

## Object Detection in Deep Learning -I

Ujjwal

INRIA

19<sup>th</sup> Jan of 2021

## 1. Introduction

What is Object Detection ?

## 2. Evaluation Metrics for Object Detection

## 3. Components of Object Detection

## 4. Object Detectors

$$f : \mathcal{X} \longrightarrow \mathcal{Y}$$

- ▶  $\mathcal{X} \subset \mathbb{R}^{m \times n}$  : Set of images.
- ▶  $\mathcal{Y} \subset T^M$  : Target set where  $M$  is the maximum possible number of detections in  $\mathcal{X}$ .  $T = C \times B$  where,
  - ▶  $C = [0, 1]^{N+1}$  : Confidence set of detections where  $N$  is total number of object classes.  $+1$  is for the background class.
    - ▶  $\tilde{c}_i \in C$  denotes the confidence of  $i^{\text{th}}$  class being present.
  - ▶  $B = \mathbb{R}^4$  : Bounding box set of detections.
- ▶  $C \times B$  denotes that object detection is a unification of classification and localization in machine learning.

## Compositional View

A mapping  $f$  depicting object detection can be written as a composition :

- ▶  $f = f_2 \circ f_1 \triangleq f_2(f_1(.))$
- ▶  $f_1 : \mathcal{X} \rightarrow C$ 
  - ▶ **Classification:** Which objects are present in  $\mathcal{X}$  ?
- ▶  $f_2 : \mathcal{X} \rightarrow B$ 
  - ▶ **Regression:** Where are the objects present in  $\mathcal{X}$  ?

## 1. Introduction

## 2. Evaluation Metrics for Object Detection

- Evaluation Fundamentals

- Precision and Recall

- Interpreting Object Detector Output

- Precision Recall Curve

## 3. Components of Object Detection

## 4. Object Detectors



## 2. Evaluation Metrics for Object Detection

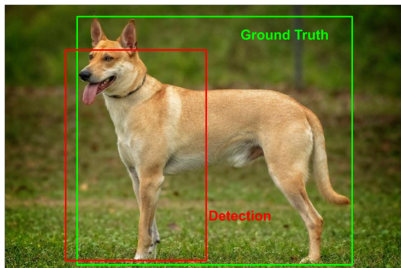
### a) Evaluation Fundamentals

# Fundamentals

---

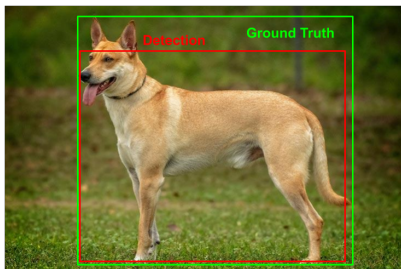
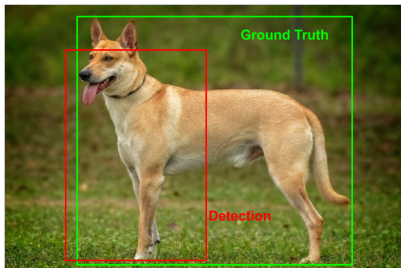
Ground Truth	Detector Output	
	Present	Absent
Present	True Positive (TP)	False Negative (FN)
Absent	False Positive (FP)	True Negative (TN)

Detection result interpretation based on an object's presence in **Ground Truth** (GT) and **Detector output**.

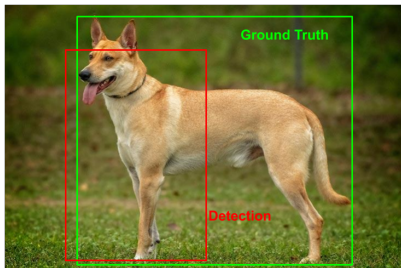


We consider two examples to demonstrate the relevance of localization in object detection.



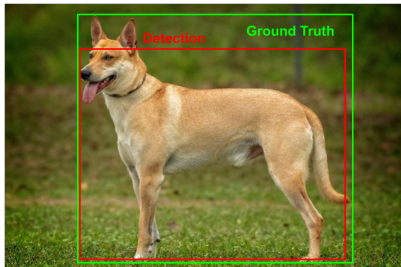


We assume that a bounding box with a label of “Dog” has been detected.



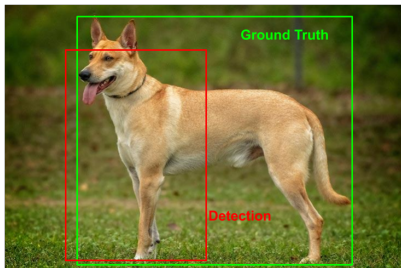
### Bad Detection

Detection overlap with GT is low.

