

Class 6: Real-Time Object Detection

Problem Addressed: Object Detection

- งานยาก เรื่องของการระบุตำแหน่งและการจำแนกประเภทของวัตถุ (locate and classify objects)
- Goal \Rightarrow fast & high accuracy สำคัญ
- นอกจาก detect เราต้อง classify and localize เช่น แมว เบ็ด หมา คอมฯ ต้องรู้ว่าภาพนี้ประกอบไปด้วยอะไรบ้าง

☐ Importance

- Visual modality is very powerful
- Responsive robot system \Rightarrow required real-time vision based object detection

☐ YOLO concept

- YOLOv1 \Rightarrow มีการแบ่งภาพ ออกเป็น grid และเมื่อ center ของวัตถุ อยู่ใน grid cell หมายความว่า เซลล์จะมีหน้าที่ detect object นั้น จากนั้น พัฒนาเป็น
- YOLOv2 \Rightarrow boundary box สามารถ detect object class > 9k class
- YOLOv3 \Rightarrow ให้ box detect sizes object \Rightarrow scale ต่างกันมากขึ้น
- YOLOv4, v5, v6, tiny, PP \Rightarrow speed & acc
- YOLOR, X \Rightarrow improved times generalize
- YOLOv7 \Rightarrow focus small optimization
- YOLOv8, YOLONAS \Rightarrow improve performance trained on coco dataset \Rightarrow 2023

YOLO Overview

\rightarrow split photo into $S \times S$ grid : each cell, detect object โดย center point ตกอยู่ใน grid \Rightarrow ช่วยในการระบุตำแหน่ง

\rightarrow each grid, predict bounding box \Rightarrow ในแต่ละ box จะมี object อยู่ในกล่องนั้น ก็คือจะทำนายตัวแปร

- x, y : coordinate of center object that are related with cell ค่าขนาดของ grid cell จะอยู่ในช่วง [0,1] เพราะผ่านการ normalized เช่น มุมซ้ายบน กับมุมขวาล่าง เท่ากับ [1,1] then center point of this cell is [0.5, 0.5]
- w, h : แสดงถึงความกว้างและความสูงของ bounding box
- Confidence: ทำหน้าที่ระบุว่า bounding box accurate แค่ไหน \Rightarrow ผ่านการทำ IOU \Rightarrow ideal = 0

\rightarrow YOLO is one of the fastest model for object detection

☐ YOLO training

- YOLO is a regression algorithm. \Rightarrow Predict bounding box from image
- Need to understand Parameters!!!
 - X : input \Rightarrow image \Rightarrow array or matrix of pixel (w*h) values (RGB values)
 - Y : output \Rightarrow tensor มีขนาด $S \times S \times (B(\text{box}) \times 5(\text{parameter}) \times C)$ แต่ละ grid ทำหน้าที่ predict class and distribution for a grid block

☐ YOLO Architecture

\rightarrow YOLO contain 7 convolution layers, ในการทำงานกับ image CNN เป็น เทคนิคที่เหมาะสม เพราะสามารถ capture ข้อมูลเชิงพื้นที่ได้ค่อนข้างดี, extract feature

\rightarrow ตัวอย่างเช่น

- input $S(\text{size}) = 7$, $B(\text{bounding box}) = 2$, $C(\text{number of class}) = 20$
- Output is $S \times S \times (5B + C) \Rightarrow 7 \times 7 \times (5 \times 2 + 20) = 7 \times 7 \times 30 \Rightarrow$ tensor size จาก YOLO

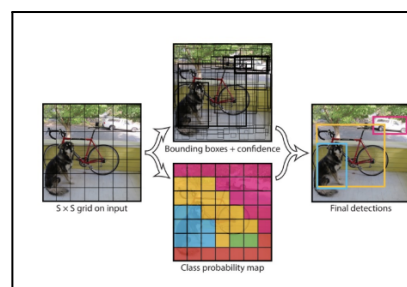
☐ Non-maximal suppression

\rightarrow case : multiple box, redundancy

\rightarrow filter bounding box and setting threshold : เอา low confidence score ออกไป \rightarrow screen

\rightarrow identify class of object : ต้องคำนวณ class score และหา argmax ช่วยระบุคลาสของ object

$$\Pr(\text{Class}_i | \text{Object}) * \Pr(\text{Object}) * \text{IOU}_{\text{pred}}^{\text{truth}} = \Pr(\text{Class}_i) * \text{IOU}_{\text{pred}}^{\text{truth}}$$



☐ YOLO Prediction

→ it's still possible to get redundant boxes but after we finalize with the high overlap (keep only highest confidence) ⇒ adds 2-3% on final MAP score

