



deeplearning.ai

Object Detection

Object localization

What are localization and detection?

Image classification



"Car"

1 object

Classification with
localization



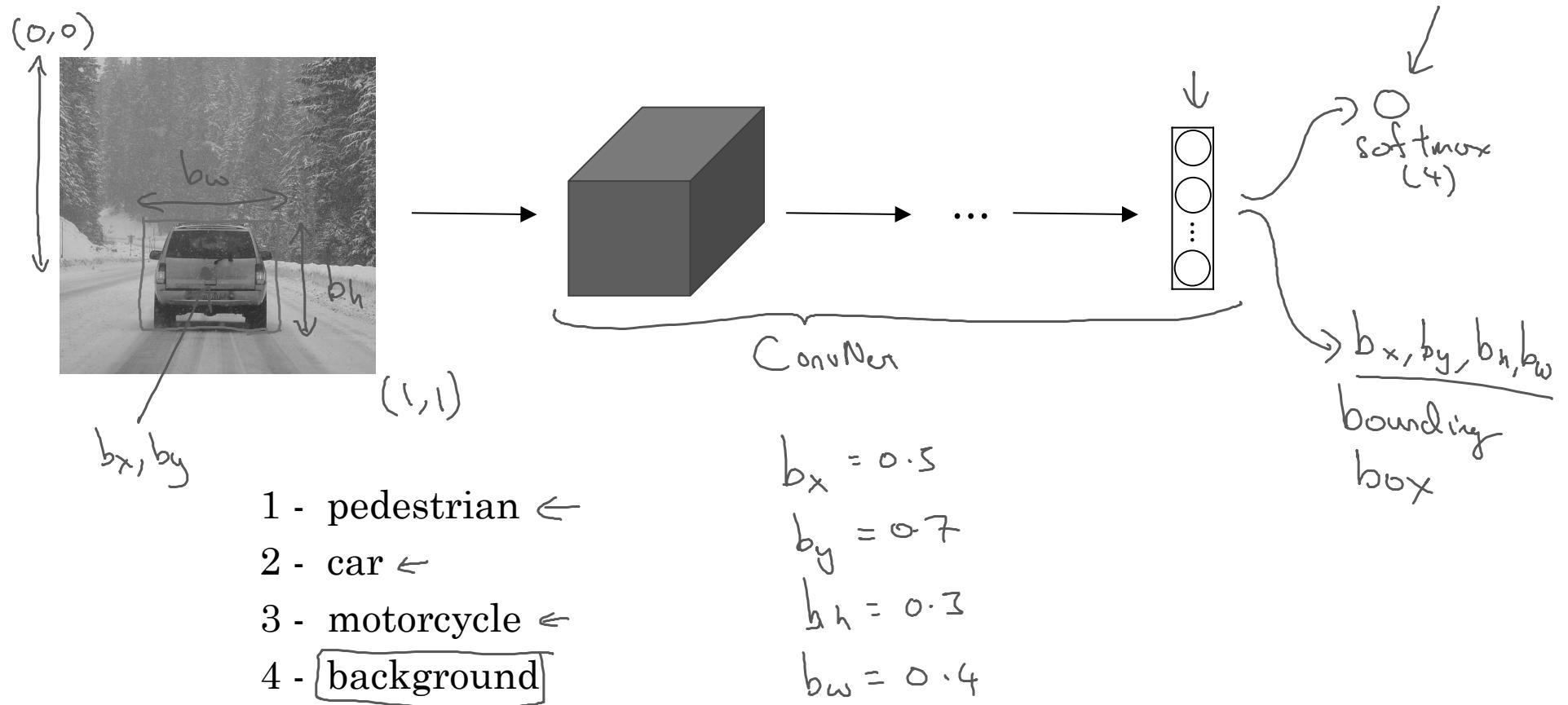
"Car"

Detection



multiple
objects

Classification with localization



Defining the target label y

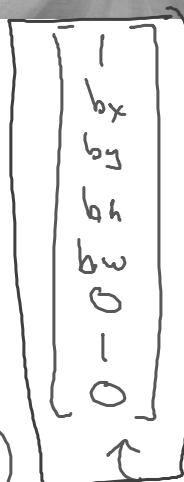
- 1 - pedestrian
- 2 - car \leftarrow
- 3 - motorcycle
- 4 - background \leftarrow

$$L(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 \\ + \dots + (\hat{y}_8 - y_8)^2 & \text{if } \underline{y_1 = 1} \\ (\hat{y}_1 - y_1)^2 & \text{if } \underline{y_1 = 0} \end{cases}$$

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

(x, y)

