



intro



deeplearning.ai

Object Detection

Object localization

An area that is just exploding

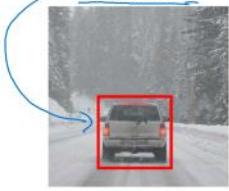
What are localization and detection?

Image classification



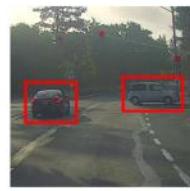
"Car"

Classification with localization



"Car"

Detection



multiple objects

Andrew Ng

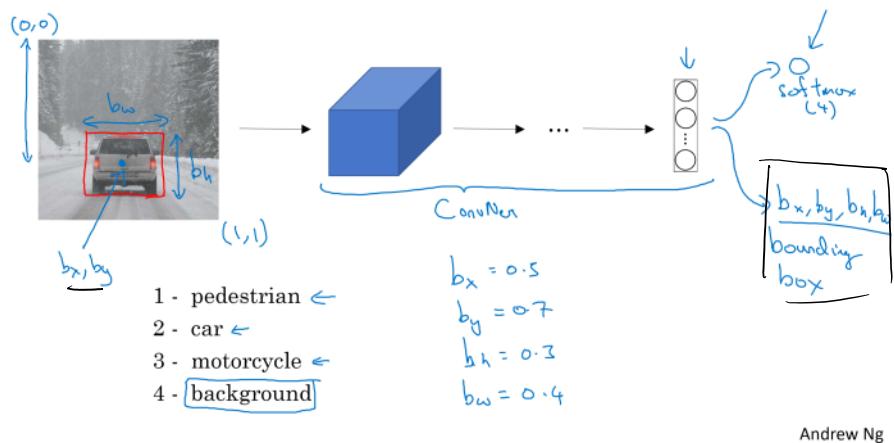
multiple objects
detect and localize

for auto-driving task
detect different types of object

classification with localization

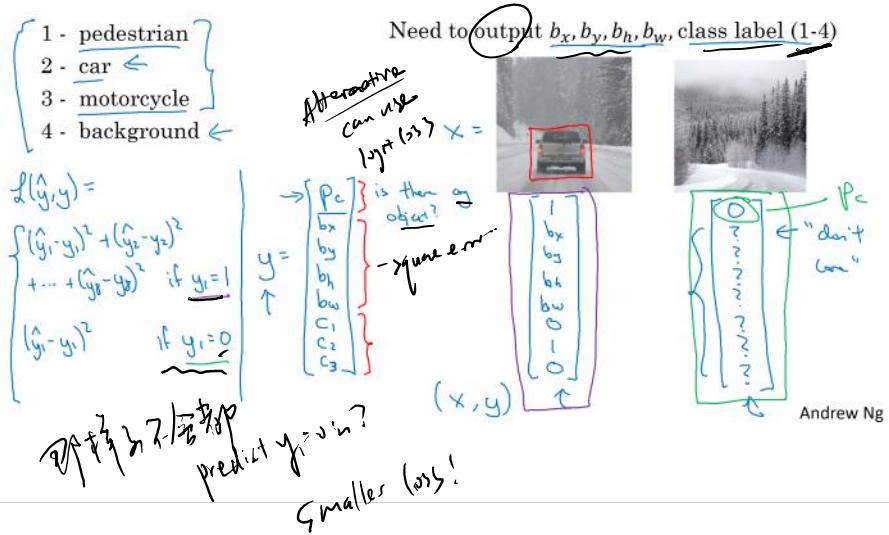
putting a bounding box around the position of the object

Classification with localization



b_x, b_y, b_h, b_w | output
 training set (class label)
 include the object type
 type and the variables
 of location.

Defining the target label y



這樣型號的標記方法
 — Loss function 計算的問題
 — $\hat{y}_c = 0$ 是 wrong
 — BP predict \hat{y}_c 不對。



deeplearning.ai

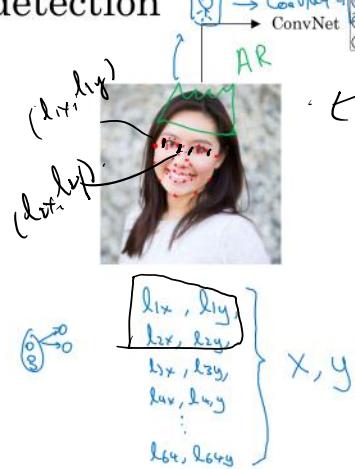
Object Detection

Landmark detection

Landmark detection



b_x, b_y, b_h, b_w



$l_{1x}, l_{1y}, l_{2x}, l_{2y}, \dots, l_{32x}, l_{32y}$

Andrew Ng

use: recognize facial expressions etc

- export important points of the image
- want: where's the corner in somewhere eyes.

