

Learning a Rotation Invariant Detector with Rotatable Bounding Box

Lei Liu, Zongxu Pan, Bin Lei
Institute of Electronics, Chinese Academy of Sciences
{lliu1, zxpan, leibin}@mail.ie.ac.cn

Abstract

Detection of arbitrarily rotated objects is a challenging task due to the difficulties of locating the multi-angle objects and separating them effectively from the background. The existing methods are not robust to angle varies of the objects because of the use of traditional bounding box, which is a rotation variant structure for locating rotated objects. In this article, a new detection method is proposed which applies the newly defined rotatable bounding box (RBox). The proposed detector (DRBox) can effectively handle the situation where the orientation angles of the objects are arbitrary. The training of DRBox forces the detection networks to learn the correct orientation angle of the objects, so that the rotation invariant property can be achieved. DRBox is tested to detect vehicles, ships and airplanes on satellite images, compared with Faster R-CNN and SSD, which are chosen as the benchmark of the traditional bounding box based methods. The results shows that DRBox performs much better than traditional bounding box based methods do on the given tasks, and is more robust against rotation of input image and target objects. Besides, results show that DRBox correctly outputs the orientation angles of the objects, which is very useful for locating multi-angle objects efficiently. The code and models are available at <https://github.com/liulei01/DRBox>.

1. Introduction

Object detection is one of the most challenging tasks in computer vision and has attracted a lot of attentions all the time. Most existing detection methods use bounding box to locate objects in images. The traditional bounding box is a rotation variant data structure, which becomes a shortcoming when the detector has to deal with orientation variations of target objects. This article discusses how to design and train a rotation invariant detector by introducing the rotatable bounding box (RBox). Unlike traditional bounding box (BBox) which is the circumscribed rectangle of a rotated object, RBox is defined to involve the orientation information into its data structure. The proposed de-

tor (Detector using RBox, DRBox) is very suitable for detection tasks that the orientation angles of the objects are arbitrarily changed.

Rotation invariant property becomes important for detection when the viewpoint of the camera moves to the top of the object thus the orientations of the objects become arbitrarily. A typical application is the object detection task in aerial and satellite images. Hundreds of remote sensing satellites are launched into space each year and generate huge amounts of images. Object detection on these images is of great importance, whereas the existing detection methods suffer from the lack of rotation invariant property. In the next paragraph, we firstly give a brief overview of the object detection methods, and then discuss how rotation invariant is taken into consideration by the recent methods.

Previous object detection methods usually use economic features and inference schemes for efficiency, and prevalent such methods include deformable part model (DPM) [8], selective search (SS) [24] and EdgeBoxes [28]. As the development of deep learning technique, deep neural networks (DNNs) have been applied for solving the object detection problem and the DNN based approaches achieve state-of-the-art detected performance. Among DNN based methods, region-based detector is widely used. Region-based convolutional neural network method (R-CNN) [10] makes use of SS for generating region proposals, and convolutional neural networks (CNN) is then employed on each proposal for detection. R-CNN is slow, in part because every proposal has to be fed into the network individually, in part because it is a multi-stage pipeline approach. Consequent approaches take much effort to integrate the detection pipelines gradually. Spatial pyramid pooling networks (SPPnets) [13] speeds up R-CNN by sharing computation. The spatial pyramid pooling is introduced to remove the fixed size input constraint and the whole image needs only pass the net once. Several innovations are employed in Fast R-CNN [9] to improve both the effect and the efficiency of detection, including the use of RoI layer, multi-task loss, and the truncated SVD. Instead of using SS to generate region proposals, Faster R-CNN applies region proposal networks (RPNs) to generate the proposals [19], making all

