

Arbitrary-Oriented Object Detection with Circular Smooth Label

Xue Yang^{1,2}, Junchi Yan^{1,2}

¹ Department of Computer Science and Engineering, Shanghai Jiao Tong University
² MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University
yangxue-2019-sjtu.edu.cn

Abstract. Arbitrary-oriented object detection has recently attracted increasing attention in vision for their importance in aerial imagery, scene text, and face etc. In this paper, we show that existing regression-based rotation detectors suffer the problem of **discontinuous boundaries**, which is directly caused by **angular periodicity** or **corner ordering**. By a careful study, we find the root cause is that the **ideal predictions** are beyond the **defined range**. We design a new rotation detection baseline, to address the boundary problem by **transforming angular prediction from a regression problem to a classification task** with little accuracy loss, whereby high-precision angle classification is devised in contrast to previous works using coarse-granularity in rotation detection. We also propose a **circular smooth label (CSL)** technique to handle the periodicity of the angle and increase the error tolerance to adjacent angles. We further introduce **four window functions** in CSL and explore the effect of **different window radius sizes** on detection performance. Extensive experiments and visual analysis on two large-scale public datasets for aerial images i.e. DOTA, HRSC2016, as well as scene text dataset ICDAR2015 and MLT, show the effectiveness of our approach. The code will be released at https://github.com/Thinklab-SJTU/CSL_RetinaNet_Tensorflow.

Keywords: Arbitrary-Oriented, Object Detection, Boundary Problem, Circular Smooth Label.

1 Introduction

Object detection is one of the fundamental tasks in computer vision. In particular, rotation detection has played a huge role in the field of aerial images [2, 4, 41, 42, 44], scene text [12, 18, 19, 24, 27, 49] and face [11, 33, 34]. The rotation detector can provide accurate orientation and scale information, which will be helpful in applications such as object change detection in aerial images and recognition of sequential characters for multi-oriented scene texts.

Recently, a line of advanced rotation detectors evolved from classic detection algorithms [3, 7, 20, 21, 32] have been proposed. Among these methods, detectors based on region regression occupy the mainstream, and the representation of multi-oriented object is achieved by rotated bounding box or quadrangles.

Although these rotation detectors have achieved promising results, there are still some fundamental problems. Specifically, we note both the five-parameter regression and the eight-parameter regression methods suffer the problem of discontinuous boundaries, as often caused by angular periodicity or corner ordering. However, the inherent reasons are not limited to the particular representation of the bounding box. In this paper, we argue that the root cause of boundary problems based on regression methods is that the ideal predictions are beyond the defined range. Thus, the model’s loss value at the boundary suddenly increase so that the model cannot obtain the prediction result in the simplest and most direct way, and additional more complicated treatment is often needed. Therefore, these detectors often have difficulty in boundary conditions. For detection using rotated bounding boxes, the accuracy of angle prediction is critical. A slight angle deviation leads to important Intersection-over-Union (IoU) drop, resulting in inaccurate object detection, especially in case of large aspect ratios.

There have been some works addressing the boundary problem. For example, IoU-smooth L1 [44] loss introduces the IoU factor, and modular rotation loss [30] increases the boundary constraint to eliminate the sudden increase in boundary loss and reduce the difficulty of model learning. However, these methods are still regression-based detection methods, and still have not solved the root cause as mentioned above.

In this paper, we are aimed to find a more fundamental rotation detection baseline to solve the boundary problem. Specifically, we consider the prediction of the object angle as a classification problem to better limit the prediction results, and then we design a circular smooth label (CSL) to address the periodicity of the angle and increase the error tolerance between adjacent angles. Although the conversion from continuous regression to discrete classification , the impact of the lost accuracy on the rotation detection task is negligible. We also introduce four window functions in CSL and explore the effect of different window radius sizes on detection performance. After a lot of experiments and visual analysis, we find that CSL-based rotation detection algorithm is indeed a better baseline choice than the angle regression-based method on different detectors and datasets. Note the regression-based and CSL-based methods mentioned in subsequent chapters are divided according to the prediction form of the angle.

In summary, the main contribution of this paper are four-folds:

- We summarize the boundary problems in different regression-based rotation detection methods [2, 4, 41, 42] and show the root cause is that the ideal predictions are beyond the defined range.
- We design a new rotation detection baseline, which transforms angular prediction from a regression problem to a classification problem. Specifically, to our best knowledge, we devise the first high-precision angle (less than 1 degree) classification based pipeline in rotation detection, in contrast to previous coarse classification granularity (around 10-degree) methods [33]. Our method has little accuracy loss compared with regression-based methods and can effectively eliminate the boundary problem.

- We also propose the circular smooth label (CSL) technique, as an independent module which can also be readily reused in existing regression based methods by replacing the regression with classification, to address angular prediction for boundary conditions and objects with large aspect ratio.
- Extensive experimental results on DOTA and HRSC2016 show the state-of-the-art performance of our detector, and the efficacy of our CSL technique as an independent component has been verified across different detectors.

2 Related Work

Horizontal region object detection. Classic object detection aims to detect general objects in images with horizontal bounding boxes, and many high-performance general-purpose object detections have been proposed. R-CNN [8] pioneers a method based on CNN detection. Subsequently, region-based models such as Fast R-CNN [7], Faster R-CNN [32], and R-FCN [3] are proposed, which improve the detection speed while reducing computational storage. FPN [20] focus on the scale variance of objects in images and propose feature pyramid network to handle objects at different scales. SSD [23], YOLO [31] and RetinaNet [21] are representative single-stage methods, and their single-stage structure allows them to have faster detection speeds. Compared to anchor-based methods, many anchor-free have become extremely popular in recent years. CornerNet [15], CenterNet [5] and ExtremeNet [48] attempt to predict some key-points of objects such as corners or extreme points, which are then grouped into bounding boxes. However, horizontal detector does not provide accurate orientation and scale information, which poses problem in real applications such as object change detection in aerial images and recognition of sequential characters for multi-oriented scene texts.

