

# Headline

---

大家好：

这是2018年度第2篇Arxiv Weekly。

本文是检测方向的文章。

## Highlight

---

利用IoU-Net提升NMS选框的性能，并利用新的方法微调选出的bbox最终提高检测性能。

## Information

---

Title

*Acquisition of Localization Confidence for Accurate Object Detection*

Link

<https://arxiv.org/pdf/1807.11590.pdf>

Source

- 北京大学 (Peking University)
- 清华大学 (Tsinghua University)
- 旷视科技 (Face++)
- 头条人工智能实验室 (Toutiao AI Lab)

## Introduction

---

现代CNN目标检测器中，一般都存在 Non-maximum suppression (NMS) 选框和bbox回归的过程。然而这里的逻辑是有一定的问题的。NMS操作的置信度本质来源于classification的label，而不是localization，这不可避免地导致了一些localization性能的降退，即使后续有进一步的回归。

本文提出了IoU-Net结构，能够预测每个候选框和对应的ground truth框之间的IoU。由于在选框的时候真正考虑了localization信息，IoU-Net能够优化NMS过程，给出更精准的预选框。另外，本文也提出了bbox refinement算法，来进一步提高最终框的精度。

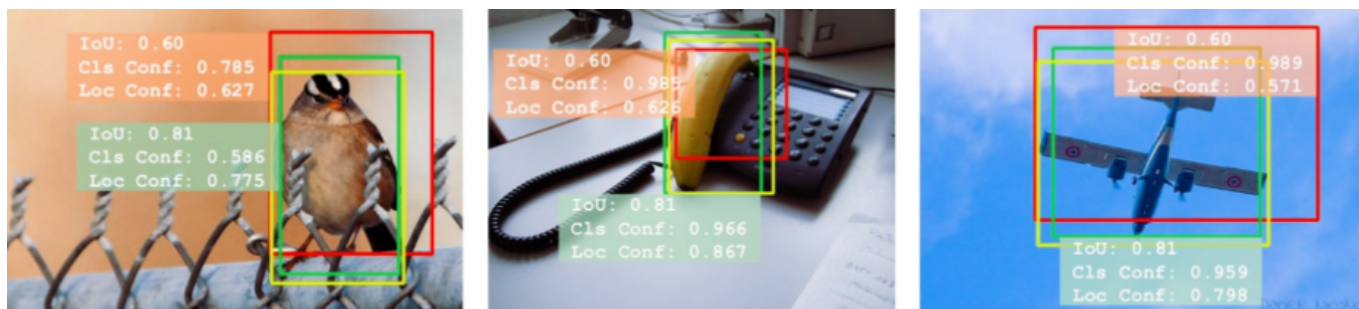
在MSCOCO上的实验表明本文提出的方法能在经典detection pipeline上取得涨点，并且具有良好的兼容性和可迁移性。