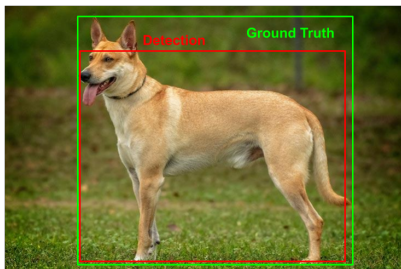


### Good Detection

Detection overlap with GT is high.



Dog is not detected. (FP)



Dog is detected. (TP)

Given  $B_1, B_2 \subset \mathbb{R}^4$ ,

- ▶ Let,  $I = \text{Area}(B_1 \cap B_2)$
- ▶ Let,  $U = \text{Area}(B_1 \cup B_2)$
- ▶  $\text{IoU} \triangleq \frac{I}{U}$
- ▶  $\text{IoU} \in [0, 1]$

$\text{IoU} \propto$  Agreement with the GT.

Given  $B_1, B_2 \subset \mathbb{R}^4$ ,

- ▶ Let,  $I = \text{Area}(B_1 \cap B_2)$
- ▶ Let,  $U = \text{Area}(B_1 \cup B_2)$
- ▶  $\text{IoU} \triangleq \frac{I}{U}$
- ▶  $\text{IoU} \in [0, 1]$

$\text{IoU} \propto$  Agreement with the GT.

$$\textit{Precision} = \frac{TP}{TP + FP} \quad (1)$$

$$Precision = \frac{TP}{TP + FP} \quad (1)$$

Of all detected objects how many are actually matching with GT

$$\textit{Precision} = \frac{TP}{TP + FP} \quad (1)$$

Of all detected objects how many are actually matching with GT

$$\textit{Recall} = \frac{TP}{TP + FN} \quad (2)$$



$$Precision = \frac{TP}{TP + FP} \quad (1)$$

Of all detected objects how many are actually matching with GT

$$Recall = \frac{TP}{TP + FN} \quad (2)$$

Of all GT objects how many are actually detected.

**High Precision** and **high recall** are desirable traits of an object detector.

Given an image  $x \in \mathcal{X}$ ,  $f(x) = \{(c_i, b_i)\}_{i=1}^M$ ,  $c_i \in C$  and  $b_i \in B$ .

Given an image  $x \in \mathcal{X}$ ,  $f(x) = \{(c_i, b_i)\}_{i=1}^M$ ,  $c_i \in C$  and  $b_i \in B$ .

$c_i \in [0, 1]^{N+1}$  denotes the probability of different object classes enclosed by  $b_i$ .

Given an image  $x \in \mathcal{X}$ ,  $f(x) = \{(c_i, b_i)\}_{i=1}^M$ ,  $c_i \in \mathcal{C}$  and  $b_i \in B$ .

$c_i \in [0, 1]^{N+1}$  denotes the probability of different object classes enclosed by  $b_i$ .

The object class enclosed by  $b_i$  is determined by

$$\mathcal{C} = \operatorname{argmax}_j \{o_j \in c_i\}$$

Given an image  $x \in \mathcal{X}$ ,  $f(x) = \{(c_i, b_i)\}_{i=1}^M$ ,  $c_i \in \mathcal{C}$  and  $b_i \in B$ .

$c_i \in [0, 1]^{N+1}$  denotes the probability of different object classes enclosed by  $b_i$ .

The object class enclosed by  $b_i$  is determined by

$$\mathcal{C} = \operatorname{argmax}_j \{o_j \in c_i\}$$

The object class enclosed by a bounding box is the class whose probability of presence is highest.

Given an image  $x \in \mathcal{X}$ , the detector output  $f(x)$  can be interpreted as follows

$$f(x) = \{(\phi_x(i), p_{\phi(i)}, b_i)\}_{i=1}^M$$

- ▶  $M$  objects have been detected.
- ▶  $\phi_x : 1 \leq i \in \mathbb{N} \leq M \longrightarrow 1 \leq j \in \mathbb{N} \leq N$ .
  - ▶  $\phi_x(i)$  is a mapping which assigns to  $i^{\text{th}}$  element of the output set for image  $x$ , a class label.
  - ▶  $\phi_x$  in general is a many-to-one mapping.
- ▶  $p_{\phi(i)}$  is the probability of the  $i^{\text{th}}$  object to have the class  $\phi(i)$ .
- ▶  $b_i$  is the bounding box representing the  $i^{\text{th}}$  object.