Arbitrary-oriented object detection. Aerial images and scene text are the main application scenarios of the rotation detector. Recent advances in multi-oriented object detection are mainly driven by adaption of classical object detection methods using rotated bounding boxes or quadrangles to represent multi-oriented objects. Due to the complexity of the remote sensing image scene and the large number of small, cluttered and rotated objects, multi-stage rotation detectors are still dominant for their robustness. Among them, ICN [2], ROI-Transformer [4], SCRDet [41], R³Det [41] are state-of-the-art detectors. Gliding Vertex [40] and RSDet [30] achieve more accurate object detection through quadrilateral regression prediction. For scene text detection, RRPN [27] employ rotated RPN to generate rotated proposals and further perform rotated bounding box regression. TextBoxes++ [18] adopts vertex regression on SSD. RRD [19] further improves TextBoxes++ by decoupling classification and bounding box regression on rotation-invariant and rotation sensitive features, respectively. Although the regression-based arbitrary-oriented object detection method occupies the mainstream, we have found that most of these methods have some boundary problems due to the situations beyond the defined range. Therefore, we design a new rotation detection baseline, which basically eliminates the bound-

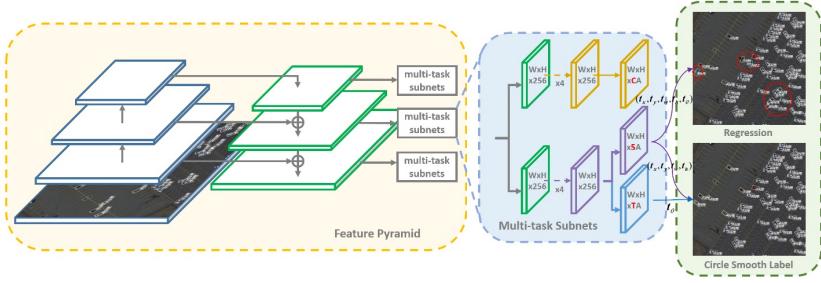


Fig. 1. Architecture of the proposed rotation detector (RetinaNet as an embodiment). ‘C’ and ‘T’ represent the number of object and angle categories, respectively.

ary problem by transforming angular prediction from a regression problem to a classification problem with little accuracy loss.

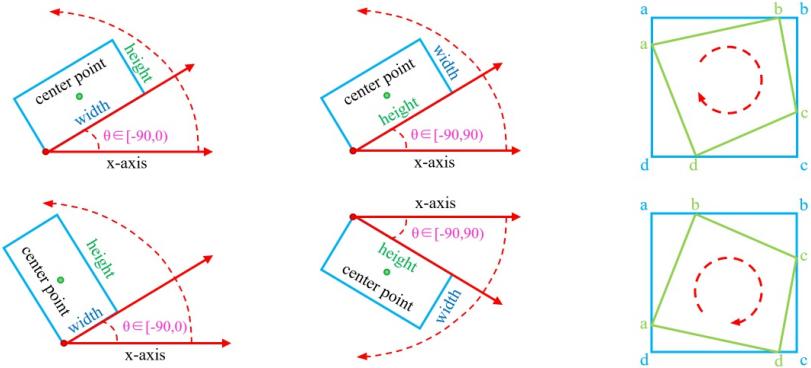
Classification for orientation information. The method of obtaining orientation information through classification is earlier used for multi-view face detection with arbitrary rotation-in-plane (RIP) angles. Divide-and-Conquer is adopted in [11], which use several small neural networks to deal with a small range of face appearance variations individually. In [33], a router network is firstly used to estimate each face candidates RIP angle. PCN [34] progressively calibrates the RIP orientation of each face candidate and shrinks the RIP range by half in early stages. Finally, PCN makes the accurate final decision for each face candidate to determine whether it is a face and predict the precise RIP angle. In other research areas, [14] adopts ordinal regression for or effective future motion classification. [43] obtains the orientation information of the ship by classifying the four sides. The above methods all obtain the approximate orientation range through classification, but cannot be directly applied to scenarios that require precise orientation information such as aerial images and scene text.

3 Proposed Method

We give an overview of our method as sketched in Figure 1. The embodiment is a single-stage rotation detector based on the RetinaNet [21]. The figure shows a multi-tasking pipeline, including regression-based prediction branch and CSL-based prediction branch, to facilitate the comparison of the performance of the two methods. It can be seen from the figure that CSL-based method is more accurate for learning the orientation and scale information of the object. It should be noted that the method proposed in this paper is applicable to most regression-based methods, which has been verified in the FPN [20] detector in experiments.

3.1 Regression-based Rotation Detection Method.

Parametric regression is currently a popular method for rotation object detection, mainly including five-parameter regression-based methods [4, 12, 27, 41,



(a) Five-parameter method with 90° angular range. (b) Five-parameter method with 180° angular range. (c) Ordered quadrilateral representation.

Fig. 2. Several definitions of bounding boxes.

42, 44] and eight-parameter regression-based methods [18, 25, 30, 40]. The commonly used five-parameter regression-based methods realize arbitrary-oriented bounding box detection by adding an additional angle parameter θ . Figure 2(a) shows one of the rectangular definition (x, y, w, h, θ) with 90° angular range [27, 41, 42, 44], θ denotes the acute angle to the x-axis, and for the other side we refer it as w . It should be distinguished from another definition (x, y, h, w, θ) illustrated in Figure 2(b), with 180° angular range [4, 27], whose θ is determined by the long side (h) of the rectangle and x-axis. The eight-parameter regression-based detectors directly regress the four corners $(x_1, y_1, x_2, y_2, x_3, y_3, x_4, y_4)$ of the object, so the prediction is a quadrilateral. The key step to the quadrilateral regression is to sort the four corner points in advance, which can avoid a very large loss even if the prediction is correct, as shown in Figure 2(c).

3.2 Boundary Problem of Regression Method.

Although the parametric regression-based rotation detection method has achieved competitive performance in different vision tasks, and has been a building block for a number of excellent detection methods, these methods essentially suffer the discontinuous boundaries problem [30, 44]. Boundary discontinuity problems are often caused by angular periodicity in the five-parameter method and corner ordering in the eight-parameter method, but there exist more fundamental root cause regardless the representation choices of the bounding box.

The boundary discontinuity problem often makes the model's loss value at the boundary suddenly increase. Thus methods have to resort to particular and often complex tricks to mitigate this issue. Therefore, these detection methods are often inaccurate in boundary conditions. We describe the boundary problem in three typical categories of regression-based methods according to their different representation forms (the first two refer to the five-parameter methods):

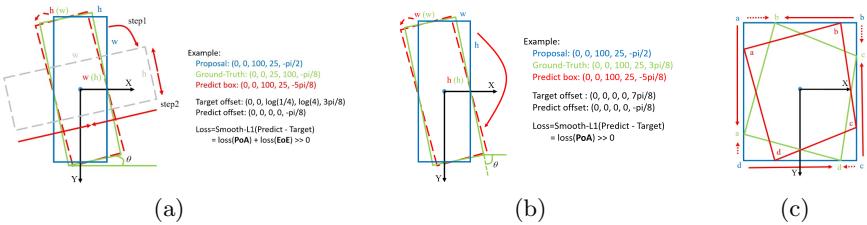


Fig. 3. The boundary problem of three categories of regression based methods. The red solid arrow indicates the actual regression process, and the red dotted arrow indicates the ideal regression process.

