deeplearning.ai

# Object Detection
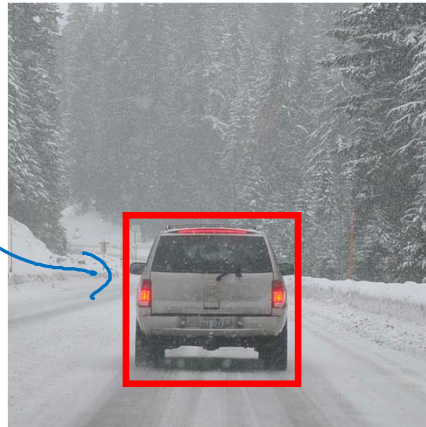
---

# Object localization

# What are localization and detection?

Image classification



"Car"

Classification with localization
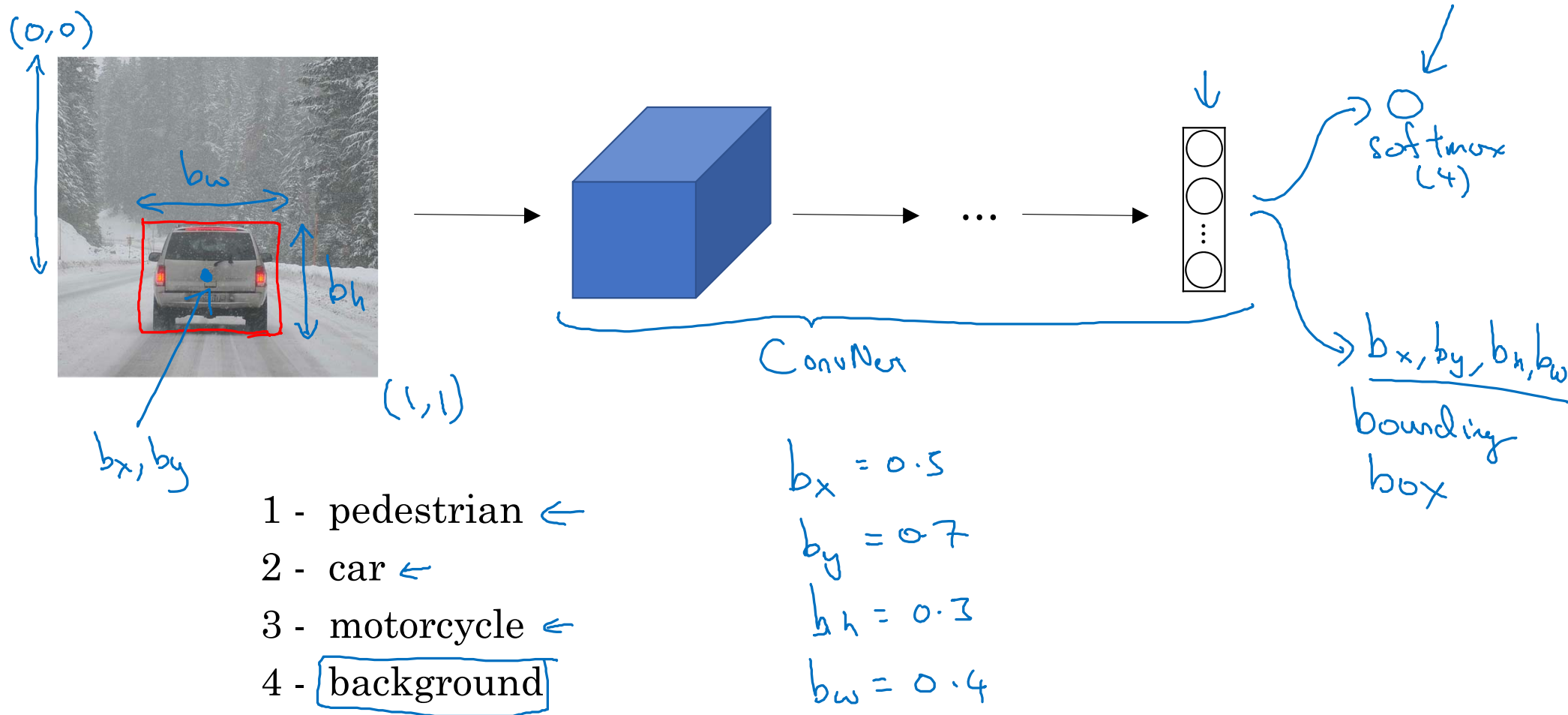


"Car"

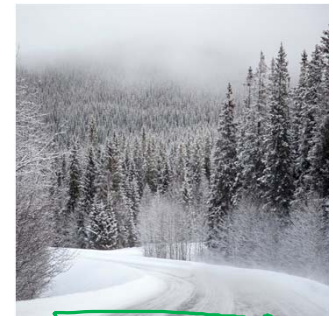Detection



1 object

multiple objects
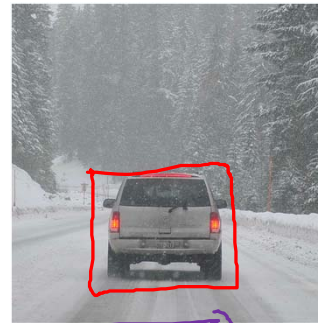
Andrew Ng

# Classification with localization



(0,0)

$b_w$

$b_h$

$b_x, b_y$

(1,1)

1 - pedestrian ←
2 - car ←
3 - motorcycle ←
4 - background

ConvNet

softmax (4)

$b_x, b_y, b_h, b_w$
bounding box

$b_x = 0.5$
$b_y = 0.7$
$b_h = 0.3$
$b_w = 0.4$

Andrew Ng

# Defining the target label y

1 - pedestrian
2 - car
3 - motorcycle
4 - background