## Keys

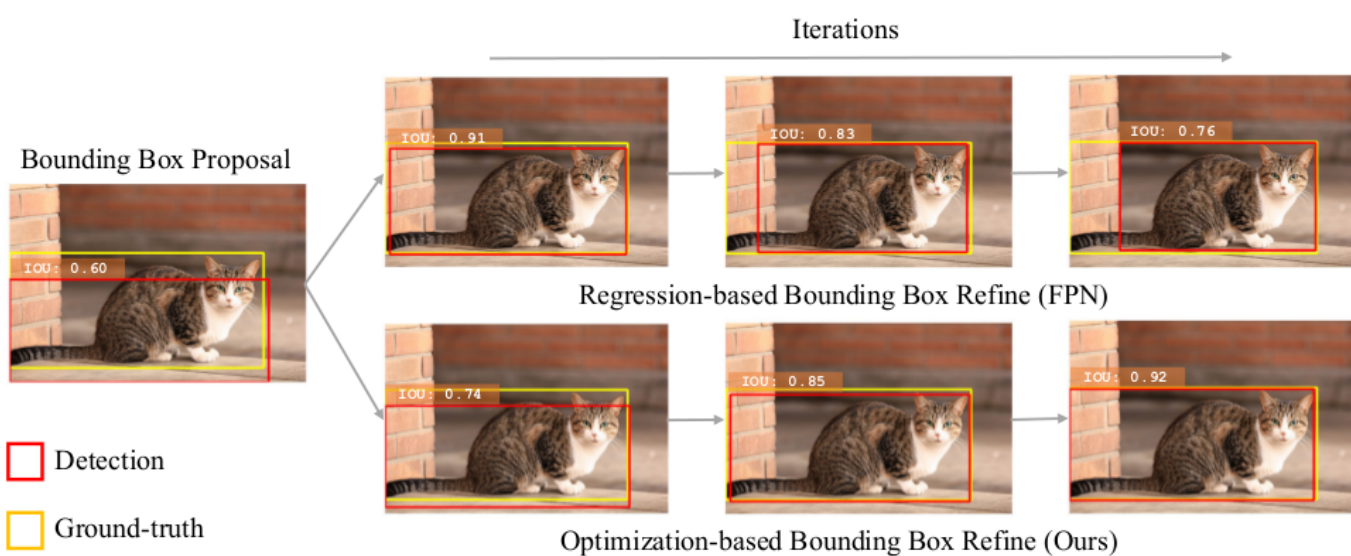
---

### 1.问题描述

- **misalignment between classification and localization accuracy.** [也即NMS会舍弃更好的框，挑选烂框来回归]



- **the absence of localization confidence makes bbox regression less interpretable.** [也即随着迭代的进行，bbox的回归结果反而恶化]



## 2. 解决方案part1：IoU-guided NMS

---

**Algorithm 1** IoU-guided NMS. Classification confidence and localization confidence are disentangled in the algorithm. We use the localization confidence (the predicted IoU) to rank all detected bounding boxes, and update the classification confidence based on a clustering-like rule.

---

**Input:**  $\mathcal{B} = \{b_1, \dots, b_n\}$ ,  $\mathcal{S}$ ,  $\mathcal{I}$ ,  $\Omega_{\text{nms}}$

$\mathcal{B}$  is a set of detected bounding boxes.

$\mathcal{S}$  and  $\mathcal{I}$  are functions (neural networks) mapping bounding boxes to their classification confidence and IoU estimation (localization confidence) respectively.

$\Omega_{\text{nms}}$  is the NMS threshold.

**Output:**  $\mathcal{D}$ , the set of detected bounding boxes with classification scores.

```

1:  $\mathcal{D} \leftarrow \emptyset$ 
2: while  $\mathcal{B} \neq \emptyset$  do
3:    $b_m \leftarrow \arg \max \mathcal{I}(b_j)$ 
4:    $\mathcal{B} \leftarrow \mathcal{B} \setminus \{b_m\}$ 
5:    $s \leftarrow \mathcal{S}(b_m)$ 
6:   for  $b_j \in \mathcal{B}$  do
7:     if  $\text{IoU}(b_m, b_j) > \Omega_{\text{nms}}$  then
8:        $s \leftarrow \max(s, \mathcal{S}(b_j))$ 
9:        $\mathcal{B} \leftarrow \mathcal{B} \setminus \{b_j\}$ 
10:    end if
11:  end for
12:   $\mathcal{D} \leftarrow \mathcal{D} \cup \{(b_m, s)\}$ 
13: end while
14: return  $\mathcal{D}$ 

