CMPE 258, Deep Learning

# Object Detection

April 5, 2018

DMH 149A

## Taehee Jeong
Ph.D., Data Scientist

# Assignment_5

## Grading policy:

The code is supposed to be executable without any extra effort and produce reasonable result within 50 minutes.

If the code cannot be executable with any error or taking more than 50 minutes, 50 points will be assigned.

If the code can be executable without any error within 50 minutes, score will be assigned as following formula.

**Score = (10 – cost ) * 10**

Re-submitting is available until March 15th, but 10 point will be deducted every re-submitting after March 8th.

If extra effort is needed to get reasonable result (whatever it is), 5 to 10 points will be deducted.

**You may use your trained weights and bias (transfer learning). In this case, please make sure to submit the trained weights and bias as one separate file (***para_yourFirstName_LastName.hdf5)*

SJSU SAN JOSÉ STATE UNIVERSITY

# Mid-term Exam_2

Start Morning on April 12$^{th}$ .

End the midnight on April 15$^{th}$

Image classification using CNN

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# Group Project Proposal

Title submission deadline: April 9th

- Project title

- List of Members

- Preferred presentation day: 4/12 or 4/24

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

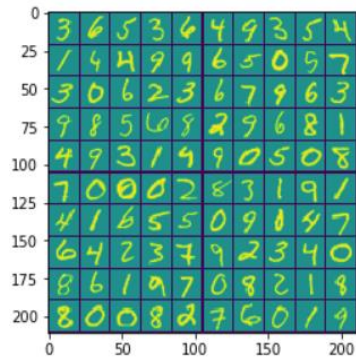# Group Project Proposal

## Content during proposal

- Justification for the project

- Background: any relevant previous work

- How to collect data set

- Which algorithms / platform will be used
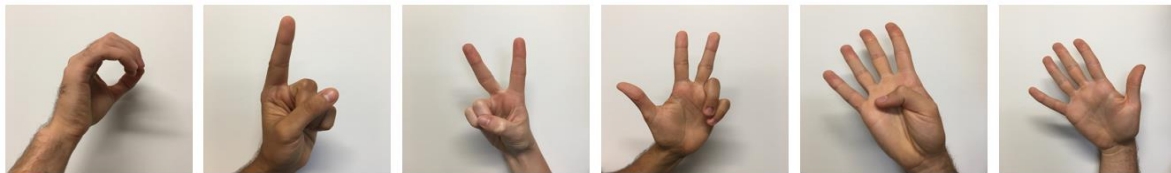
- What is the role for each team member

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# Today's lesson

Object detection

- Sliding windows
- 1 x 1 convolution
- Bounding box
- Intersection over union
- Non-max suppression

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# Image classification



Images for Hand written digits

Signs images



| y = 0 | y = 1 | y = 2 | y = 3 | y = 4 | y = 5 |

Coursera (Deep Learning specialization)

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# Image classification

Input image                                                          Output



Deep Neural Network                    softmax    Prediction:
                                                  Cat or Dog?
Convolution Neural Network

© Taehee Jeong                              SJSU SAN JOSÉ STATE UNIVERSITY

# Object localization



Where is the car located in the image?

<deep learning, Andrew Ng>

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

Pictures taken from a camera while driving around the Silicon Valley [drive.ai](https://www.drive.ai/)

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# Object localization

Bounding box

<deep learning, Andrew Ng>

© Taehee Jeong

# Classification with localization

(0,0)



$b_y$

$b_x$

$b_h$

$b_w$

(1,1)

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

<deep learning, Andrew Ng>

© Taehee Jeong

# Classification with localization

(0,0)



$b_y$

$b_x$

$b_h$

$b_w$

<deep learning, Andrew Ng>

(1,1)

Convolution    Pooling    Convolution    Pooling    Fully connected

softmax

$\hat{y}$

$b_x$

$b_y$

$b_h$

$b_w$

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# Object detection



<deep learning, Andrew Ng>

How about multiple objects in one image?