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

$x =$

$\mathcal{L}(\hat{y}, y) =$

$\begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 \\ + \cdots + (\hat{y}_8 - y_8)^2 \quad \text{if } y_1 = 1 \\ \\ (\hat{y}_1 - y_1)^2 \quad \text{if } y_1 = 0 \end{cases}$

$y = \begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$ is there any object?

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

$(x, y)$

$\begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$ ← "don't care"
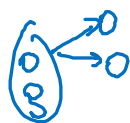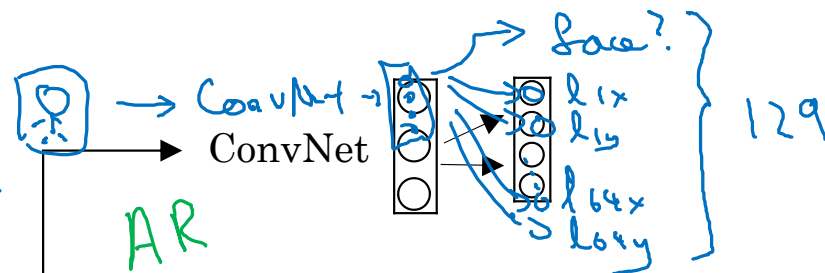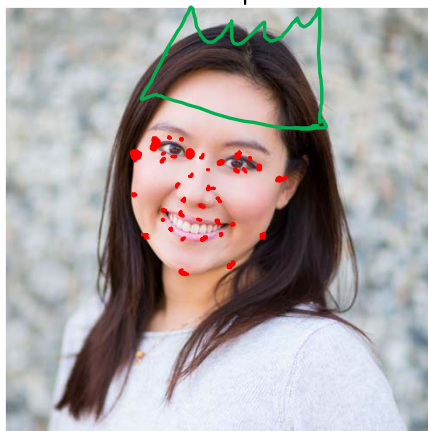
$P_c$

deeplearning.ai

# Object Detection

---

# Landmark detection

# Landmark detection



face?