- **90°-regression-based method, as sketched in Figure 3(a).** It shows that an ideal form of regression (the blue box rotates counterclockwise to the red box), but the loss of this situation is very large due to the periodicity of angular (PoA) and exchangeability of edges (EoE), see the example in Figure 3(a) and Equation 3, 4, 5 for detail. Therefore, the model has to be regressed in other complex forms (such as the blue box rotating clockwise to the gray box while scaling w and h), increasing the difficulty of regression.
- **180°-regression-based method, as illustrated in Figure 3(b).** Similarly, this method also has a problem of sharp increase of loss caused by the PoA at the boundary. The model will eventually choose to rotate the proposal a large angle clockwise to get the final predicted bounding box.
- **Point-based method, as shown in Figure 3(c).** Through further analysis, the boundary discontinuity problem still exists in the eight-parameter regression method due to the advance ordering of corner points. Consider the situation of an eight-parameter regression in the boundary case, the ideal regression process should be $\{(a \rightarrow b), (b \rightarrow c), (c \rightarrow d), (d \rightarrow a)\}$, but the actual regression process from the blue reference box to the green ground truth box is $\{(a \rightarrow g), (b \rightarrow b), (c \rightarrow c), (d \rightarrow d)\}$. In fact, this situation also belongs to PoA. By contrast, the actual and ideal regression of the blue to red bounding boxes is consistent.

Some approaches have been proposed to solve these problems based on the above analysis. For example, IoU-smooth L1 [44] loss introduces the IoU factor, and modular rotation loss [30] increases the boundary constraint to eliminate the sudden increase in boundary loss and reduce the difficulty of model learning. However, these methods are still regression-based detection methods, and no solution is given from the root cause. In this paper, we will start from a new perspective and replace regression with classification to achieve better and more robust rotation detectors. We reproduce some classic rotation detectors based on regression and compare them visually under boundary conditions, as shown in Figure 4(a) to Figure 4(e). In contrast, CLS-based methods have no boundary problem, as shown in Figure 4(i).

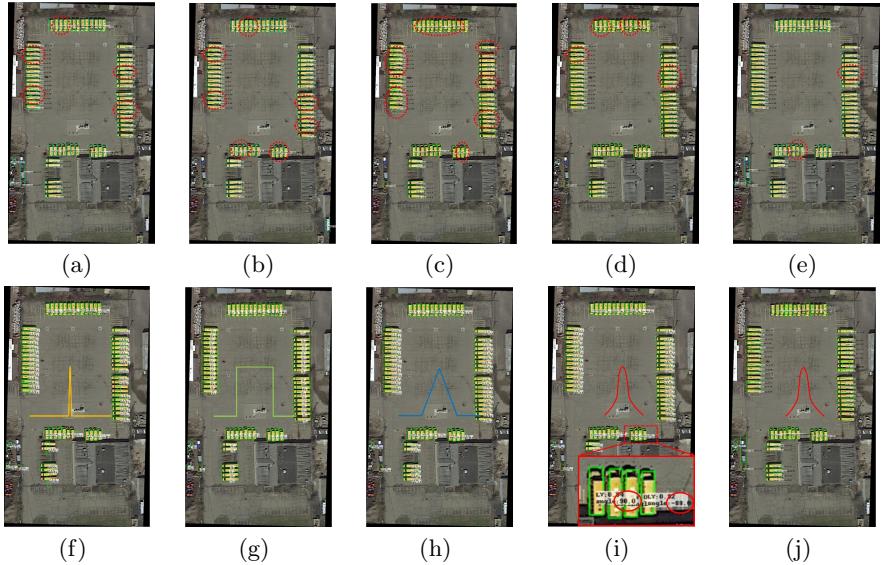


Fig. 4. Comparison of five regression-based rotation detection methods and CSL in the boundary case. (a) RetinaNet-H [41]. (b) RetinaNet-R [41]. (c) FPN-H. (d) R³Det [41]. (e) IoU-Smooth L1 [44]. (f) 180°-CSL-Pulse. (g) 180°-CSL-Rectangular. (h) 180°-CSL-Triangle. (i) 180°-CSL-Gaussian. (j) 90°-CSL-Gaussian. ‘H’ and ‘R’ represent the horizontal and rotating anchors. Red dotted circles indicate some bad cases.

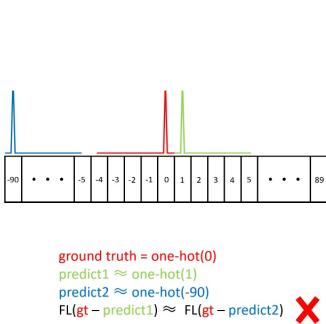
3.3 Circular Smooth Label for Angular Classification.

The main cause of boundary problems based on regression methods is that the ideal predictions are beyond the defined range. Therefore, we consider the prediction of the object angle as a classification problem to better limit the prediction results. A simple and straightforward solution is to use the object angle as its category label, and the number of categories is related to the angle range. Figure 5(a) shows the label setting for a standard classification problem (one-hot label encoding). The conversion from regression to classification can cause certain accuracy loss. Taking the five-parameter method with 180° angle range as an example, ω (default $\omega = 1^\circ$) degree per interval refers to a category. We can calculate the maximum accuracy loss $Max(loss)$ and the expected accuracy loss $E(loss)$:

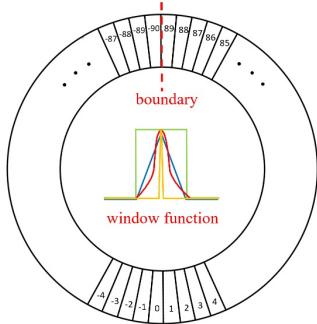
$$Max(loss) = \omega/2$$

$$E(loss) = \int_a^b x * \frac{1}{b-a} dx = \int_0^{\omega/2} x * \frac{1}{\omega/2 - 0} dx = \frac{\omega}{4} \quad (1)$$

Based on the above equations, one can see the loss is slight for a rotation detector. For example, when two rectangles with a 1 : 9 aspect ratio differ by 0.25° and 0.5° (default expected and maximum accuracy loss), the Intersection over Union (IoU) between them only decreases by 0.02 and 0.05. However, one-hot label has two drawbacks for rotation detection:



(a) One-hot label.



(b) Circle smooth label.

Fig. 5. Two kind of labels for angular classification. FL means focal loss [21].

- The EoE problem still exists when the bounding box uses the 90° -regression-based method. In addition, 90° -regression-based method has two different border cases (vertical and horizontal), while 180° -regression-based method has only vertical border cases.
- Note vanilla classification loss is agnostic to the angle distance between the predicted label and ground truth label, thus is inappropriate for the nature of the angle prediction problem. As shown in Figure 5(a), when the ground-truth is 0° and the prediction results of the classifier are 1° and -90° respectively, their prediction losses are the same, but the prediction results close to ground-truth should be allowed from a detection perspective.