```

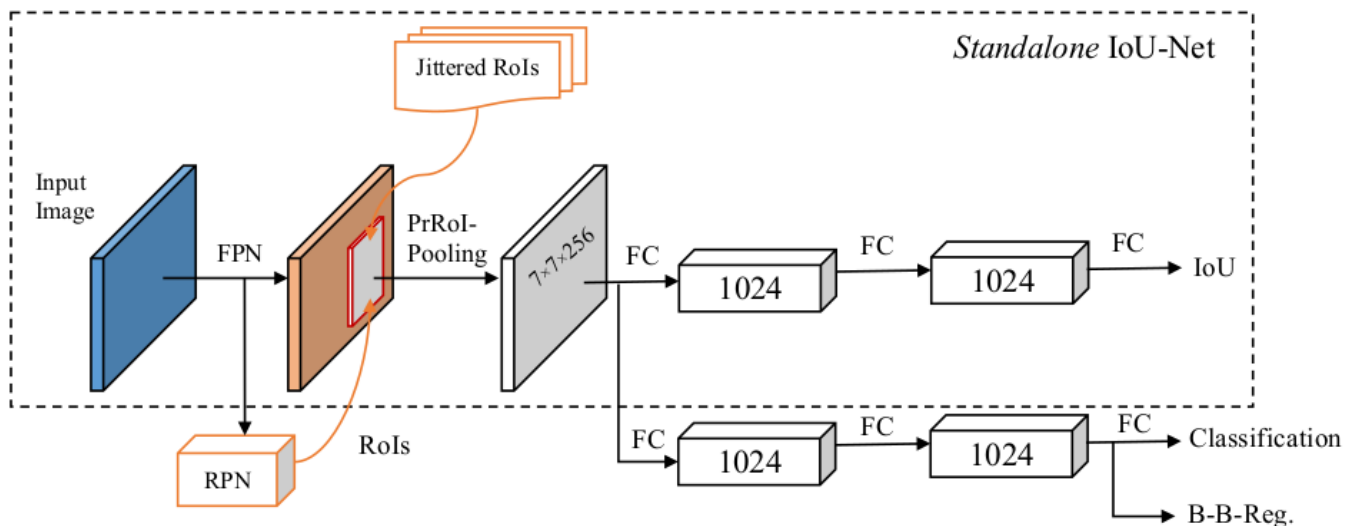
---

此为文中提出的新NMS算法，其核心为利用bbox对应区域和真值的IoU作为挑选maximum bbox的依据，而不是像传统的NMS一样直接用classification score来作为依据，这就弥合了misalignment between classification and localization accuracy。

另外值得一提的是，IoU-guided NMS算法中，在抑制那些和最大框交集超过阈值的框的同时，会参考它们的分类信息。如果它们的分类置信度高于我们选出的最大框，则会更新最大框的分类。[这个操作中，分离classification和localization的意味也很重，甚至有点localization指导优化classification的意思在里面。不过个人感觉这种情况并不常见]

剩下的一个关键问题是，给出一个bbox，我们怎么知道它和真值的bbox的IoU是多少.....

文章中采用了如下的IoU-Net网络来解决这个问题。



网络的训练数据是通过ground truth + augmentation & randomization生成的。

### 3.解决方案part2：optimization based bounding box refinement

---

#### Algorithm 2 Optimization-based bounding box refinement

---

**Input:**  $\mathcal{B} = \{b_1, \dots, b_n\}$ ,  $\mathcal{F}$ ,  $T$ ,  $\lambda$ ,  $\Omega_1$ ,  $\Omega_2$   
 $\mathcal{B}$  is a set of detected bounding boxes, in the form of  $(x_0, y_0, x_1, y_1)$ .  
 $\mathcal{F}$  is the feature map of the input image.  
 $T$  is number of steps.  $\lambda$  is the step size, and  $\Omega_1$  is an early-stop threshold and  $\Omega_2 < 0$  is an localization degeneration tolerance.  
Function PrPool extracts the feature representation for a given bounding box and function IoU denotes the estimation of IoU by the IoU-Net.  
**Output:** The set of final detection bounding boxes.

```

1:  $\mathcal{A} \leftarrow \emptyset$ 
2: for  $i = 1$  to  $T$  do
3:   for  $b_j \in \mathcal{B}$  and  $b_j \notin \mathcal{A}$  do
4:      $\mathbf{grad} \leftarrow \nabla_{b_j} \text{IoU}(\text{PrPool}(\mathcal{F}, b_j))$ 
5:      $\text{PrevScore} \leftarrow \text{IoU}(\text{PrPool}(\mathcal{F}, b_j))$ 
6:      $b_j \leftarrow b_j + \lambda * \text{scale}(\mathbf{grad}, b_j)$ 
7:      $\text{NewScore} \leftarrow \text{IoU}(\text{PrPool}(\mathcal{F}, b_j))$ 
8:     if  $|\text{PrevScore} - \text{NewScore}| < \Omega_1$  or  $\text{NewScore} - \text{PrevScore} < \Omega_2$  then
9:        $\mathcal{A} \leftarrow \mathcal{A} \cup \{b_j\}$ 
10:    end if
11:  end for
12: end for
13: return  $\mathcal{B}$ 