ConvNet → $l_{1x}$ $l_{1y}$ $l_{64x}$ $l_{64y}$ } 129

AR

$b_x, b_y, b_h, b_w$

$l_{1x}, l_{1y},$
$l_{2x}, l_{2y},$
$l_{3x}, l_{3y},$
$l_{4x}, l_{4y},$
$\vdots$
$l_{64}, l_{64y}$ } $X, Y$

$l_{1x}, l_{1y},$
$\vdots$
$l_{32x} l_{22y}$

Andrew Ng

Object Detection

deeplearning.ai

Object detection

# Car detection example

Training set:

# Sliding windows detection



$\rightarrow$ ConvNet $\rightarrow$ 0

$\rightarrow$ ConvNet

Computation cost

Andrew Ng

Object Detection

Convolutional implementation of sliding windows

deeplearning.ai

# Turning FC layer into convolutional layers



14 × 14 × 3     5 × 5     10 × 10 × 16     MAX POOL     2 × 2     5 × 5 × 16     FC     400     FC     400     y softmax (4)

14 × 14 × 3     5 × 5     10 × 10 × 16     MAX POOL     2 × 2     5 × 5 × 16     FC   5 × 5   400     1 × 1 × 400     FC   1 × 1   400     1 × 1 × 400     1×1     1 × 1 × 4

5×5×16

Andrew Ng

# Convolution implementation of sliding windows

$14 \times 14 \times 3$ $\xrightarrow{5 \times 5}$ $10 \times 10 \times 16$ $\xrightarrow[2 \times 2]{\text{MAX POOL}}$ $5 \times 5 \times 16$ $\xrightarrow[5 \times 5]{\text{FC}}$ $1 \times 1 \times 400$ $\xrightarrow[1 \times 1]{\text{FC}}$ $1 \times 1 \times 400$ $\xrightarrow[1 \times 1]{\text{FC}}$ $1 \times 1 \times 4$

$16 \times 16 \times 3$ $\xrightarrow{5 \times 5}$ $12 \times 12 \times 16$ $\xrightarrow[2 \times 2]{\text{MAX POOL}}$ $6 \times 6 \times 16$ $\xrightarrow[5 \times 5]{\text{FC}}$ $2 \times 2 \times 400$ $\xrightarrow[1 \times 1]{\text{FC}}$ $2 \times 2 \times 400$ $\xrightarrow[1 \times 1]{\text{FC}}$ $2 \times 2 \times 4$

$28 \times 28 \times 3$ $\xrightarrow{5 \times 5}$ $24 \times 24 \times 16$ $\xrightarrow[2 \times 2]{\text{MAX POOL}}$ $12 \times 12 \times 16$ $\xrightarrow{5 \times 5}$ $8 \times 8 \times 400$ $\xrightarrow{1 \times 1}$ $8 \times 8 \times 400$ $\xrightarrow{1 \times 1}$ $8 \times 8 \times 4$

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

Andrew Ng

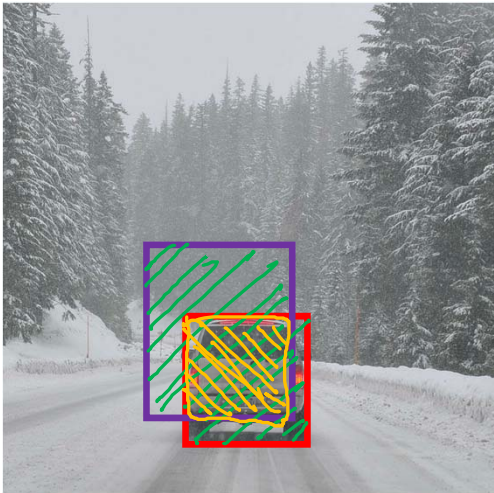# Convolution implementation of sliding windows