Therefore, we design a circular smooth label (CSL) technique to obtain more robust angular prediction through classification without suffering boundary conditions, including EoE and PoA. It can be clearly seen from Figure 5(b) that CSL involves a circular label encoding with periodicity, and the assigned label value is smooth with a certain tolerance. The expression of CSL is as follows:

$$CSL(x) = \begin{cases} g(x), & \theta - r < x < \theta + r \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

where $g(x)$ is a window function. r is the radius of the window function. θ represents the angle of the current bounding box. An ideal window function $g(x)$ is required to hold the following properties:

- **Periodicity:** $g(x) = g(x + kT)$, $k \in N$. $T = 180/\omega$ represents the number of bins into which the angle is divided, and the default value is 180.
- **Symmetry:** $0 \leq g(\theta + \varepsilon) = g(\theta - \varepsilon) \leq 1$, $|\varepsilon| < r$. θ is the center of symmetry.
- **Maximum:** $g(\theta) = 1$.
- **Monotonic:** $0 \leq g(\theta \pm \varepsilon) \leq g(\theta \pm \varsigma) \leq 1$, $|\varsigma| < |\varepsilon| < r$. The function presents a monotonous non-increasing trend from the center point to both sides

We give four efficient window functions that meet the above three properties: pulse functions, rectangular functions, triangle functions, and Gaussian

functions, as shown in Figure 5(b). Note that the label value is continuous at the boundary and there is no arbitrary accuracy loss due to the periodicity of CSL. In addition, one-hot label is equivalent to CSL when the window function is a pulse function or the radius of the window function is very small.

3.4 Loss Function

Our multi-tasking pipeline contains regression-based prediction branch and CSL-based prediction branch, to facilitate the performance comparison of the two methods on an equal footing. The regression of the bounding box is:

$$\begin{aligned} t_x &= (x - x_a)/w_a, t_y = (y - y_a)/h_a \\ t_w &= \log(w/w_a), t_h = \log(h/h_a), \\ t_\theta &= (\theta - \theta_a) \cdot \pi/180 \quad (\text{only for regression branch}) \end{aligned} \quad (3)$$

$$\begin{aligned} t'_x &= (x' - x_a)/w_a, t'_y = (y' - y_a)/h_a \\ t'_w &= \log(w'/w_a), t'_h = \log(h'/h_a), \\ t'_\theta &= (\theta' - \theta_a) \cdot \pi/180 \quad (\text{only for regression branch}) \end{aligned} \quad (4)$$

where x, y, w, h, θ denote the box's center coordinates, width, height and angle, respectively. Variables x, x_a, x' are for the ground-truth box, anchor box, and predicted box, respectively (likewise for y, w, h, θ).

The multi-task loss is used which is defined as follows:

$$\begin{aligned} L &= \frac{\lambda_1}{N} \sum_{n=1}^N obj_n \cdot \sum_{j \in \{x, y, w, h, \theta_{reg}\}} L_{reg}(v'_{nj}, v_{nj}) \\ &\quad + \frac{\lambda_2}{N} \sum_{n=1}^N L_{CSL}(\theta'_n, \theta_n) + \frac{\lambda_3}{N} \sum_{n=1}^N L_{cls}(p_n, t_n) \end{aligned} \quad (5)$$

where N indicates the number of anchors, obj_n is a binary value ($obj_n = 1$ for foreground and $obj_n = 0$ for background, no regression for background). v'_{nj} denotes the predicted offset vectors, v_{nj} is the targets vector of ground-truth. θ_n, θ'_n denote the label and predict of angle respectively. t_n represents the label of object, p_n is the probability distribution of various classes calculated by Sigmoid function. The hyper-parameter $\lambda_1, \lambda_2, \lambda_3$ control the trade-off and are set to $\{1, 0.5, 1\}$ by default. The classification loss L_{cls} and L_{CSL} is focal loss [21] or sigmoid cross-entropy loss depend on detector. The regression loss L_{reg} is smooth L1 loss as used in [7].

4 Experiments

We use Tensorflow [1] to implement the proposed methods on a server with GeForce RTX 2080 Ti and 11G memory. The experiments in this article are

initialized by ResNet50 [10] by default unless otherwise specified. We perform experiments on both aerial benchmarks and scene text benchmarks to verify the generality of our techniques. Weight decay and momentum are set 0.0001 and 0.9, respectively. We employ MomentumOptimizer over 4 GPUs with a total of 4 images per minibatch (1 images per GPU). At each pyramid level we use anchors at seven aspect ratios $\{1, 1/2, 2, 1/4, 4, 1/6, 6\}$, and the remaining anchor settings are the same as the original RetinaNet and FPN.

4.1 Benchmarks and Protocols

DOTA [39] is one of the largest aerial image detection benchmarks. There are two detection tasks for DOTA: horizontal bounding boxes (HBB) and oriented bounding boxes (OBB). DOTA contains 2,806 aerial images from different sensors and platforms and the size of image ranges from around 800×800 to $4,000 \times 4,000$ pixels. The fully annotated DOTA benchmark contains 15 common object categories and 188,282 instances, each of which is labeled by an arbitrary quadrilateral. The short names for categories are defined as (abbreviation-full name): Plane (PL), Baseball diamond (BD), Bridge (BR), Ground field track (GTF), Small vehicle (SV), Large vehicle (LV), Ship (SH), Tennis court (TC), Basketball court (BC), Storage tank (ST), Soccer-ball field (SBF), Roundabout (RA), Harbor (HA), Swimming pool (SP), and Helicopter (HC). Half of the original images are randomly selected as the training set, 1/6 as the validation set, and 1/3 as the testing set. We divide the training and validation images into 600×600 subimages with an overlap of 150 pixels and scale it to 800×800 . With all these processes, we obtain about 27,000 patches.

ICDAR2015 [13] is the Challenge 4 of ICDAR 2015 Robust Reading Competition, which is commonly used for oriented scene text detection and spotting. This dataset includes 1,000 training images and 500 testing images. In training, we first train our model using 9,000 images from ICDAR 2017 MLT training and validation datasets, then we use 1,000 training images to fine-tune our model.

ICDAR 2017 MLT [28] is a multi-lingual text dataset, which includes 7,200 training images, 1,800 validation images and 9,000 testing images. The dataset is composed of complete scene images in 9 languages, and text regions in this dataset can be in arbitrary orientations, being more diverse and challenging.

HRSC2016 [26] contains images from two scenarios including ships on sea and ships close inshore. All images are collected from six famous harbors. The training, validation and test set include 436, 181 and 444 images, respectively.