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# Example of bounding box

$$y = (p_c, b_x, b_y, b_h, b_w, c)$$

$b_w$

$b_h$

$(b_x, b_y)$
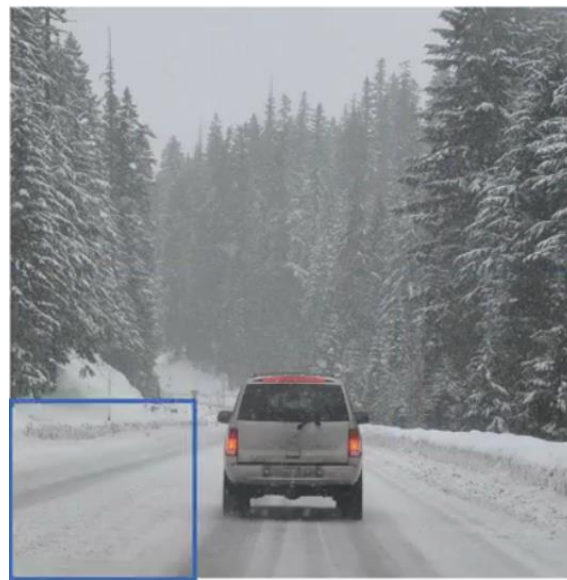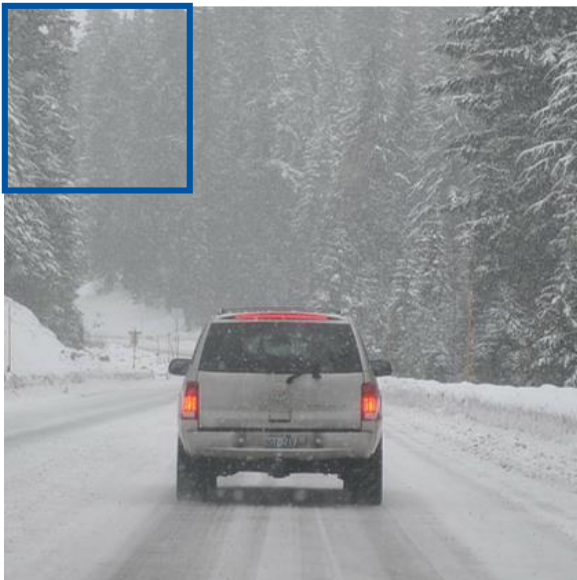
$p_c = 1$ : confidence of an object being present in the bounding box

$c = 3$ : class of the object being detected (here 3 for "car")

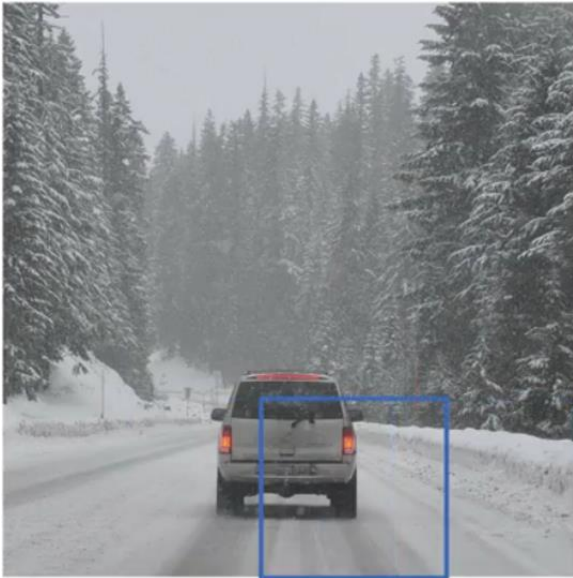<deep learning, Andrew Ng>

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# Sliding windows detection



<Deep Learning, Andrew Ng>

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# Sliding windows detection

<Deep Learning, Andrew Ng>

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# Sliding windows detection



Instead of one size,
Many different size of windows
(from small one to larger one) are
needed to apply.