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





deeplearning.ai

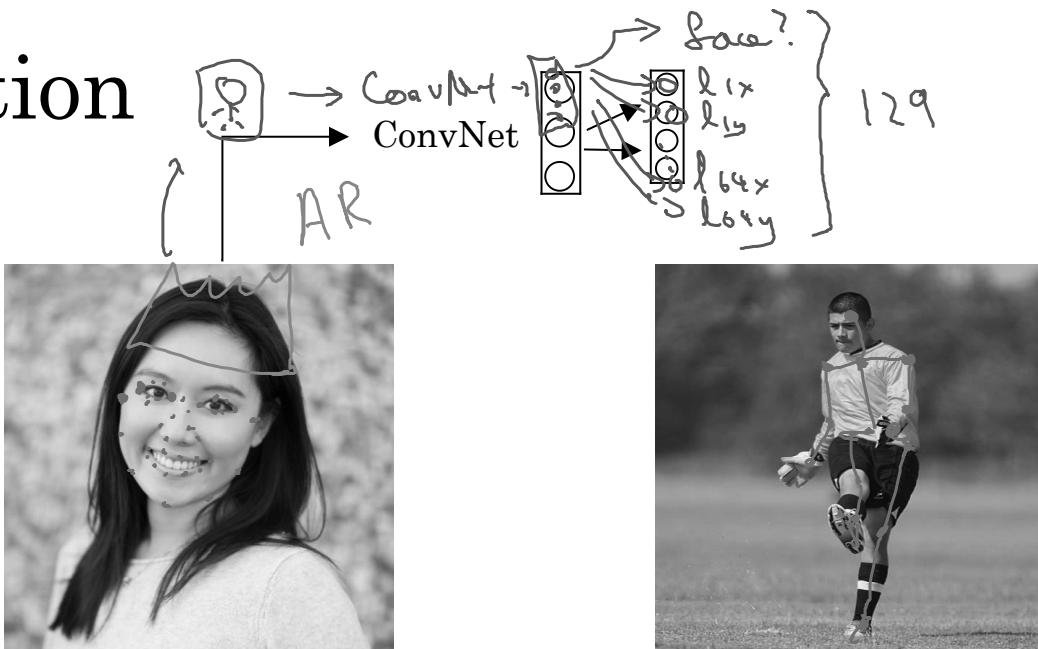
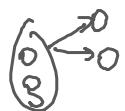
Object Detection

Landmark detection

Landmark detection



b_x, b_y, b_h, b_w



$l_{1x}, l_{1y}, \dots, l_{64x}, l_{64y}$

x, y

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



deeplearning.ai

Object Detection

Object detection

Car detection example

Training set:



x

y

1

1

1

0

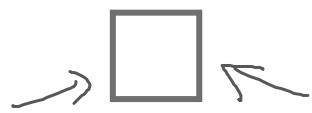
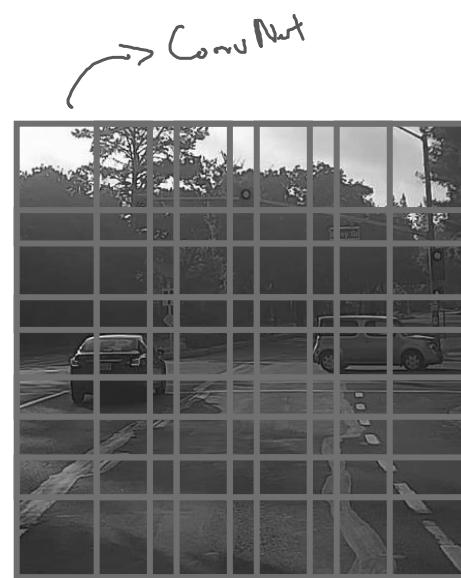
0