All datasets are trained by 20 epochs (the number of image iterations per epoch is e) in total, and learning rate was reduced tenfold at 12 epochs and 16 epochs, respectively. The initial learning rates for RetinaNet and FPN are 5e-4 and 1e-3 respectively. The value of e for DOTA, ICDAR2015, MLT and HRSC2016 are 27k, 10k, 10k and 5k, and doubled if data augmentation and multi-scale training are used.

Table 1. Comparison of four window functions on the DOTA dataset. 5-mAP refers to the mean average precision of the five categories with large aspect ratio. mAP means mean average precision of all 15 categories. EoE indicates the issue of exchangeability of edges and a tick in table means the method suffers from EoE. All the methods are free of periodicity of angular (PoA) issues.

Based Method	Angle Range	EoE	Label Mode	BR	SV	LV	SH	HA	5-mAP	mAP
RetinaNet-H (CSL-Based)	90	✓	Pulse	9.80	28.04	11.42	18.43	23.35	18.21	39.52
	90	✓	Rectangular	37.62	54.28	48.97	62.59	50.26	50.74	58.86
	90	✓	Triangle	37.25	54.45	44.01	60.03	52.20	49.59	60.15
	90	✓	Gaussian	41.03	59.63	52.57	64.56	54.64	54.49	63.51
	180		Pulse	13.95	16.79	6.50	16.80	22.48	15.30	42.06
	180		Rectangular	36.14	60.80	50.01	65.75	53.17	53.17	61.98
	180		Triangle	32.69	47.25	44.39	54.11	41.90	44.07	57.94
	180		Gaussian	41.16	63.68	55.44	65.85	55.23	56.21	64.50

Table 2. Comparison of detection results under different radius.

Based Method	Angle Range	Label Mode	r=0	r=2	r=4	r=6	r=8
RetinaNet-H(CSL-Based)	180	Gaussian	40.78	59.23	62.12	64.50	63.99
FPN-H(CSL-Based)	180	Gaussian	48.08	70.18	70.09	70.92	69.75

4.2 Ablation Study

Comparison of four window functions. Table 1 shows the performance comparison of the four window functions on the DOTA dataset. It also details the accuracy of the five categories with larger aspect ratio and more border cases in the dataset. We believe that these categories can better reflect the advantages of our method. In general, the Gaussian window function performs best, while the pulse function performs worst because it has not learned any orientation and scale information. Figures 4(f)-4(i) show the visualization of the four window functions. According to Figure 4(i)-4(j), the 180°-CSL-based method obviously has better boundary prediction due to the EoE problem still exists in the 90°-CSL-based method. The visualization results in Figure 4 are consistent with the data analysis results in Table 1.

Suitable window radius. The Gaussian window form has shown best performance, while here we study the effect of radius of the window function. When the radius is too small, the window function tends to a pulse function. Conversely, the discrimination of all predictable results becomes smaller when the radius is too large. Therefore, we choose a suitable radius range from 0 to 8, Table 2 shows the performance of the two detectors in this range. Although both detectors achieve the best performance with a radius of 6, the single-stage detection method is more sensitive to radius. We speculate that the instance-level feature extraction capability (like RoI Pooling [7] and RoI Align [9]) in the two-stage detector is stronger than the image-level in the single-stage detector. Therefore, the two-stage detection method can distinguish the difference between the two approaching angles. Figure 6 compares visualizations at different window raduis. When the radius is 0, the detector cannot learn any orientation and scale information, which is consistent with the performance of the pulse function above.

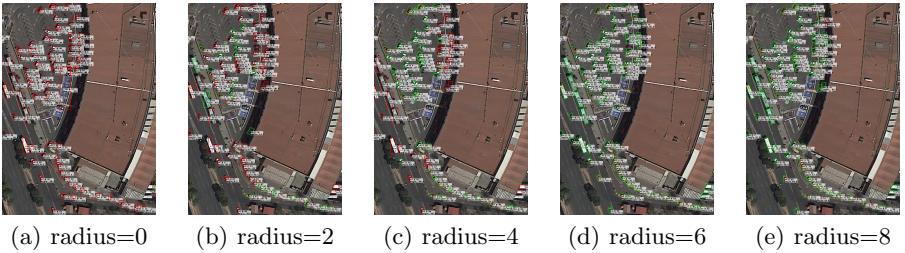


Fig. 6. Visualization of detection results (RetinaNet-H CSL-Based) under different radius. The red bounding box indicates that no orientation and scale information has been learned, and the green bounding box is the correct detection result.

Table 3. Comparison between CSL-based and regression-based methods on DOTA. Improvement by CSL-based methods have been made under the same configuration.

Based Method	Angle Range	Angle Pred.	PoA	EoE	Label Mode	BR	SV	LV	SH	HA	5-mAP	mAP
RetinaNet-H	90	regression-based	✓	✓	-	41.15	53.75	48.30	55.92	55.77	50.98	63.18
	90	CSL-based		✓	Gaussian	41.03	59.63	52.57	64.56	54.64	54.49 (+3.51)	63.51 (+0.33)
	180	regression-based	✓		-	38.47	54.15	47.89	60.87	53.63	51.00	64.10
	180	CSL-based			Gaussian	41.16	63.68	55.44	65.85	55.23	56.21 (+5.21)	64.50 (+0.40)
RetinaNet-R	90	regression-based	✓	✓	-	32.27	64.64	71.01	68.62	53.52	58.01	62.76
	90	CSL-based		✓	Gaussian	35.14	63.21	73.92	69.49	55.53	59.46 (+1.45)	65.45 (+2.69)
FPN-H	90	regression-based	✓	✓	-	44.78	70.25	71.13	68.80	54.27	61.85	68.25
	90	CSL-based		✓	Gaussian	45.46	70.22	71.96	76.06	54.84	63.71 (+1.86)	69.02 (+0.77)
	180	regression-based	✓		-	45.88	69.37	72.06	72.96	62.31	64.52	69.45
	180	CSL-based			Gaussian	47.90	69.66	74.30	77.06	64.59	66.70 (+2.18)	70.92 (+1.47)

As the radius becomes larger and optimal, the detector can learn the angle in any direction.

Classification is better than regression. Three rotation detectors in Table 3, including RetinaNet-H, RetinaNet-R and FPN-H, are used to compare the performance differences between CSL-based and regression-based methods. The former two are single-stage detectors, whose anchor format is different. The latter is a classic two-stage detection method. It can be clearly seen that CSL has better detection ability for objects with large aspect ratios and more boundary conditions. It also should be noted that CSL is designed to solve the boundary problem, whose proportion in the entire dataset is relatively small, so the overall performance (mAP) is not as obvious as the five categories listed (5-mAP). Overall, the CSL-based rotation detection algorithm is indeed a better baseline choice than the angle regression-based method.