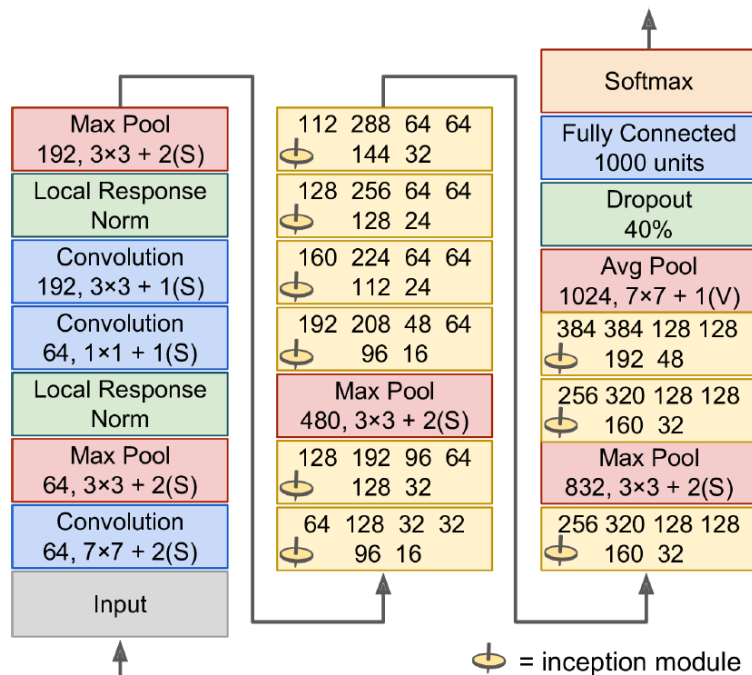<Deep Learning, Andrew Ng>

© Taehee Jeong

SJSU SAN JOSÉ STATE
UNIVERSITY

# Sliding windows detection

Implementation Problem : sliding windows detection takes long time to compute.

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# GoogLeNet

won the ILSVRC 2014 challenge with 93% accuracy.



1 x 1 convolution is used.

<Hands-on ML, Aurelien Geron>

"Going Deeper with Convolutions," C. Szegedy et al. (2015)

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# 1x 1 convolutions

In the case of one channel



| 1 | 2 | 3 | 6 | 5 | 8 |
|---|---|---|---|---|---|
| 3 | 5 | 5 | 1 | 3 | 4 |
| 2 | 1 | 3 | 4 | 9 | 3 |
| 4 | 7 | 8 | 5 | 7 | 9 |
| 1 | 5 | 3 | 7 | 4 | 8 |
| 5 | 4 | 9 | 8 | 3 | 5 |

