

DATA X

Object Detection with YOLO

Chad Wakamiya
Spring 2020

Agenda

Object Detection

Defining the object detection problem and a naive solution.

YOLO Algorithm

- YOLO algorithm steps
- Bounding boxes
- Measuring performance (UoI)
- Non-max suppression

YOLO Implementations

- Pretrained models with the COCO dataset.
- Custom trained models

The background is a solid dark blue color. It is filled with a repeating pattern of white, wireframe-style 3D cubes and icons. The icons include a hand cursor pointing at a cube, a magnifying glass with a plus sign, a funnel, and various stacked cube configurations. The text 'Object Detection' is centered in the middle of the image in a large, bold, white sans-serif font.

Object Detection

Classification vs. Object Detection

Object Detection is the problem of locating and classifying objects in an image.

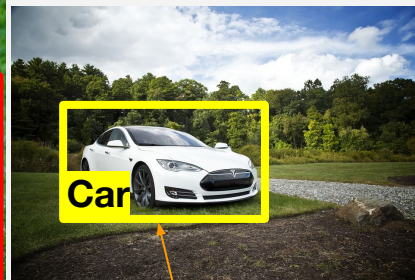
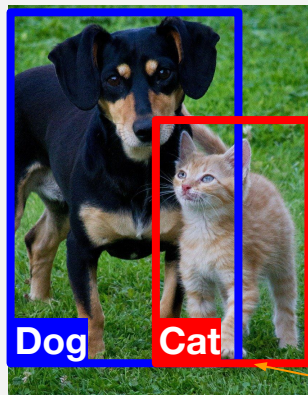
Classification

- Each image has one object
- Model predicts one label



Object Detection

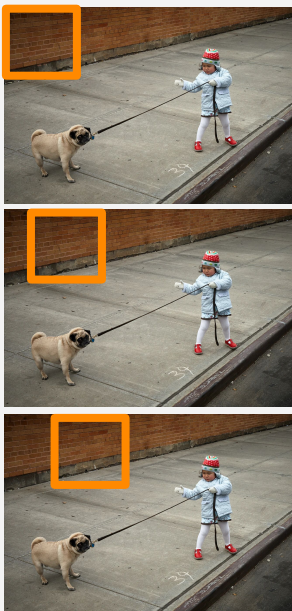
- Each image may contain multiple objects
- Model classifies objects and identifies their location.



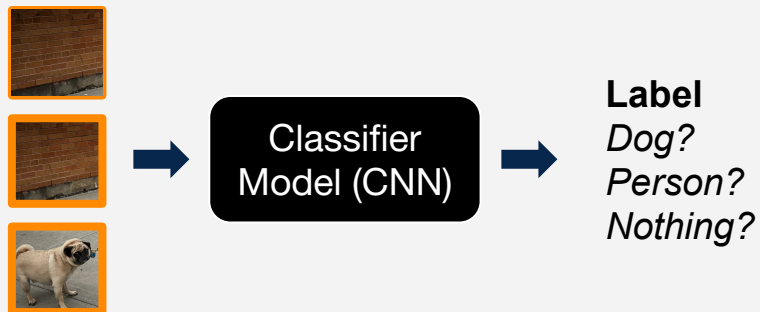
Bounding Box

Naive Approach

1. Scan the image with a sliding window.



2. Feed the images into a classifier model to predict a label for that region.



- This approach is slow since it checks many windows that don't contain anything -> Not good for real time uses.
- The Region-based Convolutional Neural Net (**R-CNN**) is an improved version that strategically selects regions that are likely to contain an object to run through the CNN.



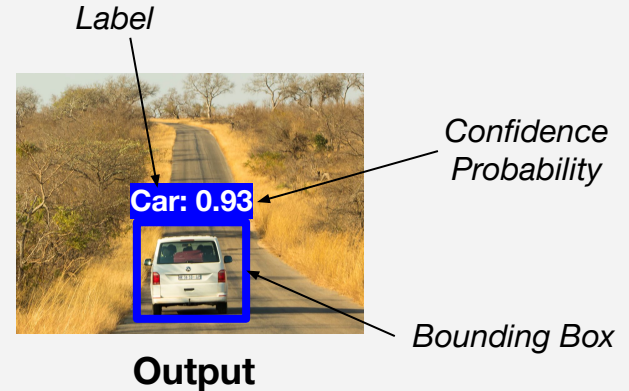
YOLO Algorithm

YOLO "You Only Look Once"

- Instead of making predictions on many regions of an image, YOLO passes the *entire* image at once into a CNN that predicts the **labels, bounding boxes, and confidence probabilities** for objects in the image.
- YOLO runs much **faster** than region based algorithms quick because requires only a single pass through a CNN.