can extract many points

post detection
e.g. midpoint of chest --

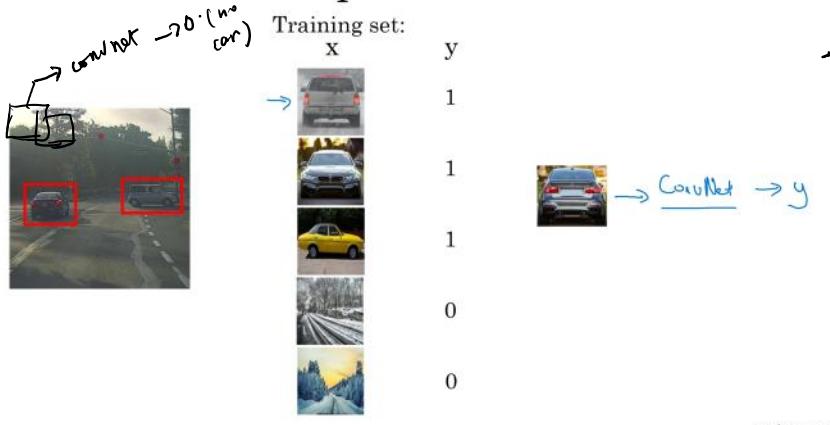


deeplearning.ai

Object Detection

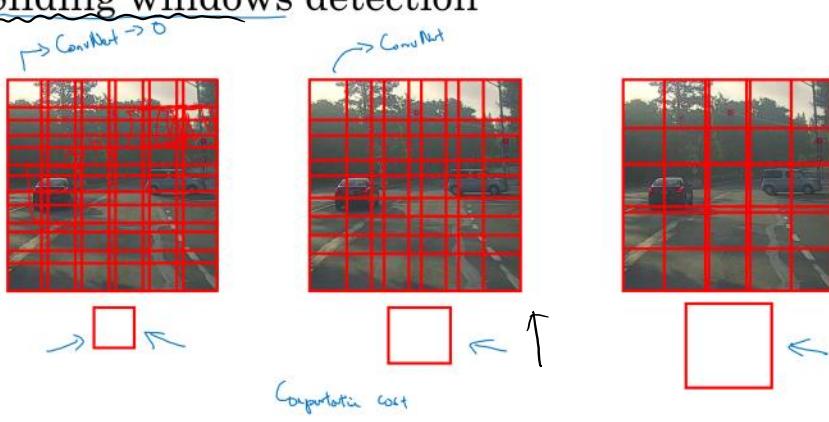
Object detection

Car detection example



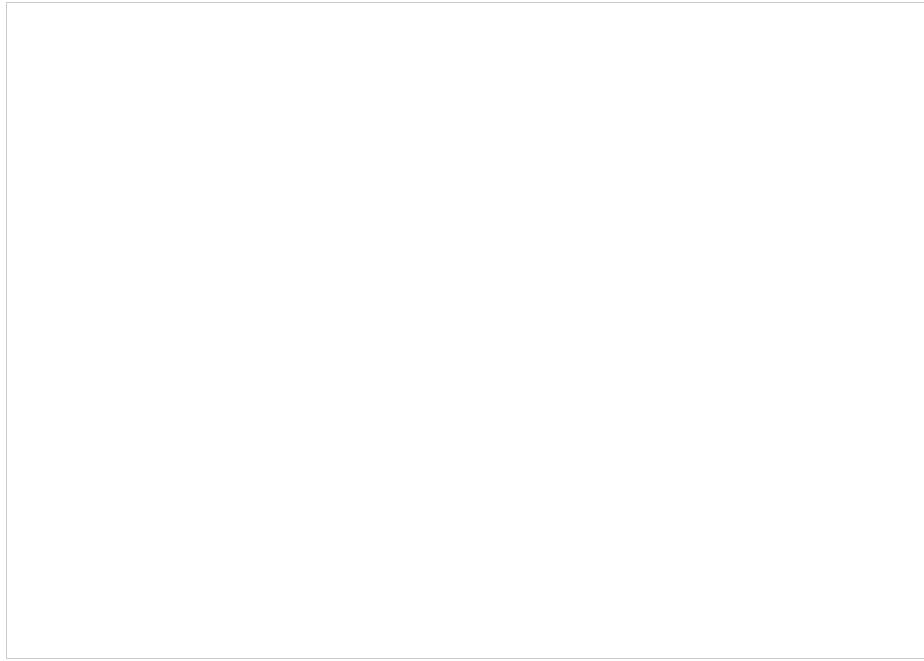
sliding windows detection algorithm
— small region
— pass many crop # of the image.
→
large region (choose a larger window)