28 × 28  →(5 × 5)→  16 × 16  →MAX POOL (2 × 2)→  12 × 12  →(5 × 5)→  8 × 8 × 400  →(1 × 1)→  8 × 8 × 400  →(1 × 1)→  8 × 8 × 4

Andrew Ng

deeplearning.ai

Object Detection
_____

Intersection
over union

# Evaluating object localization



$$\text{Intersection over Union} \quad (\text{IoU})$$

$$= \frac{\text{Size of } \boxed{\text{(yellow)}}}{\text{Size of } \boxed{\text{(green)}}}$$

"Correct" if IoU $\geq$ 0.5

0.6

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

deeplearning.ai

Object Detection

Non-max
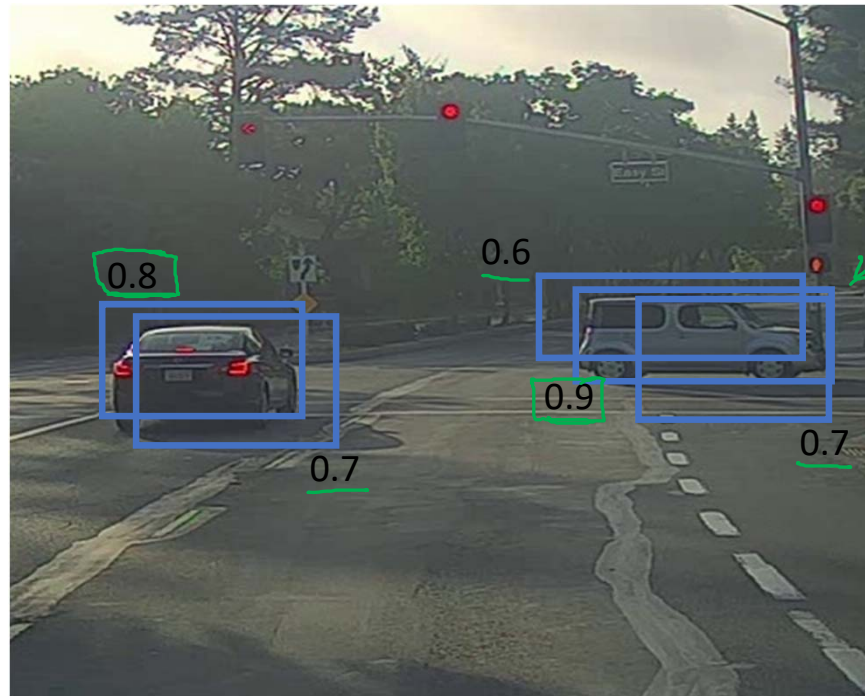suppression
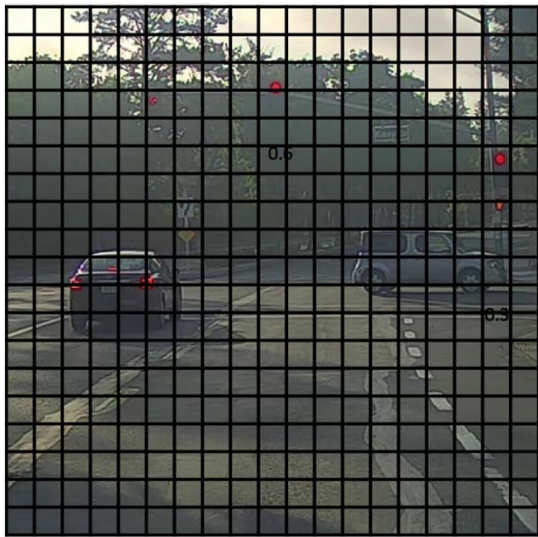
# Non-max suppression example

