictions with the two anchor boxes. And in general, you might use more anchor boxes,

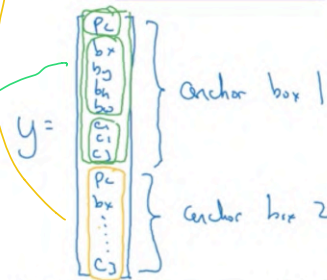


$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

Anchor box 1:



Anchor box 2:



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

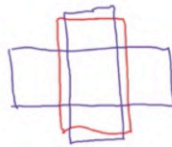
Andrew Ng

Anchor box algorithm

Previously:

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

Output y :
 $3 \times 3 \times 8$



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.

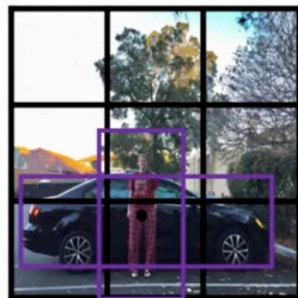
Output y :
 $3 \times 3 \times 16$
 $3 \times 3 \times 2 \times 8$

(grid cell, anchor box)

of anchor boxes

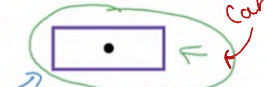
Andrew Ng

Anchor box example



Anchor box 1:

Anchor box 2:



Pedestrian

$y =$

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

$$\begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 0 \\ 0 \\ 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

car only?

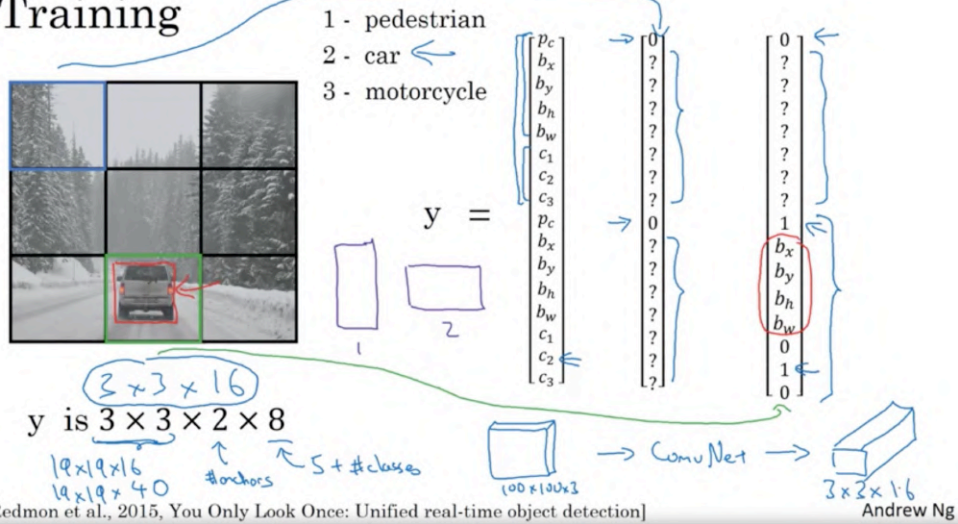
$$\begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

don't need for pedestrian
No object

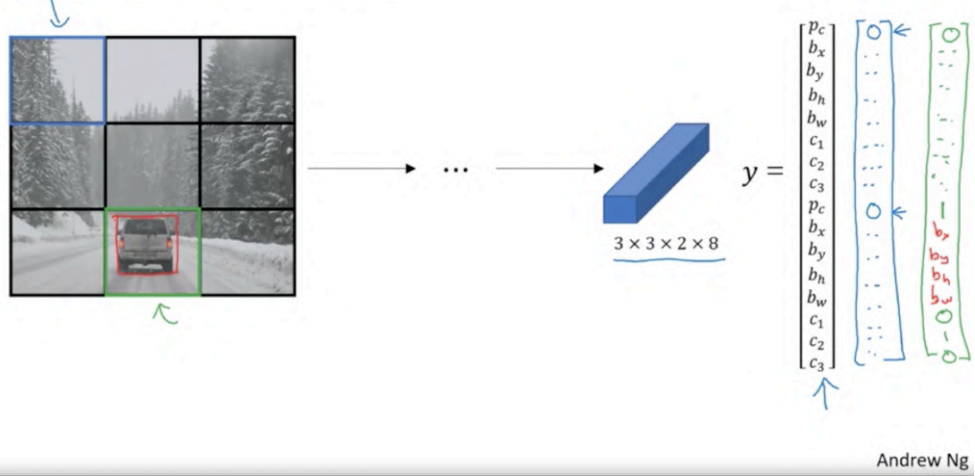
Andrew Ng

Training

What if you have two anchor boxes but three objects in the same grid cell? That's one case that this algorithm doesn't handle well. Hopefully, it won't happen. But if it does, this algorithm doesn't have a great way of handling it. I will just influence some default tiebreaker for that case. Or what if you have two objects associated with the same grid cell, but both of them have the same anchor box shape? Again, that's another case that this algorithm doesn't handle well. If you influence some default way of tiebreaking if that happens, hopefully this won't happen with your data set, it won't happen much at all. And so, it shouldn't affect performance as much.



Making predictions



Outputting the non-max suppressed outputs



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

Outputting the non-max suppressed outputs



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

Andrew Ng

Outputting the non-max suppressed outputs

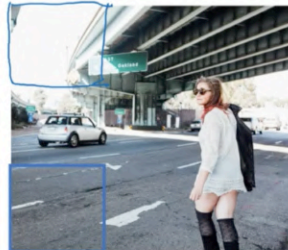


- For each grid cell, get 2 predicted bounding boxes.
- Get rid of low probability predictions.
- For each class (pedestrian, car, motorcycle) use non-max suppression to generate final predictions.

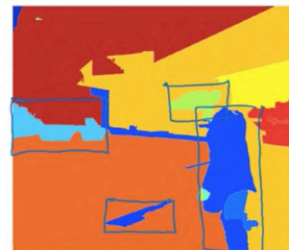
- 3 times

Andrew Ng

Region proposal: R-CNN



↖



*Segmentation algorithm
~ 2,000*

Faster algorithms

→ R-CNN: Propose regions. Classify proposed regions one at a time. Output label + bounding box. ←

Fast R-CNN: Propose regions. Use convolution implementation of sliding windows to classify all the proposed regions. ←

Faster R-CNN: Use convolutional network to propose regions.

[Girshik et. al, 2013. Rich feature hierarchies for accurate object detection and semantic segmentation]

[Girshik, 2015. Fast R-CNN]

[Ren et. al, 2016. Faster R-CNN: Towards real-time object detection with region proposal networks] Andrew Ng