```

---

核心的思路在于利用IoU Net对bbox的IoU预测能力，指导后续的bbox微调。微调的时候用最简单的梯度上升优化方法，并设定两个阈值，一个用来判定收敛；一个用来判定是否已经进入localization degeneration的阶段并予以遏制。

p.s. 我们可以看到，不同的refinement方法，都是在寻找一种方式解决如下的最优化问题：

$$c^* = \arg \min_c \text{crit}(\text{transform}(box_{\text{det}}, c), box_{\text{gt}}), \quad (1)$$

其中det是检测网络输出的bbox；gt是真值bbox；c是网络微调器transformer的参数；crit是loss函数形式[一般选取smooth-L1 distance]

### 4.关于precise ROI pooling

为了能够实现3中的算法，我们需要给出一个可导的RoI Pooling算法。RoI Align算法已经比较好地解决了RoI中misalignment的问题，因此其基本想法是应该沿用的。所需要解决的是导数求解的问题。

因此本文提出了精准RoI Pooling (PrRoI Pooling)算法，其实质就是把feature map上的点利用双线性插值(bilinear interpolation)算法转化为连续的函数，如下公式所示：

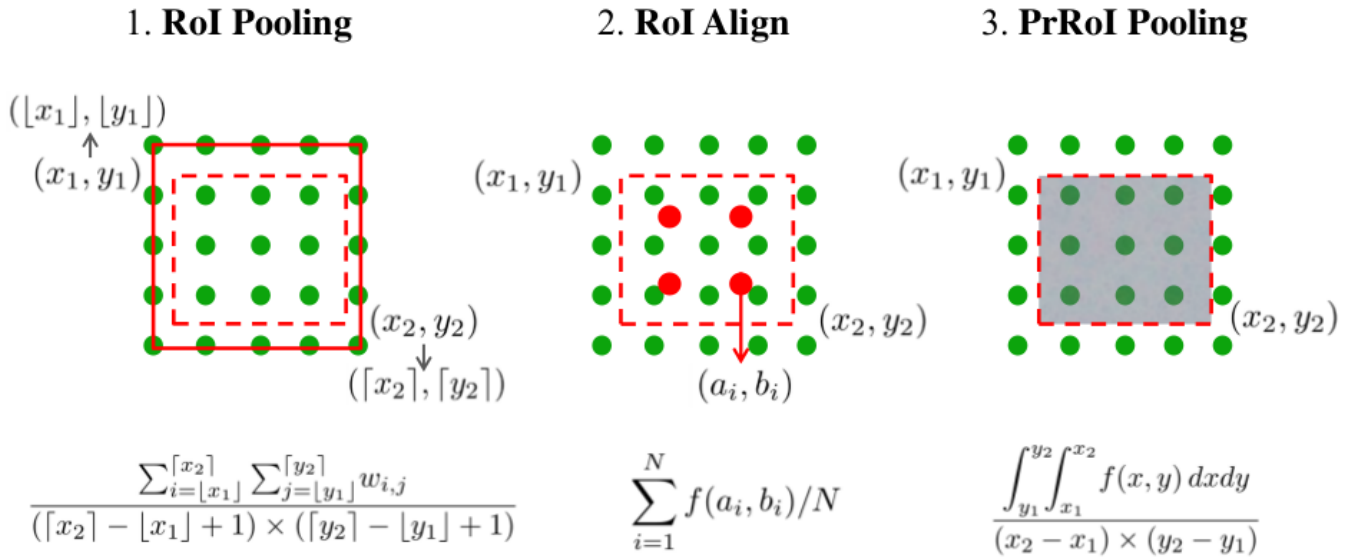
$$f(x, y) = \sum_{i,j} IC(x, y, i, j) \times w_{i,j}, \quad (2)$$

where  $IC(x, y, i, j) = \max(0, 1 - |x - i|) \times \max(0, 1 - |y - j|)$  is the interpolation coefficient. Then denote a bin of a RoI as  $bin = \{(x_1, y_1), (x_2, y_2)\}$ , where  $(x_1, y_1)$