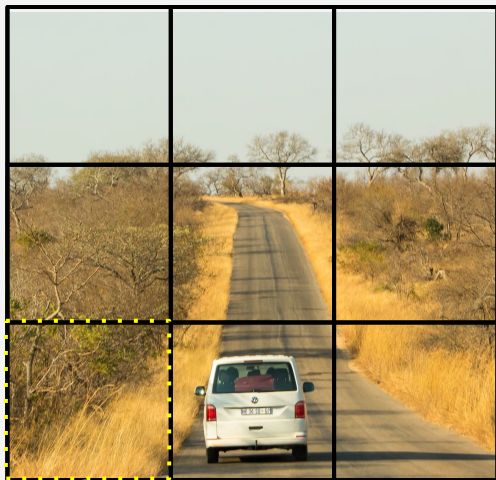
Input



Output

YOLO Steps

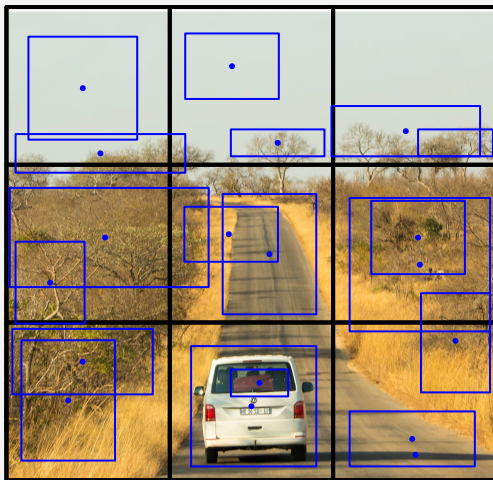
1. Divide the image into cells with an $S \times S$ grid.



$S = 3$

Cell

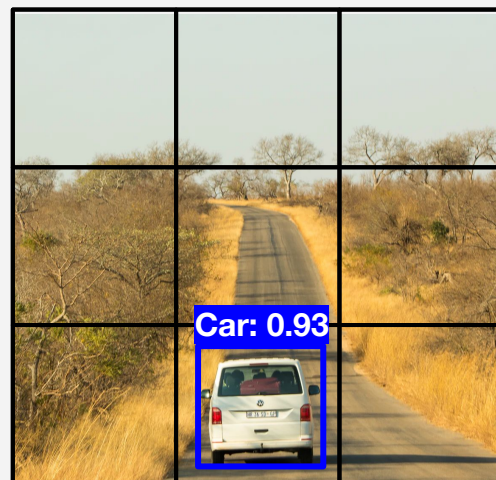
2. Each cell predicts B bounding boxes.



$B = 2$

A cell is responsible for detecting an object if the object's bounding box falls within the cell. (Notice that each cell has 2 blue dots.)