CSL performance on other datasets. In order to further verify that CSL-based method is a better baseline model, we have also verified it in other datasets, including the text dataset ICDAR2015, MLT, and another remote sensing dataset HRSC2016. These three datasets are single-class object detection datasets, whose objects have a large aspect ratio. Although boundary conditions still account for a small proportion of these data sets, CSL still shows a stronger performance advantage. As shown in Table 4, the CSL-based method is improved by 1.21%, 0.56%, and 1.29% (1.4%) respectively compared with the

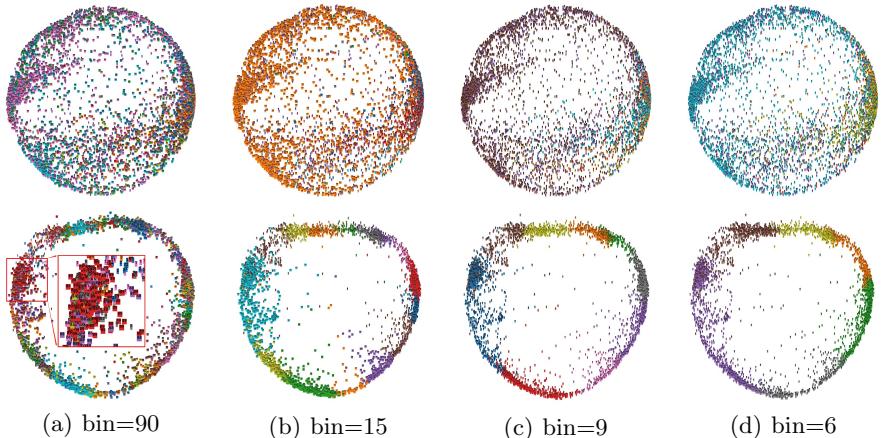


Fig. 7. Angular feature visualization of the 90-CSL-FPN detector on the DOTA dataset. First, we divide the entire angular range into several bins, and bins are different between columns. The two rows show two-dimensional feature visualizations of pulse and gaussian function, respectively. Each point represents a RoI of the test set with a index of the bin it belongs to.

Table 4. Comparison between CSL-based and regression-based methods on the text dataset ICDAR2015, MLT, and another remote sensing dataset HRSC2016. 07 or 12 means use the 2007 or 2012 evaluation metric.

Method	ICDAR2015			MLT			HRSC2016	
	Recall	Precision	Hmean	Recall	Precision	Hmean	mAP (07)	mAP (12)
FPN-regression-based	81.81	83.07	82.44	56.15	80.26	66.08	88.33	94.70
FPN-CSL-based	83.00	84.30	83.65 (+1.21)	56.72	80.77	66.64 (+0.56)	89.62 (+1.29)	96.10 (+1.40)

regression-based method under the same experimental configuration. These experimental results provide strong support for demonstrating the versatility of the CSL-based method.

Visual analysis of angular features. By zooming in on part of Figure 4(i), we find that the prediction of the boundary conditions became continuous (for example, two large vehicle in the same direction predicted 90° and -88° , respectively). This phenomenon reflects the purpose of designing the CSL: the labels are periodic (circular) and the prediction of adjacent angles has a certain tolerance. In order to confirm that the angle classifier has indeed learned this property, we visually analyze the angular features of each region of interest (RoI) in the FPN detector by principal component analysis (PCA) [38], as shown in Figure 7. The detector does not learn the orientation information of well when we use the pulse window function. It can be seen from the first row of Figure 7 that the feature distribution of RoI is relatively random, and the prediction results of some angles occupy the vast majority. For the gaussian function, the feature distribution is obvious a ring structures, and the features of adjacent angles are close to each other and have a certain overlap. It is this property

Table 5. Detection accuracy on each object (AP) and overall performance (mAP) on DOTA. Note O²-DNet uses Hourglass104 [29] as backbone.

Method	Backbone	PL	BD	BR	GTF	SV	LV	SH	TC	BC	ST	SBF	RA	HA	SP	HC	mAP
FR-O [39]	ResNet101	79.09	69.12	17.17	63.49	34.20	37.16	36.20	89.19	69.60	58.96	49.4	52.52	46.69	44.80	46.30	52.93
IENet [22]	ResNet101	64.54	39.82	32.07	49.71	65.01	52.58	81.45	44.66	78.51	46.54	56.73	64.40	64.24	36.75	57.14	
R-DFPN [42]	ResNet101	80.92	65.82	33.77	58.94	55.77	50.94	54.78	90.33	66.34	68.66	48.73	51.76	55.10	51.32	35.88	57.94
R ² CNN [12]	ResNet101	80.94	65.67	35.34	67.44	59.92	50.81	55.81	90.67	66.92	72.39	55.06	52.23	55.14	53.35	48.22	60.67
RRPN [27]	ResNet101	88.52	71.20	31.66	59.30	51.85	56.19	57.25	90.81	72.84	67.38	56.69	52.84	53.08	51.94	53.58	61.01
ICN [2]	ResNet101	81.40	74.30	47.70	70.30	64.90	67.80	70.00	90.80	79.10	78.20	53.60	62.90	67.00	64.20	50.20	68.20
RADet [17]	ResNeXt101	79.45	76.99	48.05	65.83	65.46	74.40	68.86	89.70	78.14	74.97	49.92	64.63	66.14	71.58	62.16	69.09
RoI-Transformer [4]	ResNet101	88.64	78.52	43.44	75.92	68.81	73.68	83.59	90.74	77.27	81.46	58.39	53.54	62.83	58.93	47.67	69.56
P-RSDet [47]	ResNet101	89.02	73.65	47.33	72.03	70.58	73.71	72.76	90.82	80.12	81.32	59.45	57.87	60.79	65.21	52.59	69.82
CAD-Net [45]	ResNet101	87.8	82.4	49.4	73.5	71.1	63.5	76.7	90.9	79.2	73.3	48.4	60.9	62.0	67.0	62.2	69.9
O ² -DNet [37]	Hourglass104	89.31	82.14	47.33	61.21	71.32	74.03	78.62	90.76	82.23	81.36	60.93	60.17	58.21	66.98	61.03	71.04
SCRDet [44]	ResNet101	89.98	80.65	52.09	68.36	68.36	60.32	72.41	90.85	87.94	86.86	65.02	66.68	66.25	68.24	65.21	72.61
SARD [36]	ResNet101	89.93	84.11	54.19	72.04	68.41	61.18	66.00	90.82	87.79	86.59	65.65	64.04	66.68	68.84	68.03	72.95
FADet [16]	ResNet101	90.21	79.58	45.49	76.41	73.18	68.27	79.56	90.83	83.40	84.68	53.40	65.42	74.17	69.69	64.86	73.28
R ³ Det [41]	ResNet152	89.49	81.17	50.53	66.10	70.92	78.66	78.21	90.81	85.26	84.23	61.81	63.77	68.16	69.83	67.17	73.74
RSDet [30]	ResNet152	90.1	82.0	53.8	68.5	70.2	78.7	73.6	91.2	87.1	84.7	64.3	68.2	66.1	69.3	63.7	74.1
Gliding Vertex [40]	ResNet101	89.64	85.00	52.26	77.34	73.01	73.14	86.82	90.74	79.02	86.81	59.55	70.91	72.94	70.86	57.32	75.02
Mask OBB [35]	ResNeXt-101	89.56	85.95	54.21	72.90	76.52	74.16	85.63	89.85	83.81	86.48	54.89	69.64	73.94	69.06	63.32	75.33
FFA [6]	ResNet101	90.1	82.7	54.2	75.2	71.0	79.9	83.5	90.7	83.9	84.6	61.2	68.0	70.7	76.0	63.7	75.7
APE [50]	ResNeXt-101	89.96	83.62	53.42	76.03	74.01	77.16	79.45	90.83	87.15	84.51	67.72	60.33	74.61	71.84	65.55	75.75
CSL (FPN based)	ResNet152	90.25	85.53	54.64	75.31	70.44	73.51	77.62	90.84	86.15	86.69	69.60	68.04	73.83	71.80	68.93	76.17