Sliding windows detection



large computational cost!
using sliding window detection.
(before: simple classifier + sliding window)
Now: NN + sliding → unbearable comput
Next: convolutional sliding window/
solve the
window size ↑

conv sliding windows



Turning FC layer into convolutional layers



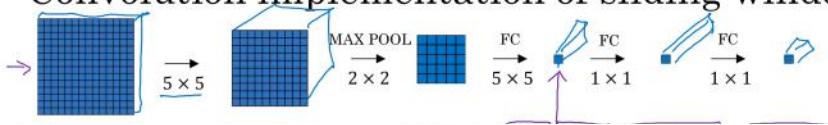
Mathematically, the same as fully connected softmax activation.

↙ output classes

filter then output
 $1 \times 1 \times n_c$

逐個把第3維拿掉

Convolution implementation of sliding windows



running on the
result as top left (4×4)

running on
 same result as top left $14 \times 14 \times 3$
 share parameters
 share computation
 running sliding window on top-left $14 \times 14 \times 3$

conv net. To 7x7
 true & false -> run

prev sequentially.

make prediction for all regions simultaneously

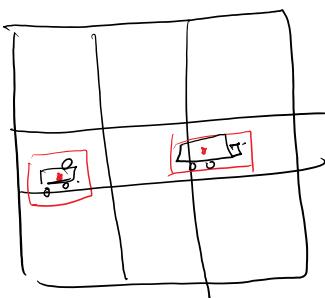
$$\begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

problem: not output the accurate bounding boxes.

Bounding box precision

YOLO Algorithm

You only looks once



Assign the car to the grid cell containing the

- Divide image into 9 grid cells
- Labels for training for each grid cell

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_w \\ b_h \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 1 \\ 2 \\ \vdots \\ 7 \end{bmatrix} \quad \begin{bmatrix} 1 \\ b_w \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

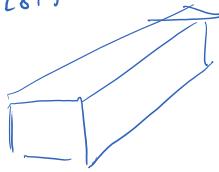
meaning of these dimensions
 See the first few slides.
 they label the location, type of object

Assign the car to the grid cell containing the mid point.

$$\begin{bmatrix} bw \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

(8,)

output



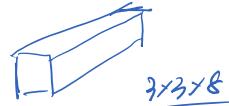
$\frac{3 \times 3 \times 8}{3 \times \text{grid}}$ length of y



→ there an object? Where it is?



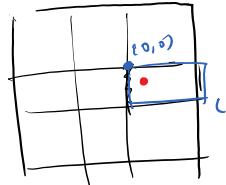
→ conv → max pool ... →



As long as ≤ 1 object each grid cell, this algorithm should work fine.

What's special You are Not implementing this algorithm q times for each grid.
You share computation among grids — efficient

one more detail, How to code the bounding box?



$$\begin{bmatrix} 1 \\ 0.4 \\ 0.3 \\ 0.9 \\ 0.5 \\ 0 \\ 0.1 \end{bmatrix}$$

there are some other parameterization

the location is specified relative to
the grid cell cell 11's in.

(Hc, Tg, ~~W~~ \Rightarrow $\frac{1}{2}, \frac{1}{2}, 1_w$).

The YOLO paper is hard to read ... [warning]

intersection union

(Evaluate) Object detection model.



deeplearning.ai

Object Detection

Intersection

(Evaluate) Object detection
model.

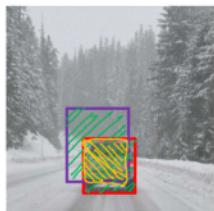


deeplearning.ai

Object Detection

Intersection over union

Evaluating object localization



Intersection over Union (IoU)

