

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/343194514>

A Survey on Performance Metrics for Object-Detection Algorithms

Conference Paper · July 2020

DOI: 10.1109/IWSSIP48289.2020

CITATIONS
113

READS
14,976

3 authors:



Rafael Padilla
Federal University of Rio de Janeiro

11 PUBLICATIONS 476 CITATIONS

[SEE PROFILE](#)



Sergio Lima Netto
Federal University of Rio de Janeiro

158 PUBLICATIONS 1,661 CITATIONS

[SEE PROFILE](#)



Eduardo A. B. da Silva
Federal University of Rio de Janeiro

295 PUBLICATIONS 3,089 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



DORIS project [View project](#)



D.Sc. Thesis: Algorithms for New Scenarios on Communications [View project](#)

A Survey on Performance Metrics for Object-Detection Algorithms

Rafael Padilla¹, Sergio L. Netto², Eduardo A. B. da Silva³
^{1,2,3}PEE, COPPE, Federal University of Rio de Janeiro, P.O. Box 68504, RJ, 21945-970, Brazil
{rafael.padilla, sergioln, eduardo}@smt.ufrj.br

Abstract—This work explores and compares the plethora of metrics for the performance evaluation of object-detection algorithms. Average precision (AP), for instance, is a popular metric for evaluating the accuracy of object detectors by estimating the area under the curve (AUC) of the precision \times recall relationship. Depending on the point interpolation used in the plot, two different AP variants can be defined and, therefore, different results are generated. AP has six additional variants increasing the possibilities of benchmarking. The lack of consensus in different works and AP implementations is a problem faced by the academic and scientific communities. Metric implementations written in different computational languages and platforms are usually distributed with corresponding datasets sharing a given bounding-box description. Such projects indeed help the community with evaluation tools, but demand extra work to be adapted for other datasets and bounding-box formats. This work reviews the most used metrics for object detection detaching their differences, applications, and main concepts. It also proposes a standard implementation that can be used as a benchmark among different datasets with minimum adaptation on the annotation files.

Keywords—object-detection metrics, average precision, object-detection challenges, bounding boxes.

I. INTRODUCTION

Object detection is an extensively studied topic in the field of computer vision. Different approaches have been employed to solve the growing need for accurate object detection models [1]. The Viola-Jones framework [2], for instance, became popular due to its successful application in the face-detection problem [3], and was later applied to different subtasks such as pedestrian [4] and car [5] detections. More recently, with the popularization of the convolutional neural networks (CNN) [6]–[9] and GPU-accelerated deep-learning frameworks, object-detection algorithms started being developed from a new perspective [10], [11]. Works as Overfeat [12], R-CNN [13], Fast R-CNN [14], Faster R-CNN [15], R-FCN [16], SSD [17] and YOLO [18]–[20] highly increased the performance standards on the field. World famous competitions such as VOC PASCAL Challenge [21], COCO [22], ImageNet Object Detection Challenge [23], and Google Open Images Challenge [24] have as their top object-detection algorithms methods inspired on the aforementioned works. Differently from algorithms such as the Viola-Jones, CNN-based detectors are flexible enough to be trained with several (hundreds or even a few thousands) classes.

A detector outcome is commonly composed of a list of bounding boxes, confidence levels and classes, as seen in Figure 1. However, the standard output-file format varies a lot for different detection algorithms. Bounding-box detections