ConvNet

$\rightarrow y$

Sliding windows detection



Computation cost

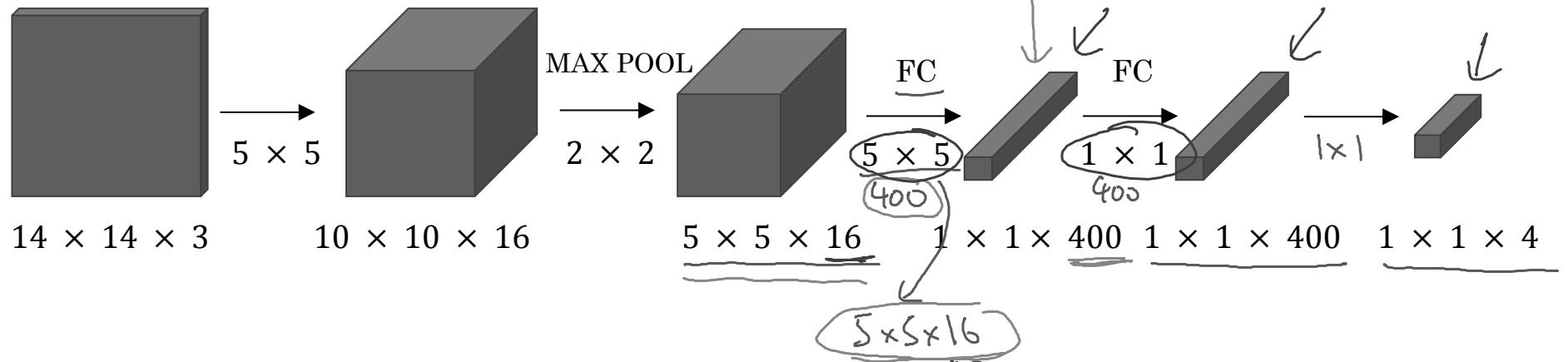
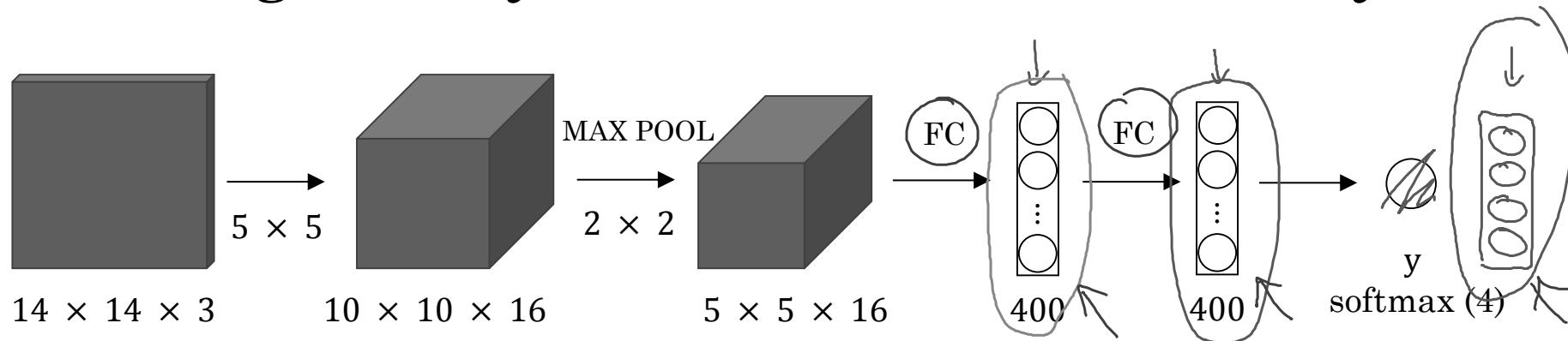