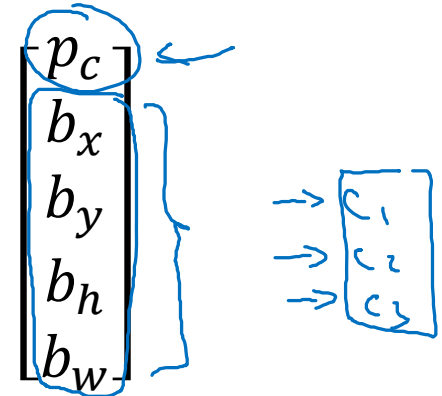# Non-max suppression example



19x19

# Non-max suppression example



Andrew Ng

# Non-max suppression algorithm



19× 19

Each output prediction is:
$$\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \end{bmatrix}$$

$c_1$
$c_2$
$c_3$

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
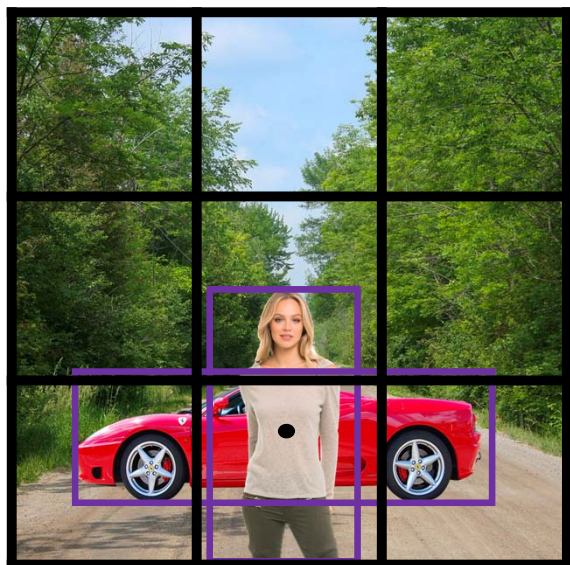  IoU $\geq 0.5$ with the box output
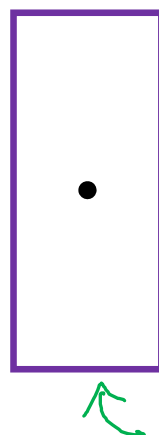  in the previous step

Andrew Ng

Object Detection
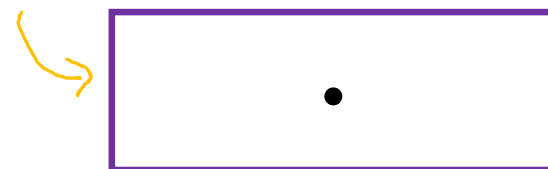
Anchor boxes

deeplearning.ai

# Overlapping objects:



Anchor box 1:

Anchor box 2:

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

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_1 \\ c_3 \\ p_c \\ b_x \\ \vdots \\ c_3 \end{bmatrix} \begin{array}{l} \text{Anchor box 1} \\ \\ \text{Anchor box 2} \end{array}$$

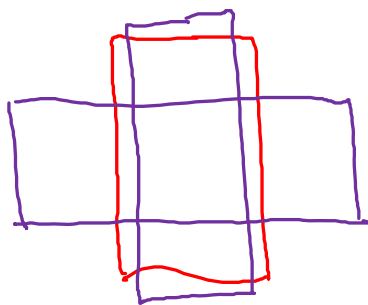[Redmon et al., 2015, You Only Look Once: Unified real-time object detection]