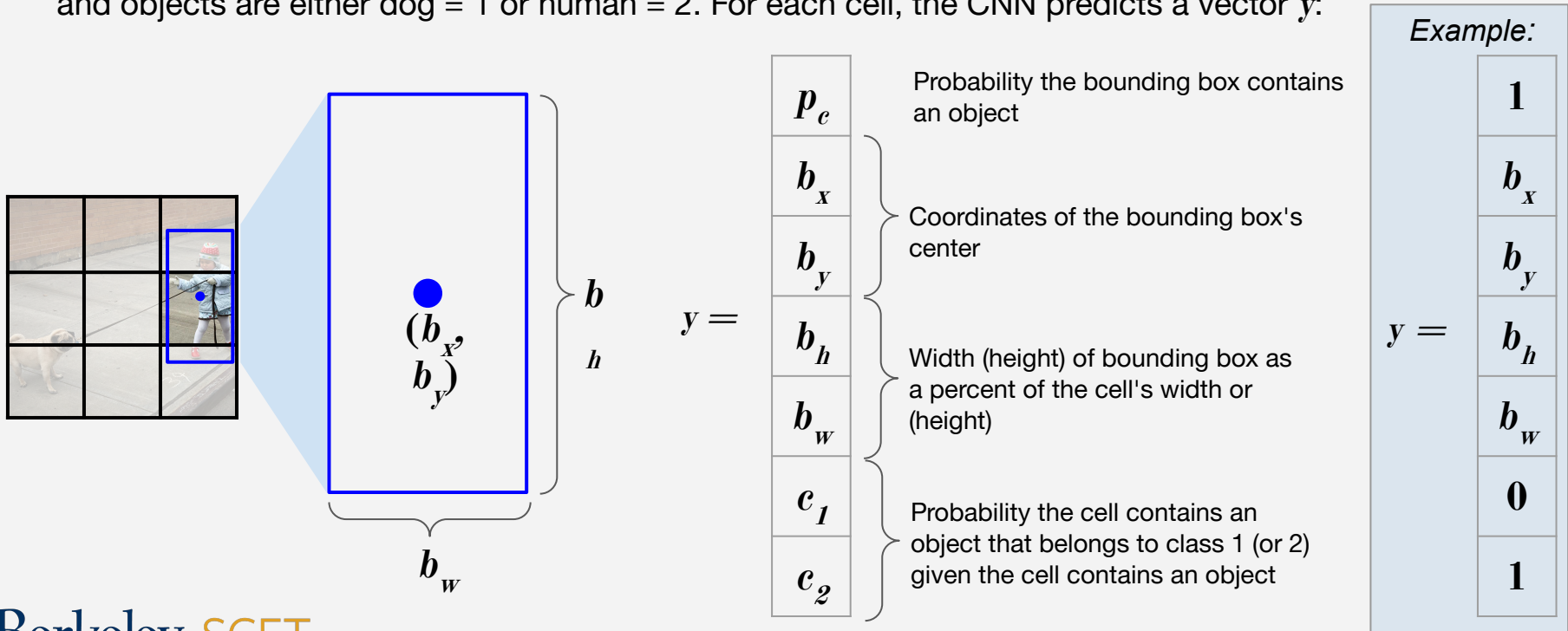
3. Return bounding boxes above confidence threshold.



All other bounding boxes have a confidence probability less than the threshold (say 0.90) so they are suppressed.

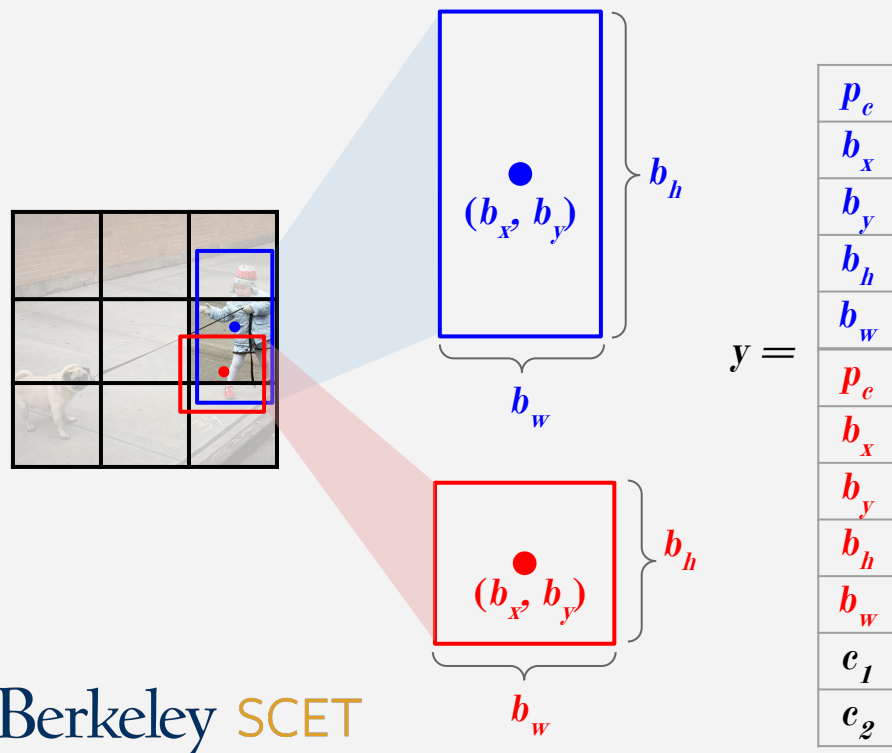
How are bounding boxes encoded?