deeplearning.ai

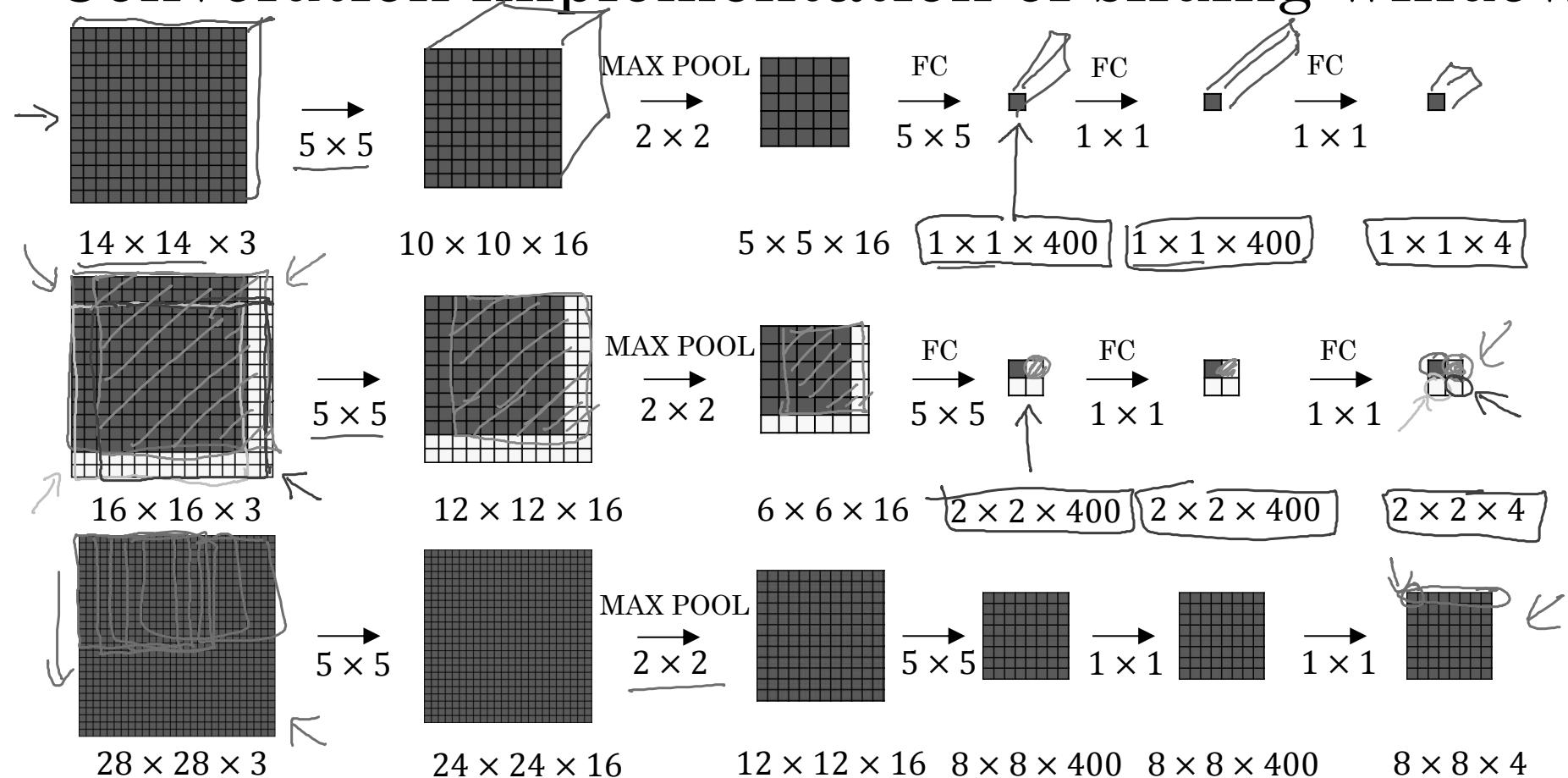
Object Detection

Convolutional implementation of sliding windows

Turning FC layer into convolutional layers

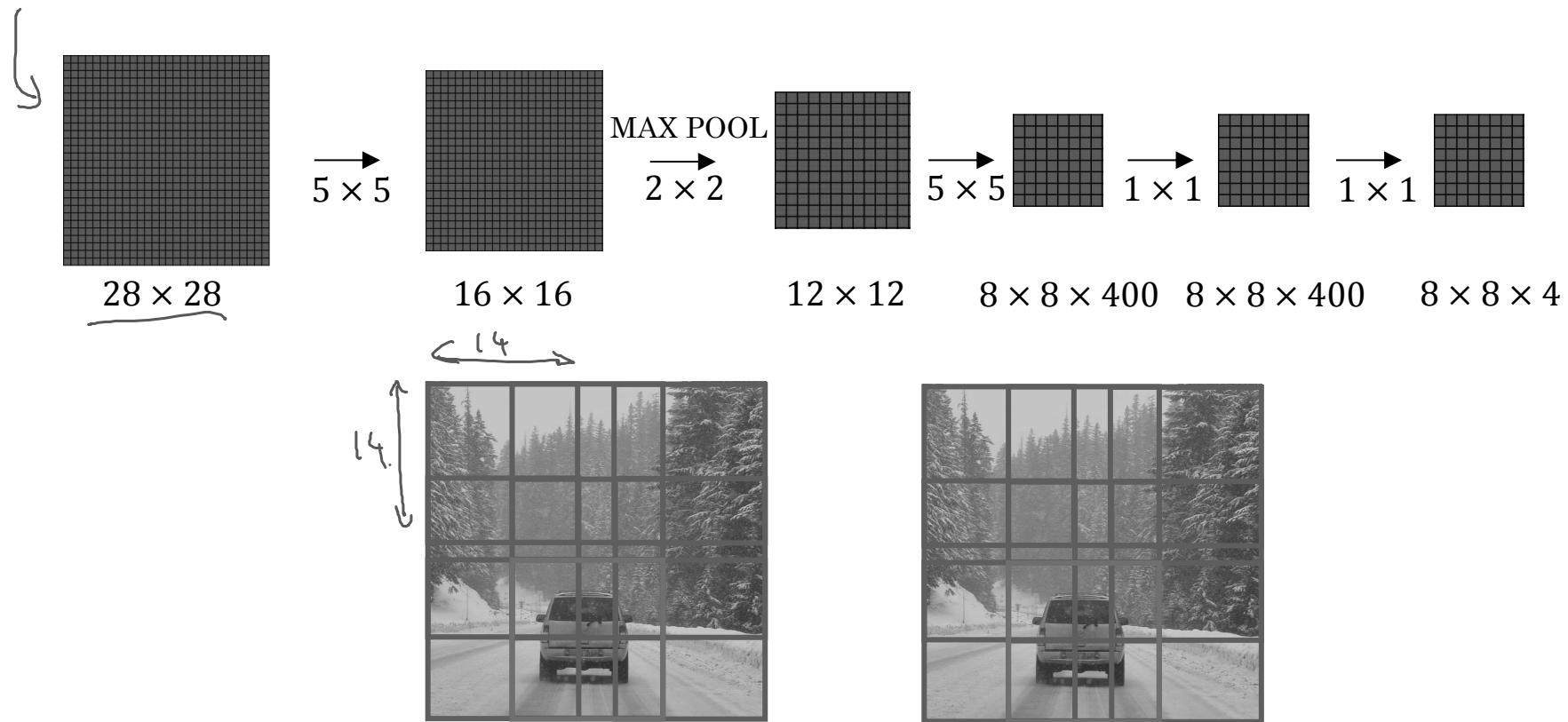


Convolution implementation of sliding windows



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

Convolution implementation of sliding windows





deeplearning.ai

Object Detection

Intersection over union

Evaluating object localization



Intersection over Union (IoU)

$$= \frac{\text{Size of } \begin{array}{c} \diagup \\ \diagdown \end{array}}{\text{Size of } \begin{array}{c} \diagup \\ \diagdown \end{array}}$$