Andrew Ng

# Anchor box algorithm

### Previously:

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

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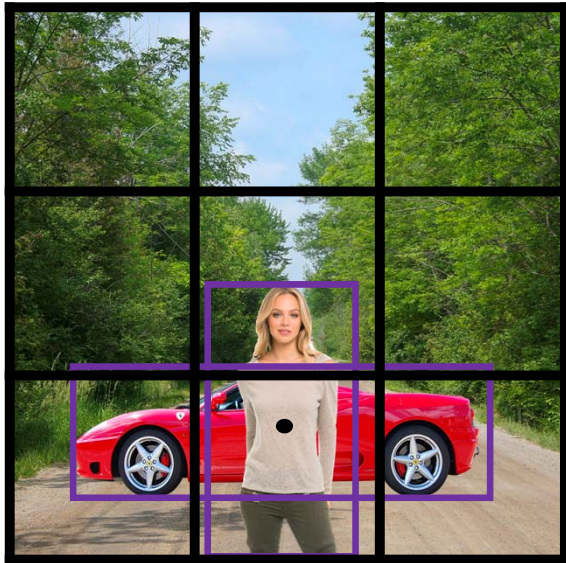
$3 \times 3 \times 8$

### 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.
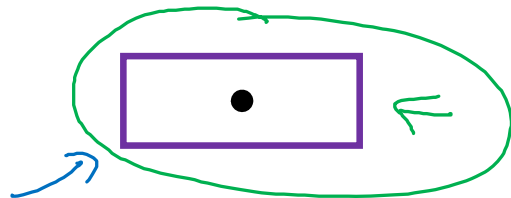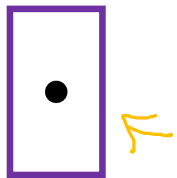
(grid cell, anchor box)

Output y:

$3 \times 3 \times 16$

$3 \times 3 \times 2 \times 8$

Andrew Ng

# Anchor box example



Anchor box 1:    Anchor box 2:

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

Handwritten annotations (orange, anchor box 1): $1$, $b_x$, $b_y$, $b_h$, $b_w$, $0$, $1$, $0$

Handwritten annotations (green, anchor box 2): $1$, $b_x$, $b_y$, $b_h$, $b_w$, $0$, $1$, $0$

Second column (car only?): $0$, ?, ?, ?, ?, ?, ?, ? — anchor box 1

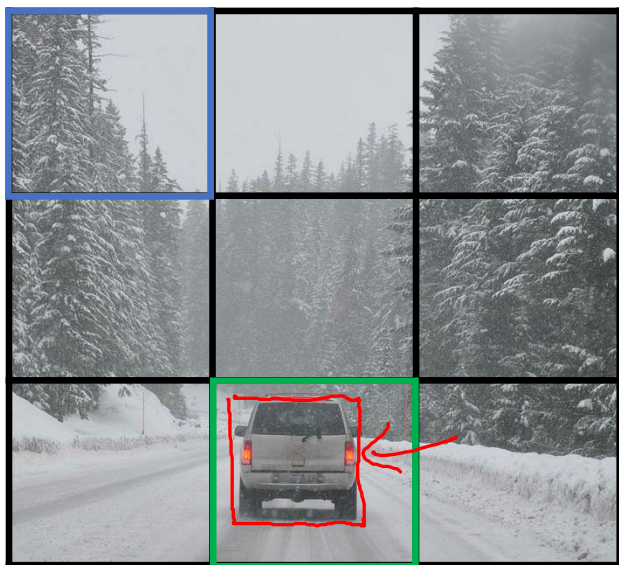green: $1$, $b_x$, $b_y$, $b_h$, $b_w$, $0$, $1$, $0$ — anchor box 2
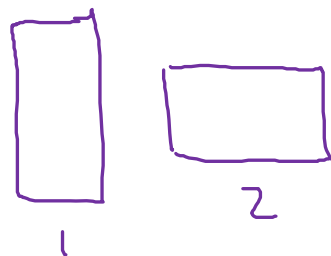
Andrew Ng

deeplearning.ai

# Object Detection

## Putting it together: YOLO algorithm

# Training

1 - pedestrian
2 - car
3 - motorcycle

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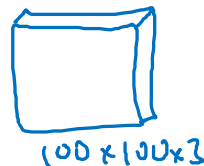\\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix} \quad \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