有了连续的feature map值，就能够把bbox以浮点数的形式直接套在feature map对应的位置上，并且进行pooling操作了。[不过由于连续性的操作，我理解这里的PrRoI Pooling更加像一种average pooling的变体]

$$\text{PrPool}(bin, \mathcal{F}) = \frac{\int_{y_1}^{y_2} \int_{x_1}^{x_2} f(x, y) dx dy}{(x_2 - x_1) \times (y_2 - y_1)}. \quad (3)$$

PrRoI Pooling和两种经典Pooling算法的对比示意图如下：



## Results

所有实验均在MSCOCO上进行。

传统NMS、Soft-NMS、IoU-NMS的对比里面，高阈值AP下IoU-NMS表现最为突出，这是因为IoU-NMS能够最大程度地优化具体的边界位置。

Method	+Soft-NMS	+IoU-NMS	AP	AP <sub>50</sub>	AP <sub>60</sub>	AP <sub>70</sub>	AP <sub>80</sub>	AP <sub>90</sub>
FPN	✓	✓	36.4	<b>58.0</b>	<b>53.1</b>	44.9	31.2	9.8
			36.8	57.5	<b>53.1</b>	<b>45.7</b>	32.3	10.3
			<b>37.3</b>	56.0	52.2	45.6	<b>33.9</b>	<b>13.3</b>
Cascade R-CNN	✓	✓	40.6	<b>59.3</b>	55.2	49.1	38.7	16.7
			<b>40.9</b>	58.2	<b>54.7</b>	<b>49.4</b>	<b>39.9</b>	17.8
			40.7	58.0	<b>54.7</b>	49.2	38.8	<b>18.9</b>
Mask-RCNN	✓	✓	37.5	<b>58.6</b>	<b>53.9</b>	46.3	33.2	10.9
			37.9	58.2	<b>53.9</b>	<b>47.1</b>	34.4	11.5
			<b>38.1</b>	56.4	52.7	46.7	<b>35.1</b>	<b>14.6</b>

将refinement操作附加在主流的pipeline上，都取得了不错的涨点。同样地，高阈值下涨点突出很多。

Method	+Refinement	AP	AP <sub>50</sub>	AP <sub>60</sub>	AP <sub>70</sub>	AP <sub>80</sub>	AP <sub>90</sub>
FPN	✓	36.4	<b>58.0</b>	<b>53.1</b>	44.9	31.2	9.8
		<b>38.0</b>	57.7	<b>53.1</b>	<b>46.1</b>	<b>34.3</b>	<b>14.6</b>
Cascade R-CNN	✓	40.6	<b>59.3</b>	55.2	49.1	38.7	16.7
		<b>41.4</b>	<b>59.3</b>	<b>55.3</b>	<b>49.6</b>	<b>39.4</b>	<b>19.5</b>
Mask-RCNN	✓	37.5	<b>58.6</b>	<b>53.9</b>	46.3	33.2	10.9
		<b>39.2</b>	57.9	53.6	<b>47.4</b>	<b>36.5</b>	<b>16.4</b>

以上方法联合使用效果更佳。

Backbone	Method	+IoU-NMS	+Refine	AP	AP <sub>50</sub>	AP <sub>60</sub>	AP <sub>70</sub>	AP <sub>80</sub>	AP <sub>90</sub>
ResNet-50	FPN	✓	✓	36.4	58.0	53.1	44.9	31.2	9.8
	IoU-Net			37.0	<b>58.3</b>	<b>53.8</b>	45.7	31.9	10.7
				37.6	56.2	52.4	46.0	34.1	14.0
				<b>38.1</b>	56.3	52.4	<b>46.3</b>	<b>35.1</b>	<b>15.5</b>
ResNet-101	FPN	✓	✓	38.5	<b>60.3</b>	<b>55.5</b>	47.6	33.8	11.3
	IoU-Net			38.9	60.2	<b>55.5</b>	47.8	34.6	12.0
				40.0	59.0	55.1	48.6	37.0	15.5
				<b>40.6</b>	59.0	55.2	<b>49.0</b>	<b>38.0</b>	<b>17.1</b>

## Insights

针对检测中通用的NMS模块提出了新的改进，并且提出了在最后对bbox进行refinement操作的想法。比较好地缓解了分类和定位标准误用以及随着训练迭代bbox反而退化的问题。

另外，提出了对ROI Pooling求导数的一种思路：PrRoI Pooling。