"Correct" if IoU ≥ 0.5 \leftarrow

0.6 \leftarrow

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



deeplearning.ai

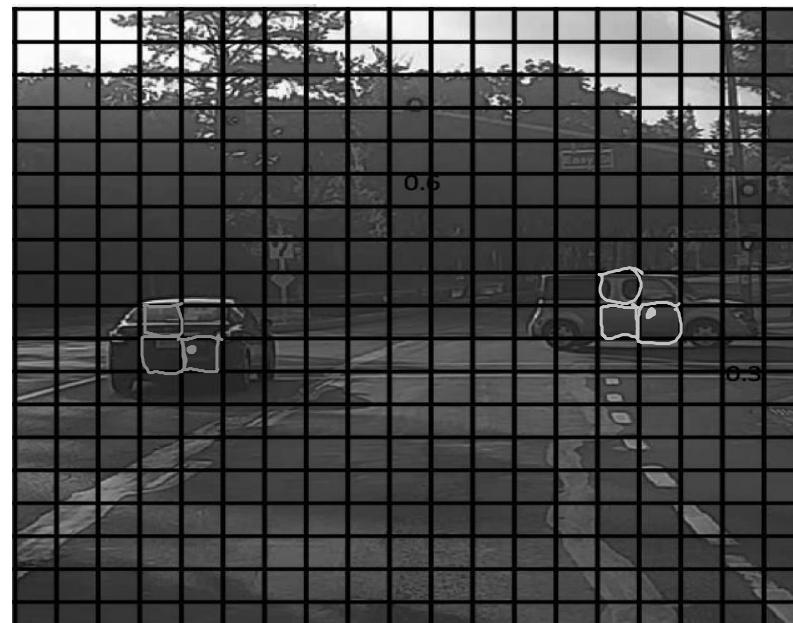
Object Detection

Non-max suppression

Non-max suppression example



Non-max suppression example

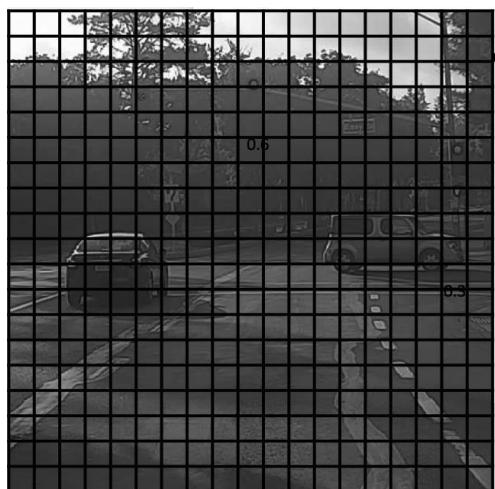


19x19

Non-max suppression example

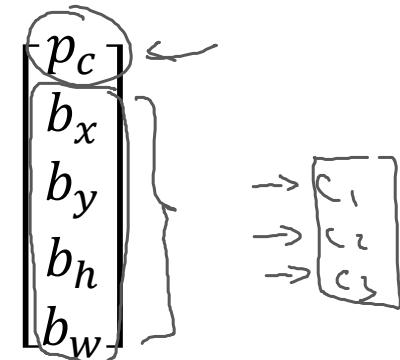


Non-max suppression algorithm



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 IoU ≥ 0.5 with the box output in the previous step