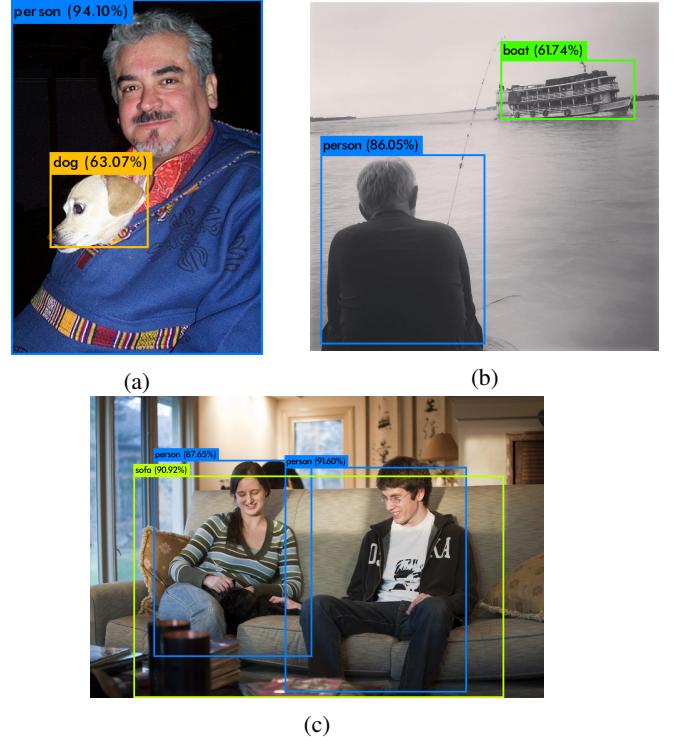


Fig. 1: Examples of detections performed by YOLO [20] in different datasets. (a) PASCAL VOC; (b) personal dataset; (c) COCO. Besides the bounding box coordinates of a detected object, the output also includes the confidence level and its class.

are mostly represented by their top-left and bottom-right coordinates $(x_{\text{ini}}, y_{\text{ini}}, x_{\text{end}}, y_{\text{end}})$, with a notable exception being the YOLO [18]–[20] algorithm, that differs from the others by outlining the bounding boxes by their center coordinates, width, and height ($\frac{x_{\text{center}}}{\text{image width}}$, $\frac{y_{\text{center}}}{\text{image height}}$, $\frac{\text{box width}}{\text{image width}}$, $\frac{\text{box height}}{\text{image height}}$).

Different challenges, competitions, and hackathons [21], [23]–[27] attempt to assess the performance of object detections in specific scenarios by using real-world annotated images [28]–[30]. In these events, participants are given a testing nonannotated image set in which objects have to be detected by their proposed works. Some competitions provide their own (or 3rd-party) source code, allowing the participants to evaluate their algorithms in an annotated validation image set before submitting their testing-set detections. In the end,

each team sends a list of bounding-boxes coordinates with their respective classes and (sometimes) their confidence levels to be evaluated.

In most competitions, the average precision (AP) and its derivations are the metrics adopted to assess the detections and thus rank the teams. The PASCAL VOC dataset [31] and challenge [21] provide their own source code to measure the AP and the mean AP (mAP) over all object classes. The City Intelligence Hackathon [27] uses the source code distributed in [32] to rank the participants also on AP and mAP. The ImageNet Object Localization challenge [23] does not recommend any code to compute their evaluation metric, but provides a pseudo-code explaining it. The Open Images 2019 [24] and Google AI Open Images [26] challenges use mAP, referencing a tool to evaluate the results [33], [34]. The Lyft 3D Object Detection for Autonomous Vehicles challenge [25] does not reference any external tool, but uses the AP averaged over 10 different thresholds, the so-called AP@50:5:95 metric.

This work reviews the most popular metrics used to evaluate object-detection algorithms, including their main concepts, pointing out their differences, and establishing a comparison between different implementations. In order to introduce its main contributions, this work is divided into the following topics: Section II explains the main performance metrics employed in the field of object detection and how the AP metric can produce ambiguous results; Section III describes some of the most known object detection challenges and their employed performance metrics, whereas Section IV presents a project implementing the AP metric to be used with any annotation format.

II. MAIN PERFORMANCE METRICS

Among different annotated datasets used by object detection challenges and the scientific community, the most common metric used to measure the accuracy of the detections is the AP. Before examining the variations of the AP, we should review some concepts that are shared among them. The most basic are the ones defined below:

- True positive (TP): A correct detection of a ground-truth bounding box;
- False positive (FP): An incorrect detection of a nonexistent object or a misplaced detection of an existing object;
- False negative (FN): An undetected ground-truth bounding box;

It is important to note that, in the object detection context, a true negative (TN) result does not apply, as there are infinite number of bounding boxes that should not be detected within any given image.