---

**Algorithm 1:** Determining a detection as TP, FP or FN.

---

**Result:** out

Given,

1. An element  $d_i \in f(x) = \{(\phi_x(i), p_{\phi_x(i)}, b_i^{det})\}_{i=1}^M$  for an image  $x$ ;
2. GT set  $G = \{\phi_g(i), 1, b_i\}_{i=1}^{\sigma(x)}$ , where  $\sigma(x)$  is number of GT objects in image  $x$ .
3. Pre-defined IoU threshold  $T$ .

Compute  $\mathbf{l}_o = \{l : l = IOU(d_i, g_i), g_i \in G\}$  and  $\mathbf{j} = \text{argmax}_j \mathbf{l}_o$ .**if**  $\mathbf{l}_o(\mathbf{j}) < T$  **then**

| out=None

**else**| **if**  $class(d_i) = class(g_j)$  **then**| | **if**  $confidence(d_i) < 0.5$  **then**

| | | out = FN

| | **else**

| | | out = TP

| | **end**| **else**

| | out = FP

| **end****end**

---

## Algorithm 2: How to write algorithms

---

**Result:** Write here the result

Given,

1. Detection set  $\mathcal{D} = f(x) = \{(\phi(i), p_{\phi(i)}, b_i^{det})\}_{i=1}^M$  for an image  $x$ ;
2. GT set.
3. Pre-defined IoU threshold  $T$ .
4.  $\tau = 0.0$

**while**  $\tau \leq 1$  **do**

    select  $\tilde{\mathcal{D}} \subseteq \mathcal{D}$  s.t.  $p_{\phi(i)} \geq \tau \forall 1 \leq i \leq |\tilde{\mathcal{D}}|$ ;  
     $\forall \tilde{d} \in \tilde{\mathcal{D}}$ ; Find  $G \subseteq GT$ , such that  $IoU(\tilde{d}, g \in G) \geq T$ ;

**end**

---

1. Introduction
2. Evaluation Metrics for Object Detection
3. Components of Object Detection  
Object Detection Components
4. Object Detectors

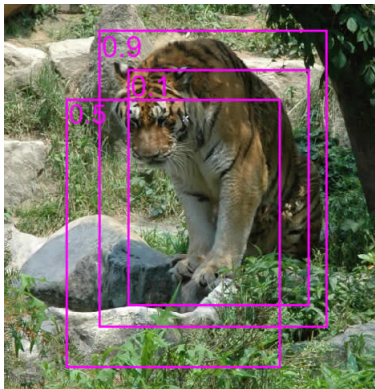
1. Pre-processor
2. Detection Head
3. Post-Processor
4. Loss Function ( During training phase )

Transforms an image into a feature map.

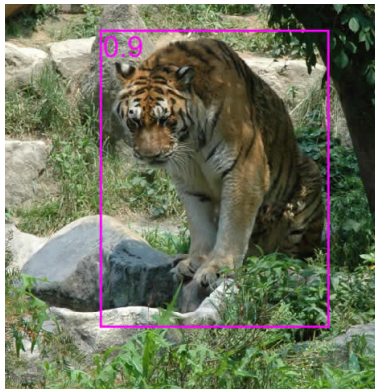
Detection head is made of two components:

1. **Classification Head:** Performs classification of a region in the feature map.
2. **Regression Head:** Performs bounding box regression over a region in the feature map.

Object detectors often give multiple detections for each object. Post-Processing aims at extracting the most relevant detection for each object.



Before post-processing.



After post-processing.

### 3. Components of Object Detection

#### a) Object Detection Components



Generally a detector loss function is of the following form:

$$f_{det} \triangleq f_{cls} + f_{reg} + f_{regularization} \quad (3)$$

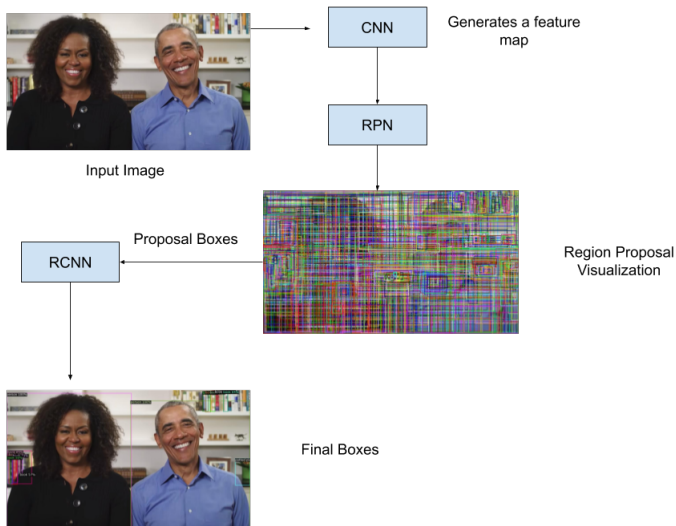
Some classification losses are:

1. Cross Entropy Loss.
2. Focal Loss.
3. Generalized Focal Loss.

Some regression losses are:

1. Smooth L1-Loss.
2. Repulsion Loss.

1. Introduction
2. Evaluation Metrics for Object Detection
3. Components of Object Detection
4. Object Detectors  
Faster-RCNN

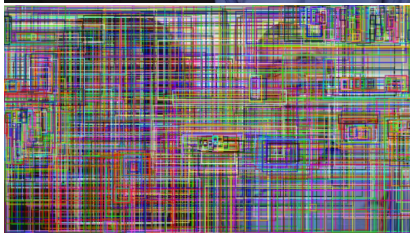


An input image is passed to a CNN to generate a feature map.

- ▶ Usually a CNN architecture trained for image classification is used.
  - ▶ This promotes faster convergence during training.
- ▶ Let  $m \times n$  be the size of an input image. If the CNN has an output stride of  $s$ , the spatial dimensions of the output feature map is  $\frac{m}{s} \times \frac{n}{s}$ .



Original Image



RPN Output

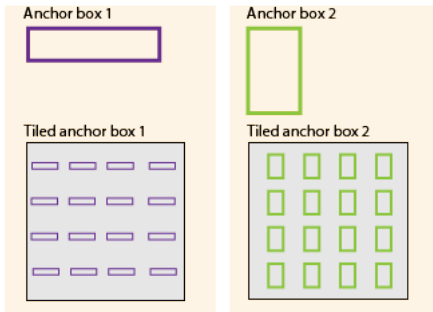
- ▶ Let us consider an image of size  $m \times n$ .
- ▶ Let us be provided with an information that all objects in the image have following constraints:
  - ▶ All objects' heights must be in the range  $[H_{min}, H_{max}]$ .
  - ▶ The aspect ratio  $r$  (i.e *height/width*) of every object must be in the range  $[r_{min}, r_{max}]$ .
- ▶ With these constraints, there is a very large number of possible regions in the image containing an object.



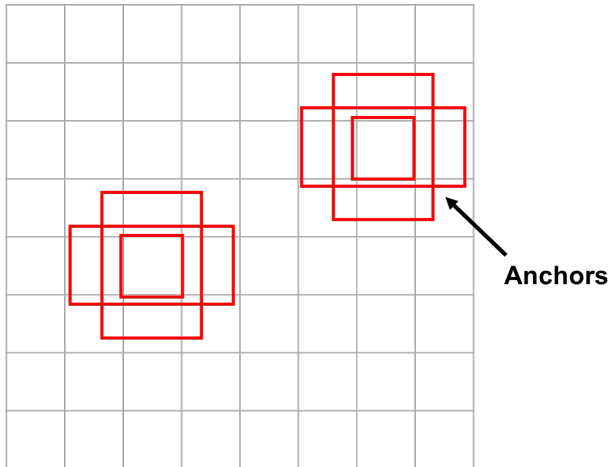
- ▶ Let an image be of size  $512 \times 512$ .
- ▶ For the height range let us consider an array  $H = [20 : 10 : 200]$  i.e minimum height is 20, maximum height is 200 with a step size of 10.
- ▶ For aspect ratio let us consider an array  $AR = [0.2 : 0.1 : 5]$ .
- ▶ Therefore for each location in the image, we have to search over  $len(H) \times len(AR) = 19 \times 49 = 931$  regions centered at it.
- ▶ Therefore total number of search regions is  $512 \times 512 \times 931 = 244056064$ .
- ▶ A RPN reduces the number of search regions.