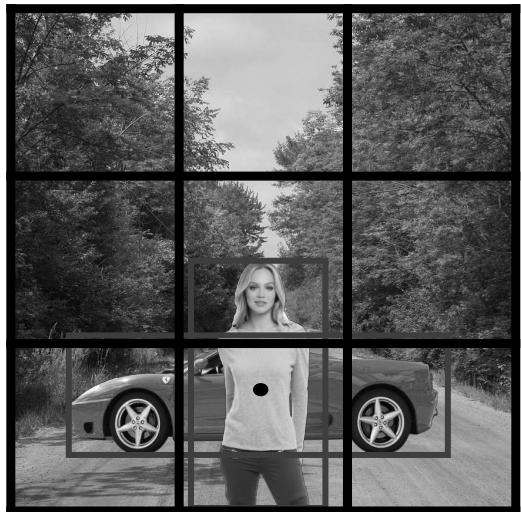


deeplearning.ai

Object Detection

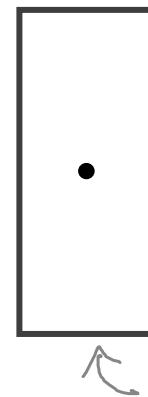
Anchor boxes

Overlapping objects:

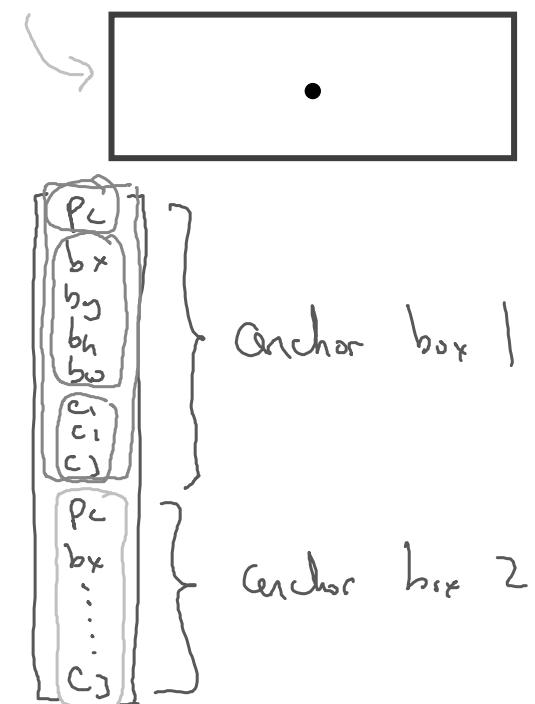


$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

Anchor box 1:



Anchor box 2:



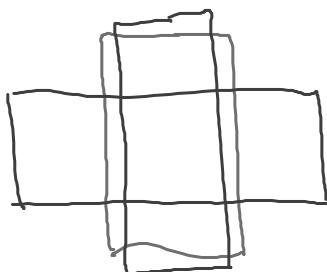
[Redmon et al., 2015, You Only Look Once: Unified real-time object detection]

Anchor box algorithm

Previously:

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

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



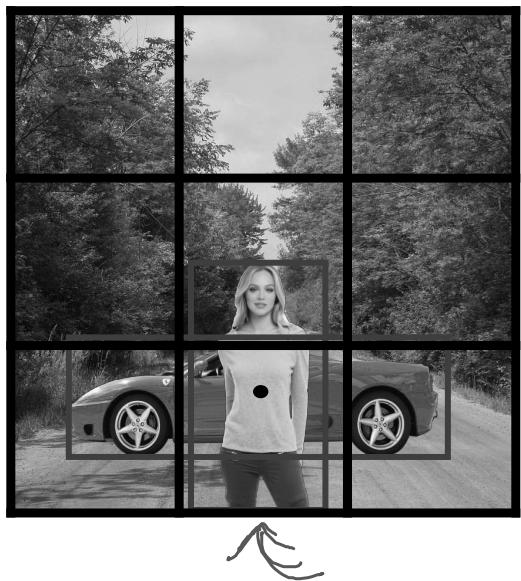
With two anchor boxes:

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.