The above definitions require the establishment of what a “correct detection” and an “incorrect detection” are. A common way to do so is using the intersection over union (IOU). It is a measurement based on the Jaccard Index, a coefficient of similarity for two sets of data [35]. In the object detection scope, the IOU measures the overlapping area between the

predicted bounding box B_p and the ground-truth bounding box B_{gt} divided by the area of union between them, that is

$$J(B_p, B_{gt}) = \text{IOU} = \frac{\text{area}(B_p \cap B_{gt})}{\text{area}(B_p \cup B_{gt})}, \quad (1)$$

as illustrated in Figure 2.

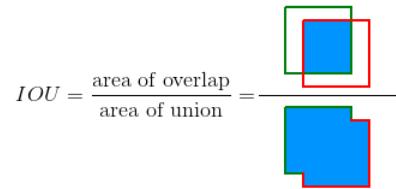


Fig. 2: Intersection Over Union (IOU).

By comparing the IOU with a given threshold t , we can classify a detection as being correct or incorrect. If $\text{IOU} \geq t$ then the detection is considered as correct. If $\text{IOU} < t$ the detection is considered as incorrect.

Since, as stated above, the true negatives (TN) are not used in object detection frameworks, one refrains to use any metric that is based on the TN, such as the TPR, FPR and ROC curves [36]. Instead, the assessment of object detection methods is mostly based on the precision P and recall R concepts, respectively defined as

$$P = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{\text{TP}}{\text{all detections}}, \quad (2)$$

$$R = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{\text{TP}}{\text{all ground truths}}. \quad (3)$$

Precision is the ability of a model to identify only relevant objects. It is the percentage of correct positive predictions. Recall is the ability of a model to find all relevant cases (all ground-truth bounding boxes). It is the percentage of correct positive predictions among all given ground truths.

The precision \times recall curve can be seen as a trade-off between precision and recall for different confidence values associated to the bounding boxes generated by a detector. If the confidence of a detector is such that its FP is low, the precision will be high. However, in this case, many positives may be missed, yielding a high FN, and thus a low recall. Conversely, if one accepts more positives, the recall will increase, but the FP may also increase, decreasing the precision. However, a good object detector should find all ground-truth objects ($FN = 0 \equiv$ high recall) while identifying only relevant objects ($FP = 0 \equiv$ high precision). Therefore, a particular object detector can be considered good if its precision stays high as its recall increases, which means that if the confidence threshold varies, the precision and recall will still be high. Hence, a high area under the curve (AUC) tends to indicate both high precision and high recall. Unfortunately, in practical cases, the precision \times recall plot is often a zigzag-like curve, posing challenges to an accurate measurement of its AUC. This is circumvented by processing the precision \times recall curve in order to remove the zigzag behavior prior to AUC estimation. There are basically

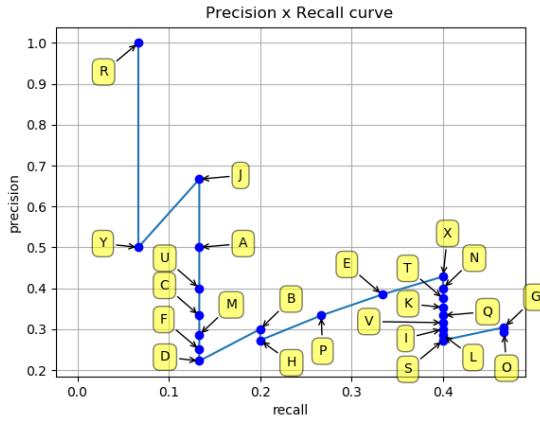


Fig. 4: Precision x Recall curve with values calculated for each detection in Table I.

and (Figure 6):

$$\begin{aligned} \text{AP}_{\text{all}} &= 1 * (0.0666 - 0) + 0.6666 * (0.1333 - 0.0666) \\ &\quad + 0.4285 * (0.4 - 0.1333) + 0.3043 * (0.4666 - 0.4) \\ \text{AP}_{\text{all}} &= 24.56\%. \end{aligned}$$