6 × 6          *          2          =

1 x 1

6 x 6

<Deep Learning, Andrew Ng>

Lin et al., 2013. Network in network

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# 1 x 1 convolutions

In the case of multiple channels



$6 \times 6 \times 32$  *  $1 \times 1 \times 32$  =  $6 \times 6 \times \#\text{ filters}$
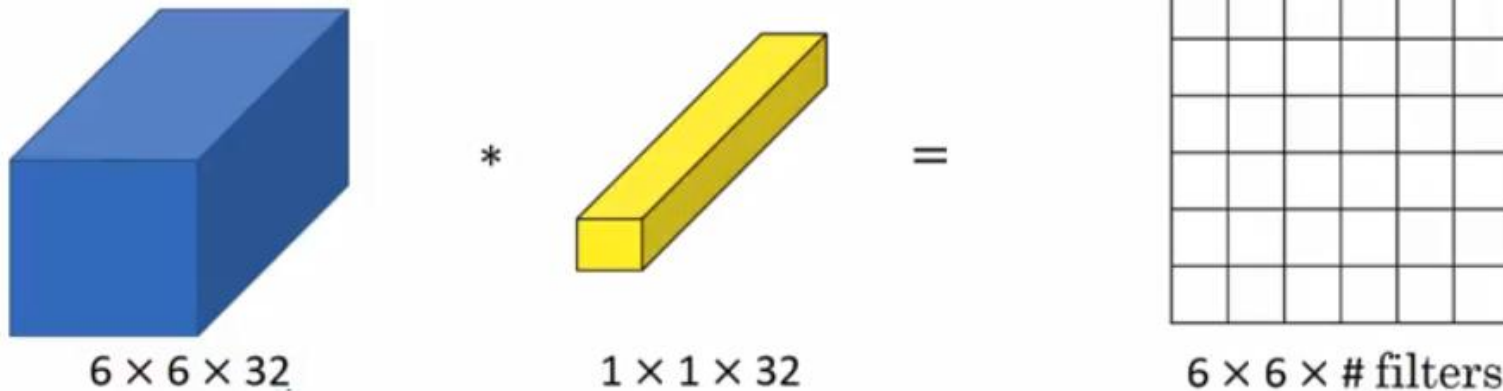
Convolution : sum of individual pixel over all channels

<Deep Learning, Andrew Ng>

Lin et al., 2013. Network in network

22          © Taehee Jeong
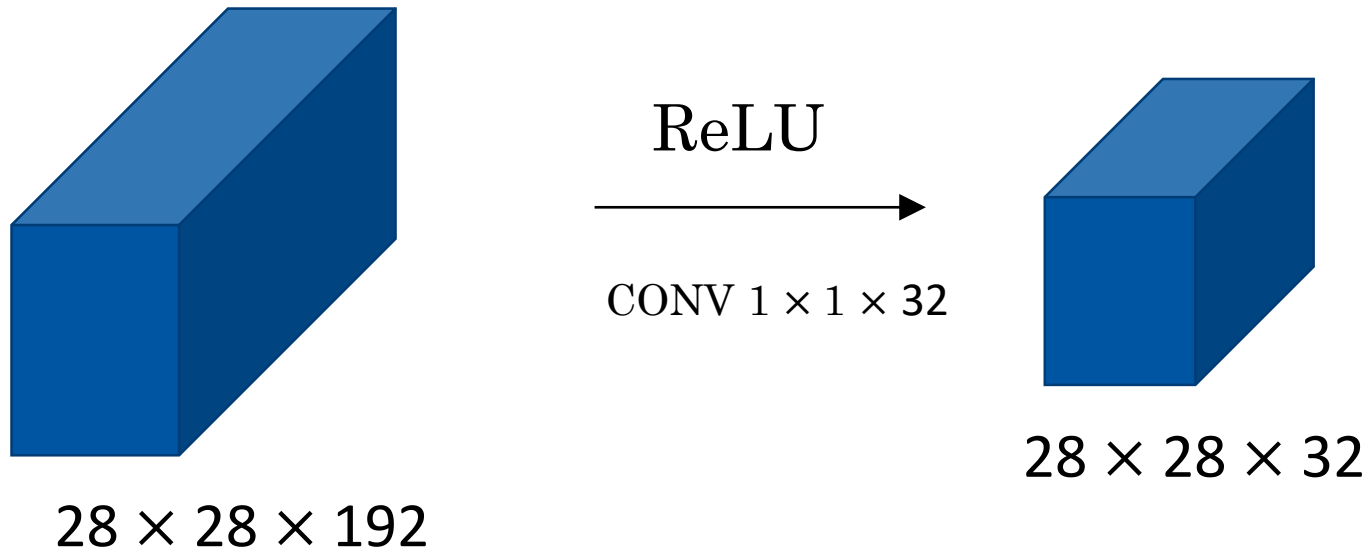
# 1 x 1 convolutions

**Network in network**



$6 \times 6 \times 32$      *      $1 \times 1 \times 32$      =      $6 \times 6 \times$ # filters