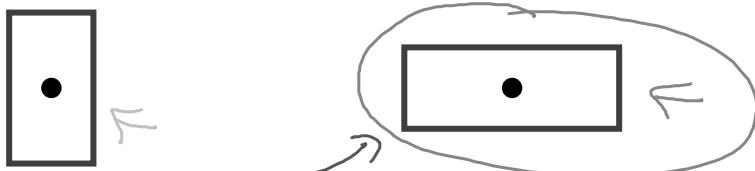
(grid cell, anchor box)

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

Anchor box example

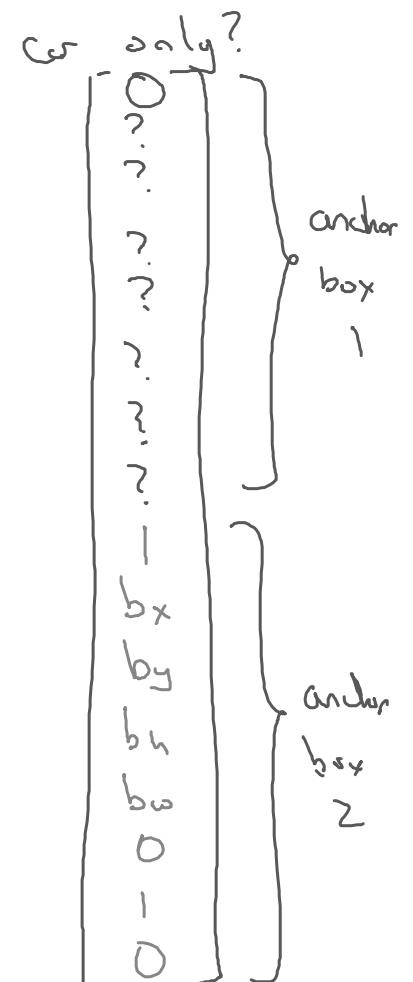


Anchor box 1: Anchor box 2:



$y =$

$$\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \quad \begin{bmatrix} 1 \\ bx \\ by \\ bh \\ bw \\ 1 \\ 0 \\ 0 \\ 1 \\ bx \\ by \\ bh \\ bw \\ 0 \\ 1 \\ 0 \end{bmatrix}$$



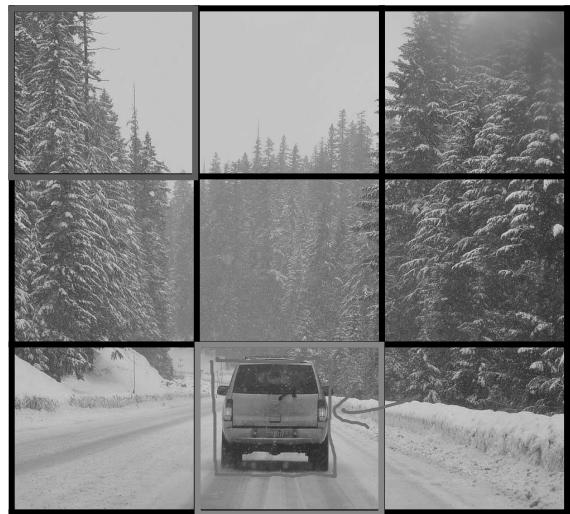


deeplearning.ai

Object Detection

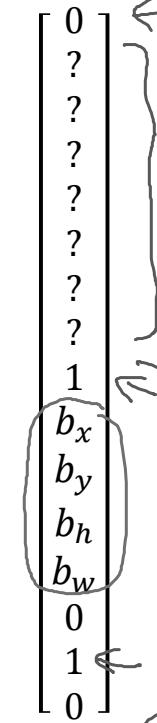
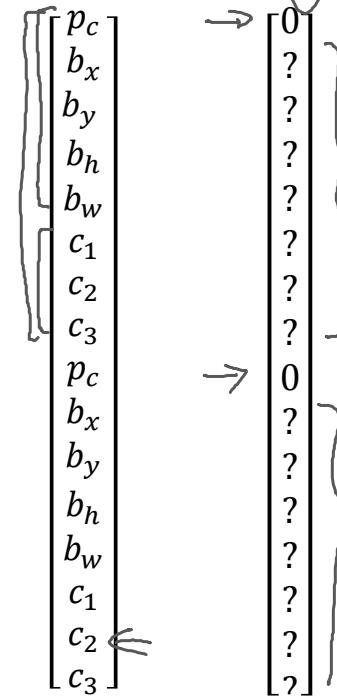
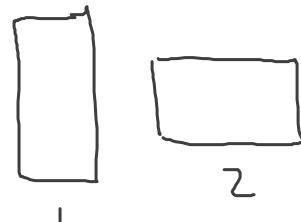
Putting it together:
YOLO algorithm