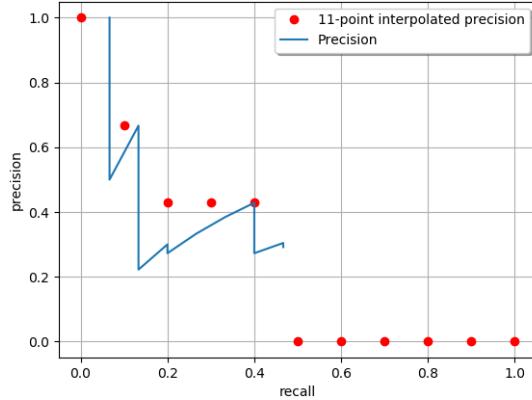


Fig. 5: Precision × Recall curves of points from Table I using the 11-point interpolation approach.

From what we have seen so far, benchmarks are not truly comparable if the method used to calculate the AP is not reported. Works found in the literature [1], [9], [12]–[20], [37] usually neither mention the method used nor reference the adopted tool to evaluate their results. This problem does not occur much often in challenges, as it is a common practice to have a reference software tool included in order for the participants to evaluate their results. Also, it is not rare to occur cases where a detector sets the same confidence level for different detections. Table I, for example, illustrates that detections R and Y obtained the same confidence level (95%). Depending on the criterion used by a certain implementation, one or other detection can be sorted as the first detection in the table, directly affecting the final result of an object-detection algorithm. Some implementations may consider the order that

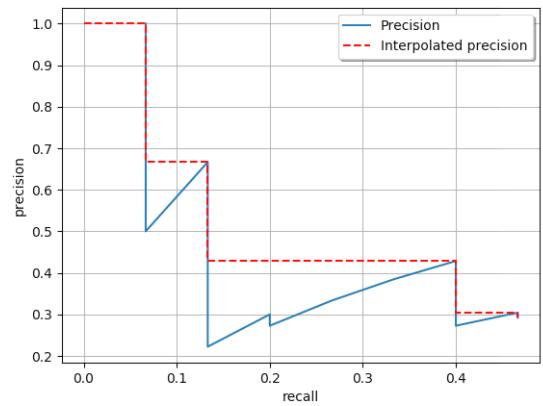


Fig. 6: Precision × Recall curves of points from Table I applying interpolation with all points.

each detection was reported as the tiebreaker (usually one or more evaluation files contain the detections to be evaluated), but in general there is no common consensus by the evaluation tools.

III. OBJECT-DETECTION CHALLENGES AND THEIR AP VARIANTS

Constantly, new techniques are being developed and new different state-of-the-art object-detection algorithms are arising. Comparing their results with different works is not an easy task. Sometimes the applied metrics vary or the implementation used by the different authors may not be the same, generating dissimilar results. This section covers the main challenges and their most popular AP variants found in the literature.

The PASCAL VOC [31] is an object-detection challenge released in 2005. From 2005 to 2012, a new version of the Pascal VOC was released with increased numbers of images and classes, starting at four classes, reaching 20 classes in its last update. The PASCAL VOC competition still accepts submissions, revealing state-of-the-art algorithms for object detections ever since. In this trial, the challenge applies the 11-point interpolated precision (see Section II) and uses the mean AP over all of its classes to rank the submission performances, as implemented by the provided development kit.

The Open Images 2019 challenge [24] in its object-detection track uses the Open Images Dataset [29] containing 12.2 M annotated bounding boxes across 500 object categories on 1.7 M images. Due to its hierarchical annotations, the same object can belong to a main class and multiple sub-classes (e.g. ‘helmet’ and ‘football helmet’). Because of that, the users should report the class and sub-classes of a given detection. If somehow only the main class is correctly reported for a detected bounding box, the unreported sub-classes affect negatively the score, as it is counted as a false negative. The metric employed by the aforementioned challenge is the mean AP over all classes using the Tensorflow Object Detection API [33].

The COCO detection challenge (bounding box) [22] is a competition which provides bounding-box coordinates of more than 200,000 images comprising 80 object categories. The