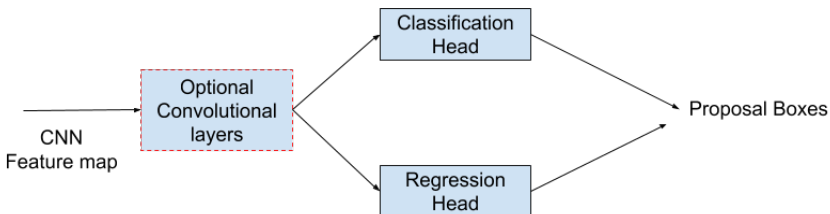
- ▶ Given an input image, a RPN provides a set of  $R$  bounding boxes known as proposals.
- ▶ Each bounding box denotes a region which the neural network thinks has a possibility of containing an object.
- ▶ This reduces the number of possible regions to be searched for objects.

- ▶ Anchors are hypothetical bounding boxes assumed to overlay a feature map.



Since à priori knowledge of object sizes is not available, a variety of anchors are assumed to overlay a feature map.





- ▶ Let us consider that there are  $R$  confocal anchors at each location.
- ▶ A classification head is a convolutional layer with  $2R$  filters.
  - ▶ The size of filters is an experimental hyperparameter. Usually  $3 \times 3$  or  $1 \times 1$  is used.
- ▶  $2R$  filters denote that for each anchor, the prediction has to be made across 2 object classes
  1. Object.
  2. Not an object.

- ▶ A regression head is a convolutional layer with  $4R$  filters.
  - ▶ The size of filters is an experimental hyperparameter. Usually  $3 \times 3$  or  $1 \times 1$  is used.
- ▶  $4R$  filters denote that for each anchor, the prediction should be for 4 quantities to predict the coordinates of the bounding box.

- ▶ During training, the groundtruth bounding boxes are known.
- ▶ the IoU is computed for each groundtruth bounding box with each anchor.
- ▶ The anchors for which  $IoU > \tau_u$  are meant to be positive anchors.
  - ▶ For positive anchors, they must be predicted as the “object” class and their regression target must be the coordinates of the groundtruth box.
- ▶ The anchors for which  $IoU < \tau_l$  are meant to be negative anchors.
  - ▶ For negative anchors, they must be predicted as the “not an object” class and their regression targets are not relevant. This will be come clearer when we discuss the loss function for RPN.
- ▶  $\tau_u$  and  $\tau_l$  are hyperparameters. In the original paper,  $\tau_u = 0.7$  and  $\tau_l = 0.3$ .
- ▶ For IoU in the range  $[0.3, 0.7]$ , the outputs are not considered during training.



$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i L_{cls}(p_i, p_i^*) \\ + \lambda \frac{1}{N_{reg}} \sum_i p_i^* L_{reg}(t_i, t_i^*).$$

- ▶ starred (\*) terms are groundtruth values. Non-starred are predicted values.
- ▶  $p_i$  is the classification head output for  $i^{th}$  anchor.
- ▶  $t_i$  is the regression head output for  $i^{th}$  anchor.
- ▶  $N_{cls}$  and  $N_{reg}$  are scalar quantities used for normalizing the values. We discuss this in the next slide.
- ▶  $\lambda$  is a regularization parameter to balance the importance of classification and regression. Usually both are equally important and so  $\lambda = 1$ .

- ▶ Total number of anchors is usually much larger than number of objects.
- ▶ As a result most anchors based on IoU values with GT boxes are negative.
- ▶ If the RPN loss function is minimized for all anchors, the negative anchors will dominate.
  - ▶ The network is learn about negative class i.e background, but will not be good at detecting objects.
- ▶ Hence, prior to loss function computation a set of anchors are chosen for which the loss function is minimized.

- ▶ We consider all anchors for which  $IoU > \tau_U$ .
  - ▶ This is because we want network to detect all objects. So, the network must learn about anchors which are positive.
- ▶ We sort the classification confidence score of all negative anchors and select those for which the classification score of being an object is highest.
  - ▶ This amounts to selecting the hardest possible examples.
- ▶ Usually the ratio of positive and negative anchors for loss function computation is considered as 1 : 3.

$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i L_{cls}(p_i, p_i^*) \\ + \lambda \frac{1}{N_{reg}} \sum_i p_i^* L_{reg}(t_i, t_i^*).$$

- ▶  $N_{cls}$  is total number of anchors considered for loss computation.
- ▶  $N_{reg}$  is total number of feature map locations i.e for a feature map of  $\frac{m}{s} \times \frac{n}{s}$ ,  $N_{reg}$  is  $\frac{mn}{s^2}$
- ▶  $p_i^*$  is 1 if the  $i^{th}$  anchor is positive else it is 0. So, regression is computed only for those anchors which are positive because for negative anchors we do not have a target i.e  $t_i^*$ .