

1 Abstract

The developer YOLO model also known as researchers implemented a new approach, YOLO, for object detection. Previously, the researchers solved the problem with multitask for object detection, but the researchers overriden it as a regression problem in YOLO. One neural network predicts class probability and bounding box by calculating the entire image only once. The bounding box is a rectangular box that tells the location of an object. The bounding box has the form of wrapping around the object it wants to detect. Class probability is the probability that the bounding box indicates which class an object belongs to. This probability value is presented as a conditional probability. YOLO is an end-to-end pipeline since it is constructed with one neural network. Finally, YOLO has faster processing speed compared to other neural networks. This speed is enough to process a video.

2 Introduction

I summarize the terms prior to the thesis summary. Detection model means to override the classifier and use it as a detector. Classification finds which class the object belongs to. For example, it is to find out whether a picture that humans think is a cat is actually a cat. In object detection, location information is required to determine. Existing models include DPM and R-CNN[4].

DPM is an abbreviation for Deformable Parts Models that detect objects in images by sliding window. R-CNN is a neural network for generating bounding boxes that enclose objects for a given whole image, using region proposal. Create a bounding box using region proposal and classify it by applying a classifier. Post-processing is applied to classified bounding box objects by adjusting, deduplication, and reassigning box scores[5]. Due to the complexity of the process, R-CNN neural networks have relatively slow computational rates. Furthermore, there exist difficulties in optimization. YOLO replaced the previous two neural networks with one regression problem. This means that the image pixels can be computed as a regression problem, the position of the bounding box, class probabilities.

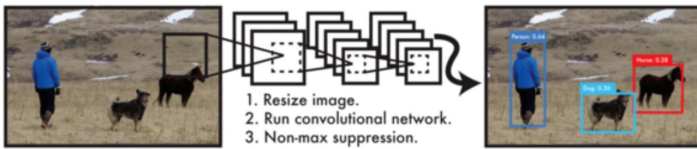


Figure 1: The YOLO described in figure 1 is as follows. The size of the input image is 448×448 . This image enters a single convolutional network to find multiple bounding boxes at the same time and calculates the corresponding class probabilities.

Because it is a single neural network, detection performance optimization is easily achieved. The advantages of YOLO are as follows. First, YOLO has a fast computational speed. The speed of YOLO is fast because it has been changed to a regression problem. YOLO also does not need a complex pipeline. In the Titan X GPU written in the paper, 45 frames per second are processed without batch processing. This is the level at which video processing is possible. Second, YOLO simultaneously predicts multiple bounding boxes for the entire image. This is a differentiated approach from DPM[2] and R-CNN. Because it detects the entire image, it also learns the surrounding information. This is a way to reduce background errors. Finally, YOLO achieves high performance when testing pictorial images after natural image learning because it does the learning through generalization of objects. However, despite the three advantages, there are disadvantages. Instead of detecting it quickly, it is less accurate than other models. Accuracy is called mAP.

3 Single Network for Single Detection

For detection, the YOLO neural network truncates the input image to $S \times S$ grid. For confidence scores, grid cells predict from B bounding boxes. For their information, a score that indicates how accurate the bounding box is and exactly contains the object that they want to detect is called confidence score. The following definitions are defined: $Pr(\text{object}) * IOU_{\text{pred}}^{\text{truth}}$

IOU is the same one used in the previous experiment, meaning Intersection over Union. This refers to the intersection of the ground truth bounding box of an object and the predicted bounding box.

$$IOU = \frac{\text{Real Bounding Box} \cap \text{Predicted Bounding Box}}{\text{Real Bounding Box} \cup \text{Predicted Bounding Box}} \quad (1)$$

$Pr(\text{object}) = 0$ if there is no object in the grid cell. Therefore, confidence score = 0. If there is an object in the grid cell, $Pr(\text{object}) = 1$. It is ideal if the confidence score is equal to the IOU. The bounding box consists of x, y, w, h , and confidence. The relative position within the bounding box is represented by $0 \leq (x, y) \leq 1$. The width and height of the bounding box are called $0 \leq (w, h) \leq 1$ when the width and height of the input image are considered as 1. Use this parameter to predict the conditional class probabilities of the grid cell. The equation is as follows: $C(\text{conditional class probabilities}) = Pr(\text{Class}_i | \text{Object})$

Obtain 1 class probabilities per grid cell regardless of the bounding box in the grid cell. For reference, one grid cell predicts B bounding boxes. In the test, the equation (4) can obtain a class specific confidence score.

$$\begin{aligned} &\text{Class specific Confidence Score} \\ &= Pr(\text{Class}_i | \text{Object}) \cdot Pr(\text{Object}) \cdot IOU_{\text{pred}}^{\text{truth}} \\ &= Pr(\text{Class}_i) \cdot IOU_{\text{pred}}^{\text{truth}} \end{aligned} \quad (2)$$

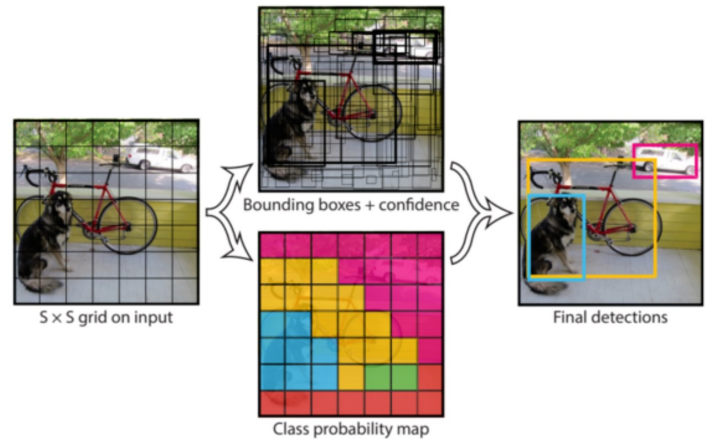


Figure 2: The researchers experimented with the PASCAL VOC dataset[1][1], which was set to $S = 7, B = 5$, with a total of 20 labeling classes present. That is, $C = 20$. Divide one image into seven grids and predict two bounding boxes from one grid cell. The size of the final tensor is $S \times S \times (B * 5 + C)$.

Network Design The layer in front of the YOLO model is the convolutional layer. The fully-connected layer is located behind the model. Perform feature vector of images through convolutional layers, and predict the coordinates of class probabilities and bounding boxes through fully-connected layers. This structure is a form of fine tuning on GoogLeNet[8], consisting of 24 convolutional layers and two fully-connected layers. The difference from GoogLeNet used 1×1 reduction layer and 3×3 convolutional layer[6]. the researchers can see that the final output is $7 \times 7 \times 30$.

6 Reference

- [1] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111(1):98–136, 2015.
- [2] Pedro F Felzenszwalb, Ross B Girshick, David McAllester, and Deva Ramanan. Object detection with discriminatively trained part-based models. *IEEE transactions on pattern analysis and machine intelligence*, 32(9):1627–1645, 2009.
- [3] Shiry Ginosar, Daniel Haas, Timothy Brown, and Jitendra Malik. Detecting people in cubist art. In *European Conference on Computer Vision*, pages 101–116. Springer, 2014.
- [4] RB Girshick. Fast r-cnn. corr, abs/1504.08083, 2015.
- [5] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014.
- [6] Min Lin, Qiang Chen, and Shuicheng Yan. Network in network. corr abs/1312.4400 (2013). *arXiv preprint arXiv:1312.4400*, 2013.
- [7] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- [8] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott E Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. corr abs/1409.4842 (2014). *arXiv preprint arXiv:1409.4842*, 2014.