1 x 1 x 32 x $n_c$ (# filters)

**Convolution Activation**

<Deep Learning, Andrew Ng>

Lin et al., 2013. Network in network

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# 1×1 convolutions

ReLU

$\longrightarrow$

CONV $1 \times 1 \times 32$

$28 \times 28 \times 192$

$28 \times 28 \times 32$

<Deep Learning, Andrew Ng>

Lin et al., 2013. Network in network

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# Turning Fully connected layer into convolutional layers

Fully connected layers



Convolutional layers



<Deep Learning, Andrew Ng>

© Taehee Jeong

# Convolution implementation of sliding windows



14 × 14 × 3 → 5 × 5 → 10 × 10 × 16 → MAX POOL 2 × 2 → 5 × 5 × 16 → 1 x 1 → 1 × 1 × 400 → 1 × 1 (400) → 1 × 1 × 400 → 1 × 1 → 1 × 1 × 4

Sliding windows : 14 x 14, stride: 2



28 x 28 x 3 → 5 × 5 → 24 x 24 x 16 → MAX POOL 2 × 2 → 12 x 12 x 16 → 5 × 5 → 8 × 8 × 400 → 1 × 1 → 8 × 8 × 400 → 1 × 1 → 8 × 8 × 4
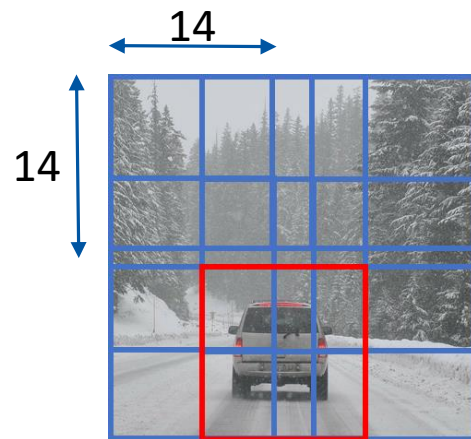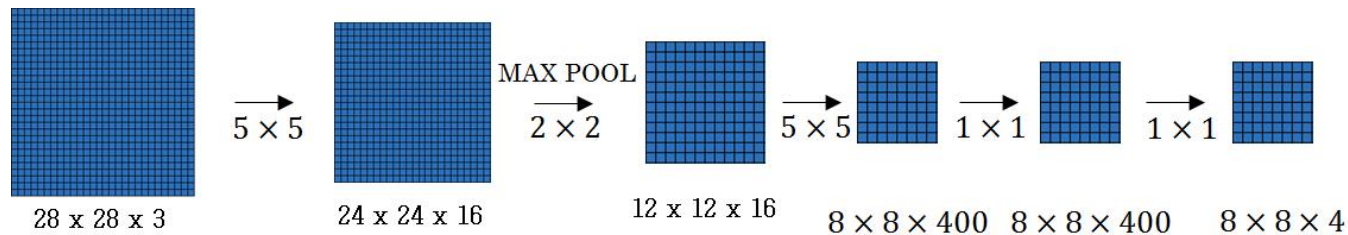
<Deep Learning, Andrew Ng>