$$= \frac{\text{Size of } \cap}{\text{Size of } \cup}$$

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

$\xleftarrow{\text{if overlap perfectly}} \text{IoU} = 1$

~~I predict for 3 for background~~
~~1/3~~

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

Andrew Ng



non max suppression



deeplearning.ai

Object Detection

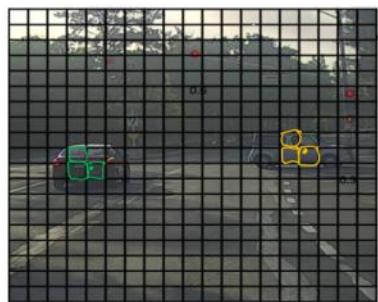
Non-max suppression

Non-max suppression example



Andrew Ng

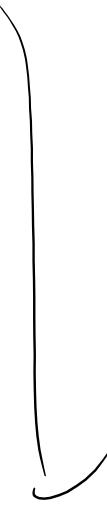
Non-max suppression example



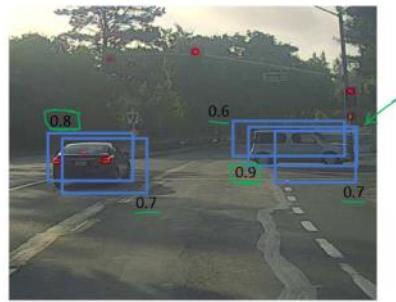
19x19

The ~~one~~ car's size span through grids

Andrew Ng



Non-max suppression example

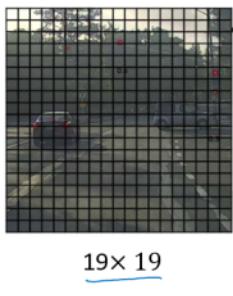


multiple detection of ~~one~~ car (object)

Pc
takes

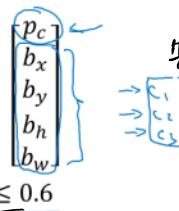
Andrew Ng

Non-max suppression algorithm



19 × 19

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

object detection -

if ≥ 3 classes.

carry out non-max suppression
 \geq times

What's going on?

Why? Get

anchor boxes



deeplearning.ai

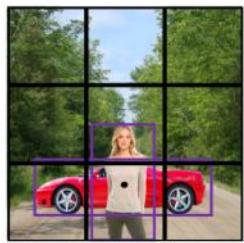
Object Detection

Anchor boxes

Solve problem

multiple objects one call

Overlapping objects:



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

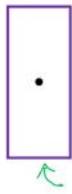
[Redmon et al., 2015, You Only Look Once: Unified real-time object detection]

?

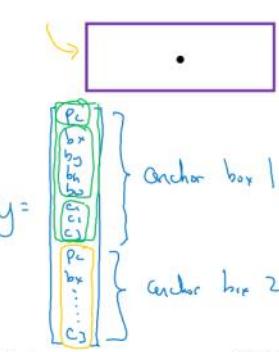
pre-defined two boxes

\Rightarrow associate the prediction with two boxes

Anchor box 1:



Anchor box 2:



pre-defined shape? How?

after getting bx, by, bw, bh

Th. ah. is it for two anchor box

\Rightarrow Th. of Th. anchor box

from label 2(?) is 2(?)

or n(?)

Andrew Ng

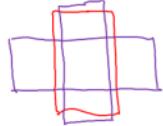
Anchor box algorithm

Previously:

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

Output y:

$$\underline{2 \times 2 \times 8}$$



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.

(grid cell, anchor box)

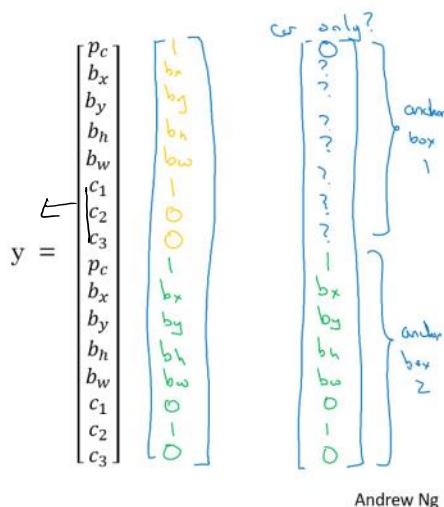
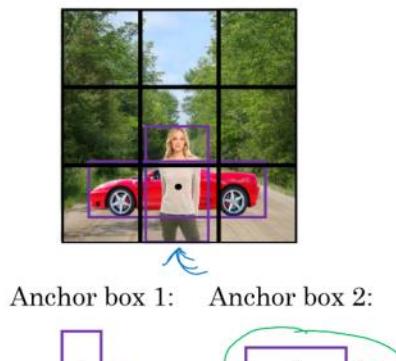
Output y:

$$\underline{3 \times 3 \times 16}$$

$$\underline{3 \times 3 \times 2 \times 8}$$

Andrew Ng

Anchor box example



每种东西 -> 4类型
anchor box

other problem
same type
two objects in a grid that have the same anchor box
→ algorithm to fix it

Anchor box Allow algorithm to specialize better.