Table 6. Detection accuracy on HRSC2016 dataset.

Method	R ² CNN [12]	RC1 & RC2 [26]	RRPN [27]	R ² PN [46]	RetinaNet-H [41]	RRD [19]
mAP (07)	73.07	75.7	79.08	79.6	82.89	84.30
Method	RoI-Transformer [4]	RSDet [30]	Gliding Vertex [40]	RetinaNet-R [41]	R ³ Det [41]	FPN-CSL-based
mAP (07)	86.20	86.5	88.20	89.18	89.33	89.62

that helps CSL-based detectors to eliminate boundary problems and accurately obtain the orientation and scale information of the object.

4.3 Comparison with the State-of-the-Art

Results on DOTA. Although CSL is only a theoretical improvement on the original regression-based rotation detection method, it can still show competitive performance through data augmentation and multi-scale training and test that are widely used. We chose DOTA as the main validation dataset due to the complexity of the remote sensing image scene and the large number of small, cluttered and rotated objects. Our data augmentation methods mainly include random horizontal, vertical flipping, random graying, and random rotation. Training and testing scale set to [400, 600, 720, 800, 1000, 1100]. As shown in Table 5, FPN-CSL-based method shows competitive performance, at 76.17%.

Results on HRSC2016. The HRSC2016 contains lots of large aspect ratio ship instances with arbitrary orientation, which poses a huge challenge to the positioning accuracy of the detector. Experimental results show that our model achieves state-of-the-art performances, about 89.62%.

5 Conclusions

In this paper, we summarize **boundary problems** on different regression-based rotation detection methods. The main cause of boundary problems based on

理想的预测结果超出了定义域

regression methods is that the ideal predictions are beyond the defined range. Therefore, consider the prediction of the object angle as a classification problem to better limit the prediction results, and then we design a circular smooth label (CSL) to adapt to the periodicity of the angle and increase the tolerance of classification between adjacent angles with little accuracy loss. We also introduce four window functions in CSL and explore the effect of different window radius sizes on detection performance. Importantly, angle high-precision classification is also the first application in rotation detection. Extensive experiments and visual analysis on different detectors and datasets show that CSL-based rotation detection algorithm is indeed an effective baseline choice.

References

1. Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al.: Tensorflow: a system for large-scale machine learning. In: OSDI. vol. 16, pp. 265–283 (2016)
2. Azimi, S.M., Vig, E., Bahmanyar, R., Körner, M., Reinartz, P.: Towards multi-class object detection in unconstrained remote sensing imagery. In: Asian Conference on Computer Vision. pp. 150–165. Springer (2018)
3. Dai, J., Li, Y., He, K., Sun, J.: R-fcn: Object detection via region-based fully convolutional networks. In: Advances in neural information processing systems. pp. 379–387 (2016)
4. Ding, J., Xue, N., Long, Y., Xia, G.S., Lu, Q.: Learning roi transformer for oriented object detection in aerial images. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
5. Duan, K., Bai, S., Xie, L., Qi, H., Huang, Q., Tian, Q.: Centernet: Keypoint triplets for object detection. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 6569–6578 (2019)
6. Fu, K., Chang, Z., Zhang, Y., Xu, G., Zhang, K., Sun, X.: Rotation-aware and multi-scale convolutional neural network for object detection in remote sensing images. ISPRS Journal of Photogrammetry and Remote Sensing **161**, 294–308 (2020)
7. Girshick, R.: Fast r-cnn. In: Proceedings of the IEEE international conference on computer vision. pp. 1440–1448 (2015)
8. Girshick, R., Donahue, J., Darrell, T., Malik, J.: Rich feature hierarchies for accurate object detection and semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 580–587 (2014)
9. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: Proceedings of the IEEE international conference on computer vision. pp. 2961–2969 (2017)
10. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
11. Huang, C., Ai, H., Li, Y., Lao, S.: High-performance rotation invariant multiview face detection. IEEE Transactions on pattern analysis and machine intelligence **29**(4), 671–686 (2007)
12. Jiang, Y., Zhu, X., Wang, X., Yang, S., Li, W., Wang, H., Fu, P., Luo, Z.: R2cnn: rotational region cnn for orientation robust scene text detection. arXiv preprint arXiv:1706.09579 (2017)
13. Karatzas, D., Gomez-Bigorda, L., Nicolaou, A., Ghosh, S., Bagdanov, A., Iwamura, M., Matas, J., Neumann, L., Chandrasekhar, V.R., Lu, S., et al.: Icdar 2015 competition on robust reading. In: 2015 13th International Conference on Document Analysis and Recognition (ICDAR). pp. 1156–1160. IEEE (2015)
14. Kim, K.R., Choi, W., Koh, Y.J., Jeong, S.G., Kim, C.S.: Instance-level future motion estimation in a single image based on ordinal regression. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 273–282 (2019)
15. Law, H., Deng, J.: Cornernet: Detecting objects as paired keypoints. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 734–750 (2018)
16. Li, C., Xu, C., Cui, Z., Wang, D., Zhang, T., Yang, J.: Feature-attentioned object detection in remote sensing imagery. In: 2019 IEEE International Conference on Image Processing (ICIP). pp. 3886–3890. IEEE (2019)