Let's use a simple example where there are 3x3 cells ($S=3$), each cell predicts 1 bounding box ($B=1$), and objects are either dog = 1 or human = 2. For each cell, the CNN predicts a vector y :

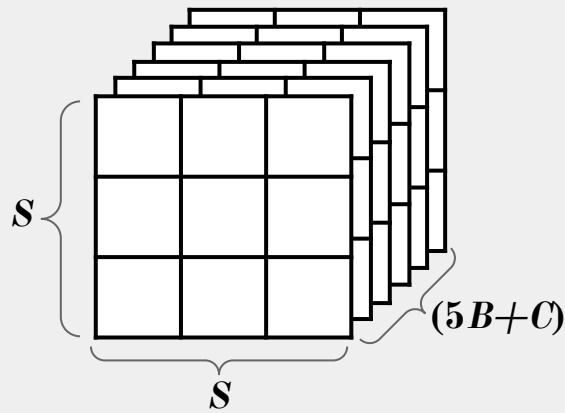


Encoding Multiple Bounding Boxes

What happens if we predict multiple bounding boxes per cell ($B > 1$)? We simply augment y .



The CNN will predict a y for each cell, so the size of the output tensor (multidimensional "matrix") should be: $S \times S \times (5B + C)$



Notice that y has $5B + C$ elements (C is the number of classes).

YOLO Overview

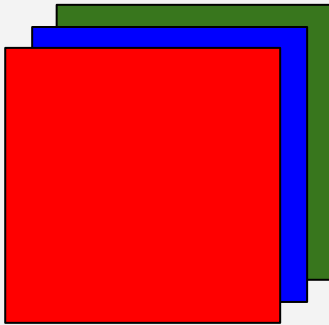
Input



Convolutional Neural Net



Output

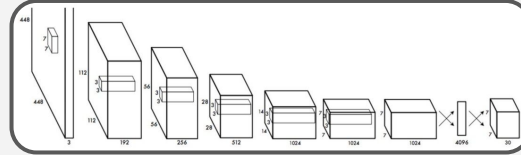


$W \times H \times 3$

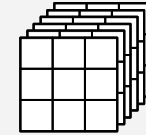
W : Width of image in pixels

H : Height of image in pixels

3 : Number of color channels in RGB



Series of convolutional and pooling layers.

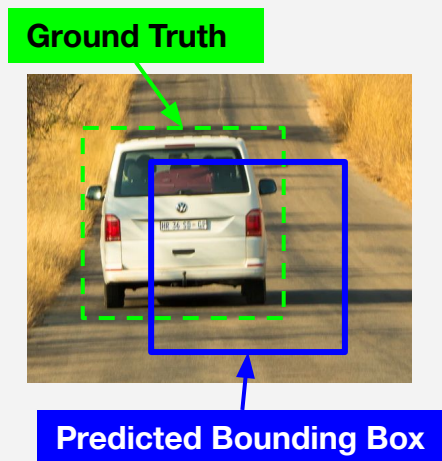


$S \times S \times (5B + C)$

A tensor that specifies the bounding box locations and class probabilities.

Measuring Performance with UoI

- **Union over Intersection (UoI)** measures the overlap between two bounding boxes.
- During training, we calculate the UoI between a predicted bounding box and the ground truth (the pre-labeled bounding box we aim to match)

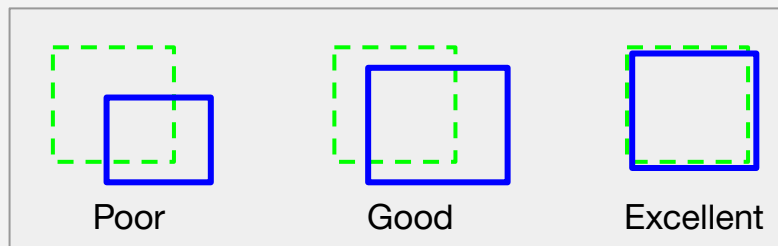
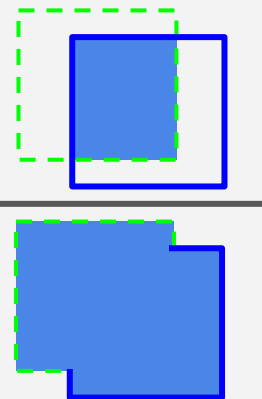


**Union over
Intersection**

=

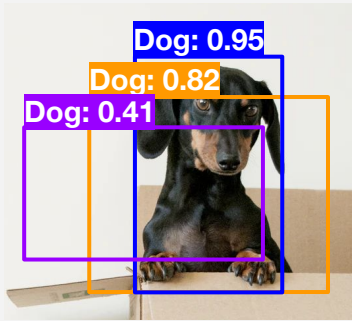
Area of Intersection

Area of Union



Double Counting Objects (Non-Max Suppression)

- When predicting more than 2 bounding boxes per cell, sometimes the same object will be detected multiple times (overlapping boxes with the same label)
- **Non-max suppression** solves multiple counting by removing the box with the lower confidence probability when the UoI between 2 boxes with the same label is above some threshold.