Training



- 1 - pedestrian
- 2 - car
- 3 - motorcycle

$$y =$$



$3 \times 3 \times 16$

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

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

#anchors

$5 + \# \text{classes}$

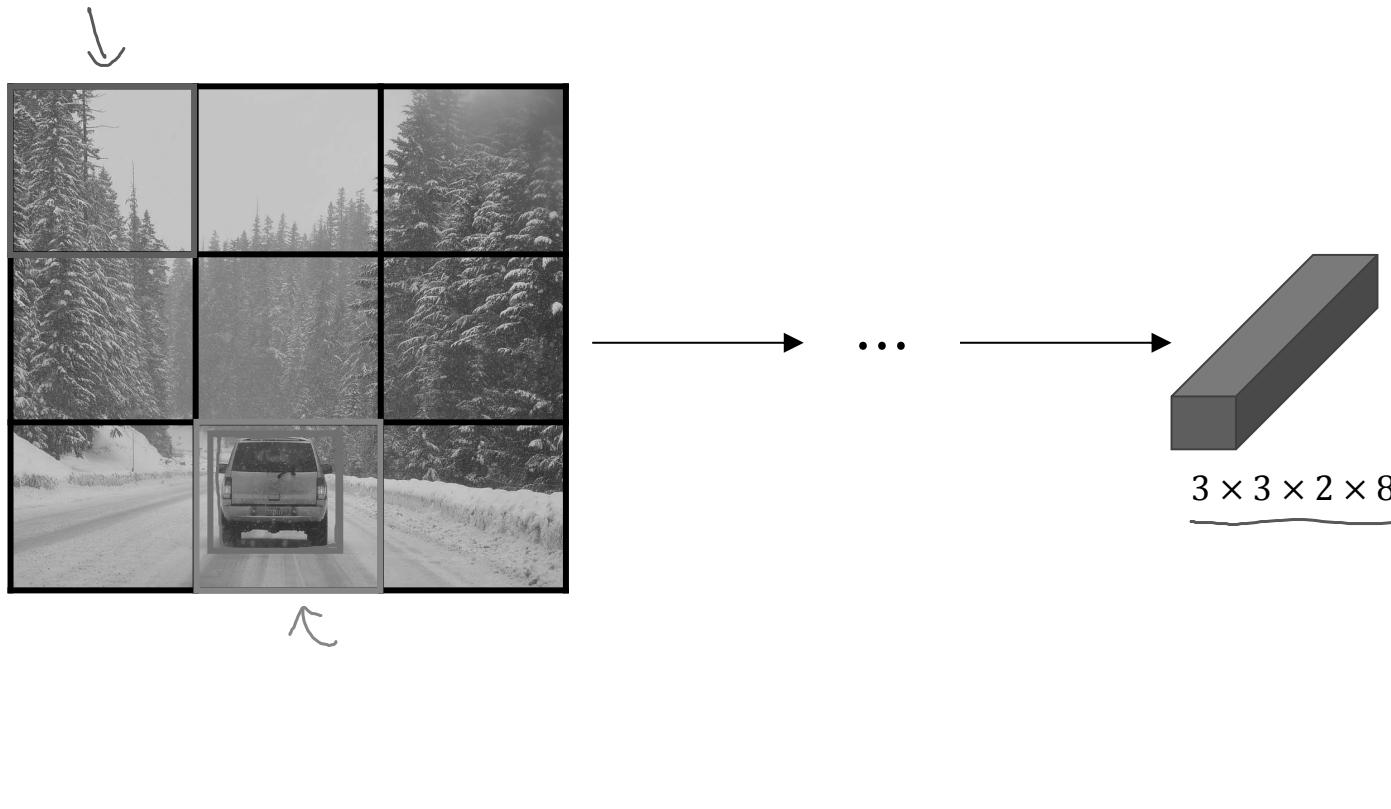
$100 \times 100 \times 3$

\rightarrow ConvNet \rightarrow

$3 \times 3 \times 16$

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

Making predictions



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

Annotations on the right side of the equation:

- An arrow points from the first p_c entry to a vertical vector $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$.
- An arrow points from the second b_x entry to a vertical vector $\begin{bmatrix} 1 \\ b_x \end{bmatrix}$.
- An arrow points from the third b_y entry to a vertical vector $\begin{bmatrix} 1 \\ b_y \end{bmatrix}$.
- An arrow points from the fourth b_h entry to a vertical vector $\begin{bmatrix} 1 \\ b_h \end{bmatrix}$.
- An arrow points from the fifth b_w entry to a vertical vector $\begin{bmatrix} 1 \\ b_w \end{bmatrix}$.
- An arrow points from the sixth c_1 entry to a vertical vector $\begin{bmatrix} 1 \\ c_1 \end{bmatrix}$.
- An arrow points from the seventh c_2 entry to a vertical vector $\begin{bmatrix} 1 \\ c_2 \end{bmatrix}$.
- An arrow points from the eighth c_3 entry to a vertical vector $\begin{bmatrix} 1 \\ c_3 \end{bmatrix}$.

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.



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]

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]