7 How to choose the anchor boxes? by hand

Advanced: use the k-means algorithm to select Anchor box (?)
later YOLO papers ...

yolo

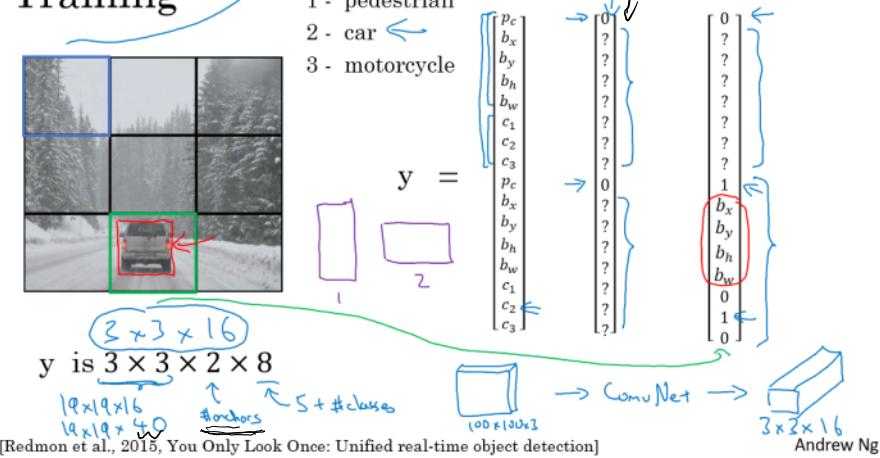


deeplearning.ai

Object Detection

Putting it together: YOLO algorithm

Training



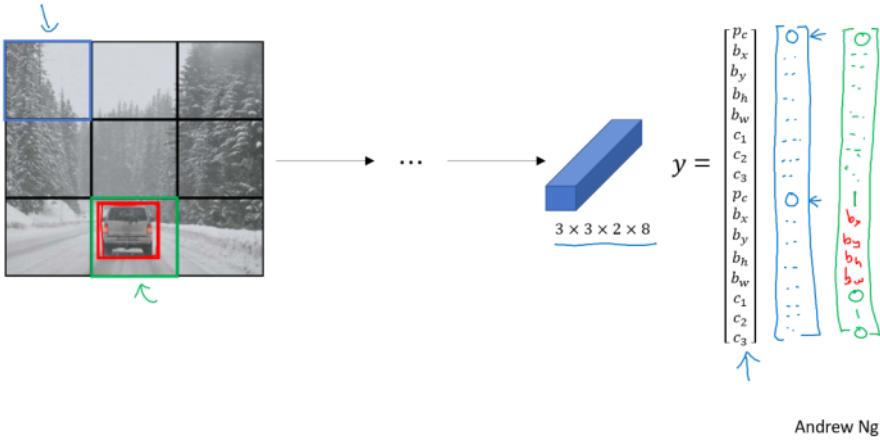
$c_1, c_2, c_3 \rightarrow$ 哪个最好?

anchor box 在 $c_1, c_2, c_3 \rightarrow$ 哪一个最好?

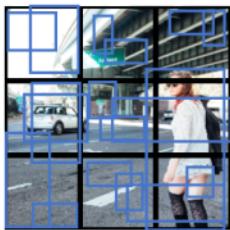
哪个最好呢? 我们怎么知道?



Making predictions



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.

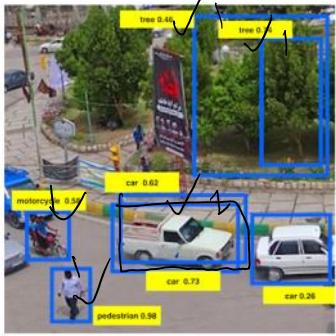
↓
Some bounding boxes
Can go outside the grid — ? how
width, height of grid ~~is 19x19~~.
Andrew Ng ~~19x19~~ ~~19x19~~ ~~19x19~~ ~~19x19~~ ~~19x19~~
like 70%

- ✗ 7. In the YOLO algorithm, at training time, only one cell ---the one containing the center/midpoint of an object--- is responsible for detecting this object.

0 / 1 points

Yes (I don't understand what this means).

- ✗ 9. Suppose you run non-max suppression on the predicted boxes above. The parameters you use for non-max suppression are that boxes with probability ≤ 0.4 are discarded, and the IoU threshold for deciding if two boxes overlap is 0.5. How many boxes will remain after non-max suppression?



- ✗ 10. Suppose you are using YOLO on a 19×19 grid, on a detection problem with 20 classes, and with 5 anchor boxes. During training, for each image you will need to construct an output volume y as the target value for the neural network; this corresponds to the last layer of the neural network. (y may include some "0"s, or "don't cares"). What is the dimension of this output volume?

0 / 1 points

$$19 \times 19 \times 5 + 25 = 1$$

5 + num of classes