Sermanet et al., 2014, OverFeat: Integrated recognition, localization and detection using convolutional networks
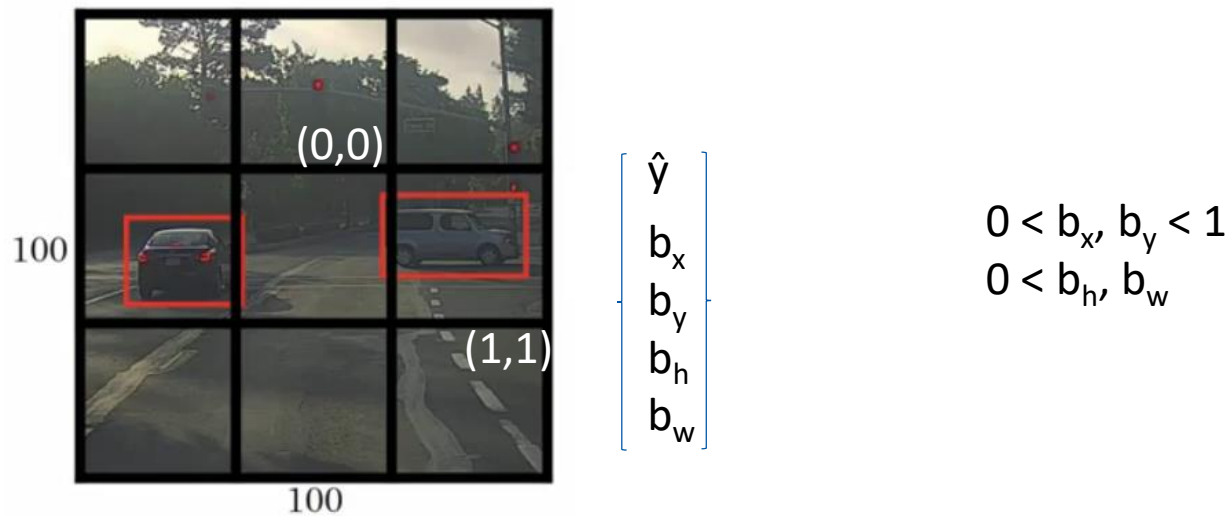
26      © Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# Convolution implementation of sliding windows

Sliding windows : 14 x 14, stride: 2



14

14

<Deep Learning, Andrew Ng>

© Taehee Jeong

SAN JOSÉ STATE UNIVERSITY

# Specify bounding boxes



$$\begin{bmatrix} \hat{y} \\ b_x \\ b_y \\ b_h \\ b_w \end{bmatrix}$$

$0 < b_x, b_y < 1$
$0 < b_h, b_w$

<Deep Learning, Andrew Ng>

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

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY
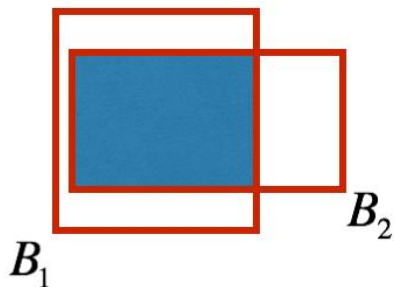
# YOLO (you only look once)



For each grid cell:

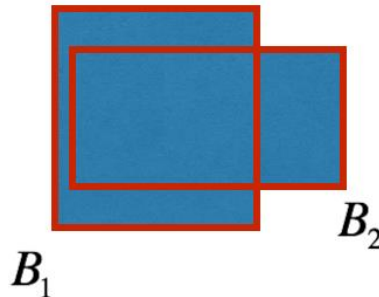$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

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

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# Intersection over Union (IoU)

**Intersection**



$B_2$

$B_1$

**Union**



$B_2$

$B_1$

**Intersection over Union**

$$IoU = \frac{B_1 \cap B_2}{B_1 \cup B_2} = \frac{\;\;\;\;\;\;\;\;\;\;\;}{\;\;\;\;\;\;\;\;\;\;\;} = P_c$$