Non-Max Suppression

The diagram illustrates the Non-Max Suppression process in three stages:

- 1. Identify the box with the highest confidence.** The image shows the three bounding boxes from the previous step. The blue box (0.95) is the highest confidence.
- 2. Calculate the UoI between the highest confidence box each of the other boxes.** Two small icons represent the UoI calculations: an orange box overlapping the blue box with 'UoI: 0.62', and a purple box overlapping the blue box with 'UoI: 0.47'.
- 3. Suppress boxes with UoI above a selected threshold (usually 0.3).** The final image shows only the blue box (0.95) remaining, as the other two boxes were suppressed because their UoI values (0.62 and 0.47) were above the 0.3 threshold.



Implementing YOLO

Pretrained Models

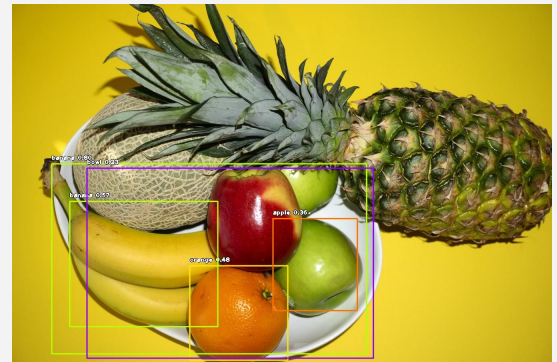
- Training a YOLO model requires images labeled with bounding boxes. These datasets may take time to label, so readily available pre-labeled images are often used to train models.
- A common dataset for image classification/detection/segmentation is the [COCO \(Common Objects in Context\)](#), a database of images with 80 labelled classes.
- Popular pretrained YOLO models with COCO:
 - [ImageAI](#) (easy-to-use, lightweight YOLO implementation)
 - [Darknet](#) (trained by the author of YOLO)



YOLO Implementation
(CNN)



Pretrained Model
with COCO



Pineapples and cantaloupes are not in COCO so they are not recognized.

COCO Pretrained Labels

Applications built with COCO trained models will only be able to identify these objects!

person	fire hydrant	elephant	skis	wine glass	broccoli	diningtable	toaster
bicycle	stop sign	bear	snowboard	cup	carrot	toilet	sink
car	parking meter	zebra	sports ball	fork	hot dog	tvmonitor	refrigerator
motorbike	bench	giraffe	kite	knife	pizza	laptop	book
aeroplane	bird	backpack	baseball bat	spoon	donut	mouse	clock
bus	cat	umbrella	baseball glove	bowl	cake	remote	vase
train	dog	handbag	skateboard	banana	chair	keyboard	scissors
truck	horse	tie	surfboard	apple	sofa	cell phone	teddy bear
boat	sheep	suitcase	tennis racket	sandwich	pottedplant	microwave	hair drier
traffic light	cow	frisbee	bottle	orange	bed	oven	toothbrush

Custom Models

- If your use case only uses objects in COCO → you can use a pretrained model.
- Otherwise you will need to train your own YOLO model. This will require:
 1. Finding images of the objects to recognize.
 2. Label bounding boxes.
 3. Train your YOLO model. There are 2 options:
 - a. Implement your own model using OpenCV, Tensorflow/Keras
 - b. Use [ImageAI](#)'s custom training methods.

References/Further Reading

- **YOLO**
 - [://towardsdatascience.com/you-only-look-once-yolo-implementing-yolo-in-less-than-30-lines-of-python-code-97fb9835bfd2](https://towardsdatascience.com/you-only-look-once-yolo-implementing-yolo-in-less-than-30-lines-of-python-code-97fb9835bfd2)
- **R-CNN**
 - <https://towardsdatascience.com/r-cnn-fast-r-cnn-faster-r-cnn-yolo-object-detection-algorithms-36d53571365e>
- **CNN**
 - <https://www.coursera.org/lecture/convolutional-neural-networks/optional-region-proposals-aCYZv>
- **YOLO**
 - <https://hackernoon.com/understanding-yolo-f5a74bbc7967>
 - <https://www.analyticsvidhya.com/blog/2018/12/practical-guide-object-detection-yolo-framework-python/>
- **Intersection Over Union**
 - <https://www.pyimagesearch.com/2016/11/07/intersection-over-union-iou-for-object-detection/>