17. Li, Y., Huang, Q., Pei, X., Jiao, L., Shang, R.: Radet: Refine feature pyramid network and multi-layer attention network for arbitrary-oriented object detection of remote sensing images. *Remote Sensing* **12**(3), 389 (2020)
18. Liao, M., Shi, B., Bai, X.: Textboxes++: A single-shot oriented scene text detector. *IEEE transactions on image processing* **27**(8), 3676–3690 (2018)
19. Liao, M., Zhu, Z., Shi, B., Xia, G.s., Bai, X.: Rotation-sensitive regression for oriented scene text detection. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 5909–5918 (2018)
20. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 2117–2125 (2017)
21. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: *Proceedings of the IEEE international conference on computer vision*. pp. 2980–2988 (2017)
22. Lin, Y., Feng, P., Guan, J.: Ienet: Interacting embranchment one stage anchor free detector for orientation aerial object detection. *arXiv preprint arXiv:1912.00969* (2019)
23. Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C.: Ssd: Single shot multibox detector. In: *European conference on computer vision*. pp. 21–37. Springer (2016)
24. Liu, X., Liang, D., Yan, S., Chen, D., Qiao, Y., Yan, J.: Fots: Fast oriented text spotting with a unified network. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 5676–5685 (2018)
25. Liu, Y., Zhang, S., Jin, L., Xie, L., Wu, Y., Wang, Z.: Omnidirectional scene text detection with sequential-free box discretization. *arXiv preprint arXiv:1906.02371* (2019)
26. Liu, Z., Yuan, L., Weng, L., Yang, Y.: A high resolution optical satellite image dataset for ship recognition and some new baselines. In: *Proc. ICPRAM*. vol. 2, pp. 324–331 (2017)
27. Ma, J., Shao, W., Ye, H., Wang, L., Wang, H., Zheng, Y., Xue, X.: Arbitrary-oriented scene text detection via rotation proposals. *IEEE Transactions on Multimedia* (2018)
28. Nayef, N., Yin, F., Bizid, I., Choi, H., Feng, Y., Karatzas, D., Luo, Z., Pal, U., Rigaud, C., Chazalon, J., et al.: Icdar2017 robust reading challenge on multi-lingual scene text detection and script identification-rrc-mlt. In: *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*. vol. 1, pp. 1454–1459. IEEE (2017)
29. Newell, A., Yang, K., Deng, J.: Stacked hourglass networks for human pose estimation. In: *European conference on computer vision*. pp. 483–499. Springer (2016)
30. Qian, W., Yang, X., Peng, S., Guo, Y., Yan, C.: Learning modulated loss for rotated object detection. *arXiv preprint arXiv:1911.08299* (2019)
31. Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: Unified, real-time object detection. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 779–788 (2016)
32. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis & Machine Intelligence* **(6)**, 1137–1149 (2017)
33. Rowley, H.A., Baluja, S., Kanade, T.: Rotation invariant neural network-based face detection. In: *Proceedings. 1998 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (Cat. No. 98CB36231)*. pp. 38–44. IEEE (1998)

34. Shi, X., Shan, S., Kan, M., Wu, S., Chen, X.: Real-time rotation-invariant face detection with progressive calibration networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2295–2303 (2018)
35. Wang, J., Ding, J., Guo, H., Cheng, W., Pan, T., Yang, W.: Mask obb: A semantic attention-based mask oriented bounding box representation for multi-category object detection in aerial images. *Remote Sensing* **11**(24), 2930 (2019)
36. Wang, Y., Zhang, Y., Zhang, Y., Zhao, L., Sun, X., Guo, Z.: Sard: Towards scale-aware rotated object detection in aerial imagery. *IEEE Access* **7**, 173855–173865 (2019)
37. Wei, H., Zhou, L., Zhang, Y., Li, H., Guo, R., Wang, H.: Oriented objects as pairs of middle lines. arXiv preprint arXiv:1912.10694 (2019)
38. Wold, S., Esbensen, K., Geladi, P.: Principal component analysis. *Chemometrics and intelligent laboratory systems* **2**(1-3), 37–52 (1987)
39. Xia, G.S., Bai, X., Ding, J., Zhu, Z., Belongie, S., Luo, J., Datcu, M., Pelillo, M., Zhang, L.: Dota: A large-scale dataset for object detection in aerial images. In: Proc. CVPR (2018)
40. Xu, Y., Fu, M., Wang, Q., Wang, Y., Chen, K., Xia, G.S., Bai, X.: Gliding vertex on the horizontal bounding box for multi-oriented object detection. arXiv preprint arXiv:1911.09358 (2019)
41. Yang, X., Liu, Q., Yan, J., Li, A., Zhang, Z., Yu, G.: R3det: Refined single-stage detector with feature refinement for rotating object. arXiv preprint arXiv:1908.05612 (2019)
42. Yang, X., Sun, H., Fu, K., Yang, J., Sun, X., Yan, M., Guo, Z.: Automatic ship detection in remote sensing images from google earth of complex scenes based on multiscale rotation dense feature pyramid networks. *Remote Sensing* **10**(1), 132 (2018)
43. Yang, X., Sun, H., Sun, X., Yan, M., Guo, Z., Fu, K.: Position detection and direction prediction for arbitrary-oriented ships via multitask rotation region convolutional neural network. *IEEE Access* **6**, 50839–50849 (2018)
44. Yang, X., Yang, J., Yan, J., Zhang, Y., Zhang, T., Guo, Z., Sun, X., Fu, K.: Scrdet: Towards more robust detection for small, cluttered and rotated objects. In: The IEEE International Conference on Computer Vision (ICCV) (October 2019)
45. Zhang, G., Lu, S., Zhang, W.: Cad-net: A context-aware detection network for objects in remote sensing imagery. *IEEE Transactions on Geoscience and Remote Sensing* **57**(12), 10015–10024 (2019)
46. Zhang, Z., Guo, W., Zhu, S., Yu, W.: Toward arbitrary-oriented ship detection with rotated region proposal and discrimination networks. *IEEE Geoscience and Remote Sensing Letters* **15**(11), 1745–1749 (2018)
47. Zhou, L., Wei, H., Li, H., Zhang, Y., Sun, X., Zhao, W.: Objects detection for remote sensing images based on polar coordinates. arXiv preprint arXiv:2001.02988 (2020)
48. Zhou, X., Zhuo, J., Krahenbuhl, P.: Bottom-up object detection by grouping extreme and center points. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 850–859 (2019)
49. Zhou, X., Yao, C., Wen, H., Wang, Y., Zhou, S., He, W., Liang, J.: East: An efficient and accurate scene text detector. In: Proceedings of the IEEE conference on computer vision and pattern recognition (2017)
50. Zhu, Y., Wu, X., Du, J.: Adaptive period embedding for representing oriented objects in aerial images. arXiv preprint arXiv:1906.09447 (2019)