1
2

$3 \times 3 \times 16$

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

$19 \times 19 \times 16$
$19 \times 19 \times 40$

#anchors
$5 + $#classes

$100 \times 100 \times 3 \rightarrow$ ComvNet $\rightarrow 3 \times 3 \times 16$

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

Andrew Ng

# Making predictions



$$3 \times 3 \times 2 \times 8$$

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

Andrew Ng

# Outputting the non-max supressed 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.
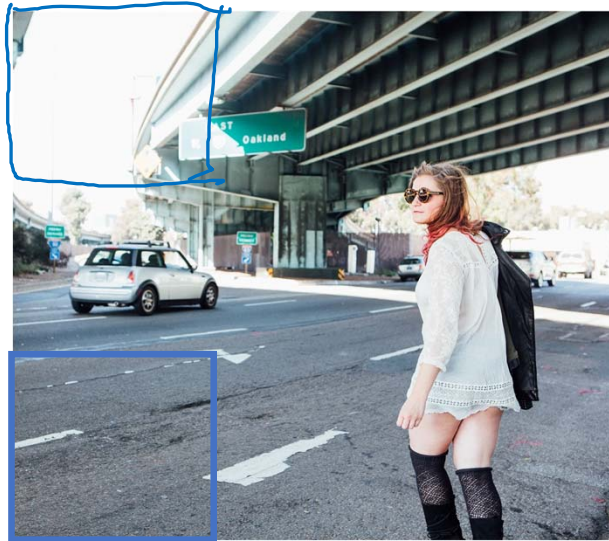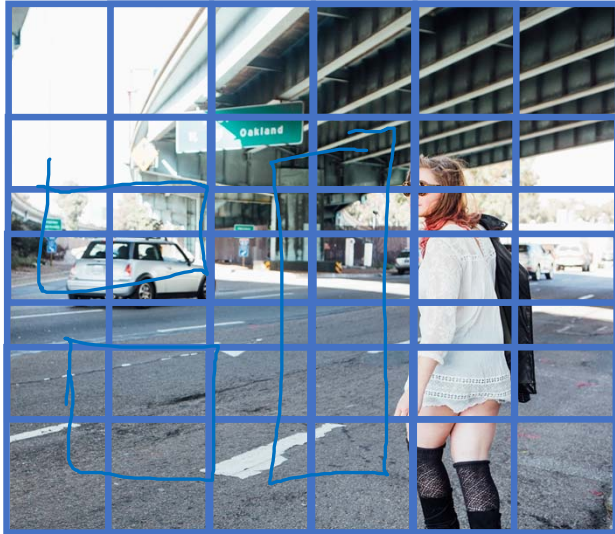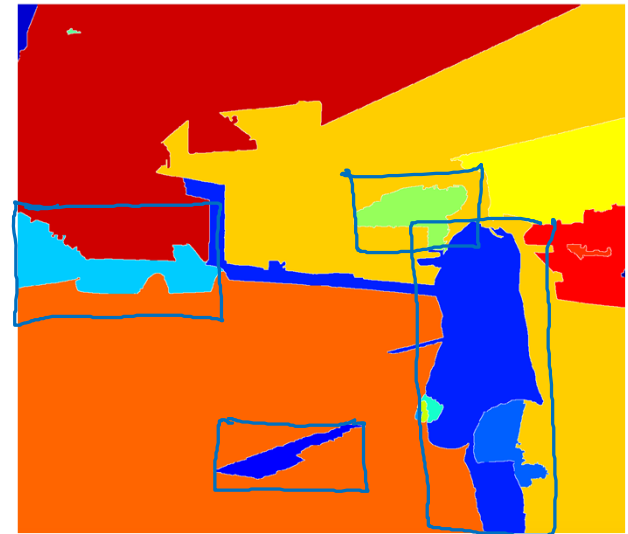
deeplearning.ai

Object Detection

Region proposals
(Optional)

# Region proposal: R-CNN



Segmentation algorithm

~ 2,000

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

# 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