"Correct" if IoU $\geq$ 0.5

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

<Deep Learning, Andrew Ng>

30          © Taehee Jeong
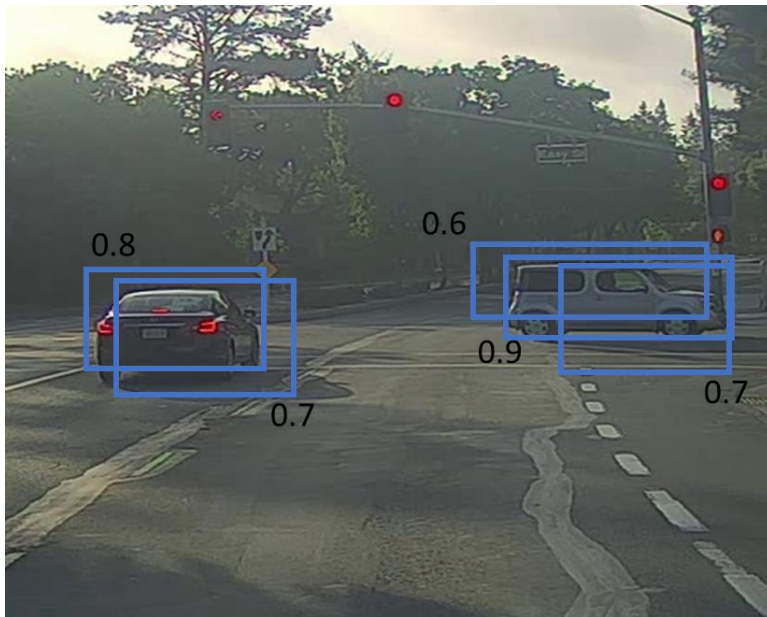
# Non-max suppression

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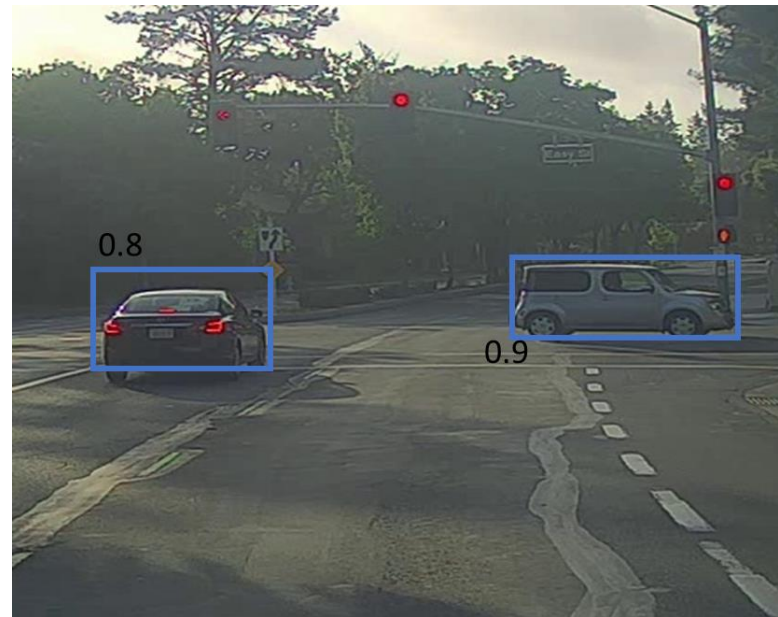
$$\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \end{bmatrix}$$

Discard all boxes with $p_c \leq 0.6$
Among remaining boxes,
Pick the box with the largest $p_c$.
Output that as a prediction.

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# Non-max suppression

Before non-max suppression



After non-max suppression



<Deep Learning, Andrew Ng>

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# Encoding architecture for YOLO

preprocessed image
(608, 608, 3)

encoding
(19,19, 5, 85)

Deep CNN

reduction
factor: 32

19

19

$p_c$  $b_x$  $b_y$  $b_h$  $b_w$          80 class probabilities

box 1          ...
box 2          ...
box 3          ...
box 4          ...
box 5          ...

<Deep Learning, Andrew Ng>

33          © Taehee Jeong

SAN JOSÉ STATE
UNIVERSITY

# Predictions of YOLO model <Deep Learning, Andrew Ng>



Pictures taken from a camera while driving around the Silicon Valley [drive.ai](https://www.drive.ai/)

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY

# Summary

## Object detection

- Sliding windows
- 1 x 1 convolution
- Bounding box
- Intersection over union
- Non-max suppression

© Taehee Jeong

SJSU SAN JOSÉ STATE UNIVERSITY