YOLO Objective Function

☐ Localization Loss

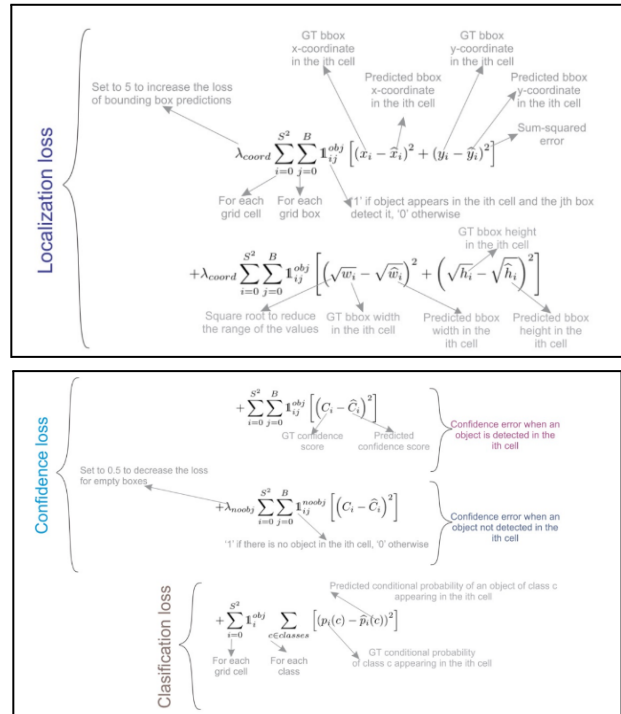
- คำนึงถึง ขนาดของ bounding box
- ensure bounding box from object predict มัน match กับ box ground truth (GT) หา error
- คำนึงถึง position (x,y): mean square error loss
- size (w,h): mean square root error loss
- อยากให้ loss ต่ำ ค่าใกล้เคียง GT

☐ Classification Loss

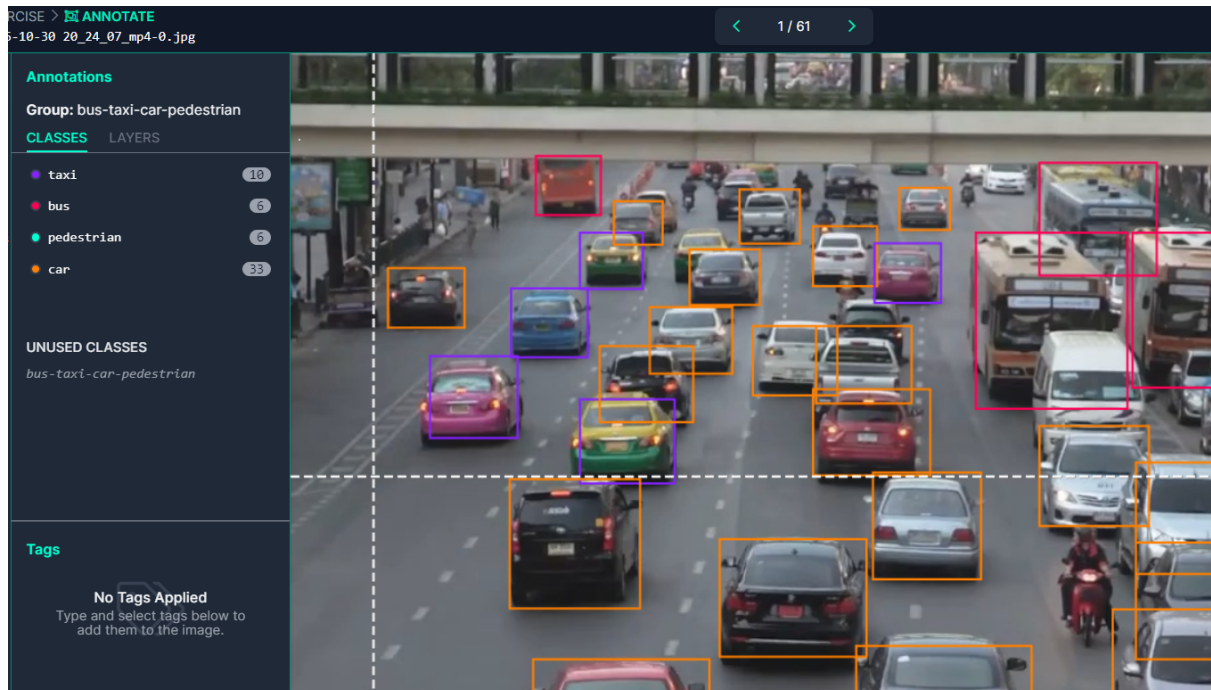
- คำนึงถึง การ identify object class

☐ Confidence Loss

- คำนึงถึง การ predict bounding box



☐ Exercise



YOLOv8