detection steps be integrated in an unified network, which achieves the best detected result at that time. Anchor boxes with multiple scales and aspect ratios are of the essence in Faster R-CNN, which locate the position of candidate objects. In Fast/Faster R-CNN, the network can be partitioned into a fully convolutional subnetwork which is independent of proposals and shares the computation, and a per-proposal subnetwork without sharing the computation. It is noticed that the per-proposal subnetwork is inefficient, and a region-based fully convolutional network (R-FCN) is proposed in [6] to remedy that issue by using FCN to share almost all computation on the entire image. By adding a branch to predict the mask of objects upon Faster R-CNN, Mask R-CNN can simultaneously detect the object and generate the segmentation mask of the object [12]. There is another kind of DL-based object detection approach that does not rely on region proposals, such as you only look once detector (YOLO) [18] and single shot multibox detector (SSD) [16]. We refer to YOLO and SSD as box-based methods since they generate several boxes for detecting objects in the image according to certain rules. These boxes are called as prior boxes and objects are supposed to locate in or near certain prior box. SSD attempts to use pyramidal feature hierarchy computed from the convolutional net, however without reusing the higher resolution maps in the pyramidal feature hierarchy. The feature pyramid network (FPN) [15] better utilizes the pyramidal feature hierarchy through introducing lateral connections which merge high-level semantic feature maps with coarse resolution and low-level semantic feature maps with refined resolution.

Many recent machine learning based methods have been used for remote sensing object detection. Many machine learning methods extract the candidate object feature, such as histogram of oriented gradients (HOG) [23, 4], texture [3], bag-of-words (BoW) [22], regional covariance descriptor (RCD) [1], and interest point [25], followed by certain classifier, for example, sparse representation classifier (SRC) [4, 3] and support vector machine (SVM) [1]. As the development of deep learning technique, DNNs have been successfully employed to solve object detection problem in remote sensing images. A CNN based pioneering study upon vehicle detection in satellite images is presented in[2], followed by many approaches that focus on different types of geospatial object detection problems [14, 20, 21, 5, 26, 17]. As another effective feature extraction tool, deep Boltzmann machine is also used as feature extractor in geo-spatial object detection problems [11, 7].

To make the approach insensitive to objects in-plane rotation, some efforts are made either adjusting the orientation, or trying to extract rotation insensitive features. For example, to deal with the rotation variation of geo-spatial objects issue, a rotation invariant regularization term is in-

troduced which enforces the samples and their rotated versions share the similar features [5]. Unlike these methods which try to eliminate the effect of rotation on the feature level, we prefer to make the rotation information useful for feature extraction so that the detection results involve the angle information of the objects. Therefore, the detection results is rotatable, whereas the performance of the detector is rotation invariant.

Zhou *et al.* proposed a network named ORN [27] which shows similar character with our method. ORN is used to classify images while extracting the orientation angle of the whole image. However, ORN can not be applied straightforward as a detection network, which needs to detect orientation locally on each object. In our method, the angle estimation is associated to plenty of prior boxes, so that the rotation of an object can be realized by the corresponding prior boxes, while other prior boxes are still available for other objects. Additionally, our method effectively separates the object proposals with its background pixels, so the angle estimation can concentrate on the object without the interfere of background.

The next section explains why rotatable bounding box (RBox) is better than BBox for rotation invariant detection. In section 3 we discuss how RBox takes place of BBox to form the newly designed detection method, DRBox. DRBox is tested on ship, airplane and vehicle detection in remote sensing images, exhibiting definitely superiority compared with traditional BBox based methods.

2. Rotatable bounding box

Most object Detection methods use bounding box (BBox) to locate target objects in images. A BBox is a rectangle parameterized by four variables: the center point position (two variables), the width, and the height. BBox meets difficulties locating objects with different orientation angles. In this situation, BBox cannot provide the exact sizes of objects, and is very difficult to distinguish dense objects. Table 1 lists the three main disadvantages of BBox and their corresponding examples. As shown in this table, a ship target in remote sensing image is aligned to BBoxes with different sizes and aspect ratios because of rotation. As the result, the width and height of BBox have no relationship with the physical size of the rotated target. Another significant disadvantage is the discrimination between object and background. In the given example, about 60% region inside the BBox belongs to background pixels when the orientation angle of the ship is near 45 degrees. The situation becomes more difficult when target objects are distributed dense, in which case the objects are hard to be separated by BBoxes.

In this article, RBox is defined to overcome the above difficulties. RBox is a rectangle with a angle parameter to define its orientation. An RBox needs five parameters to de-

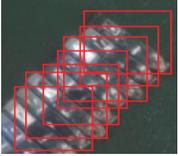
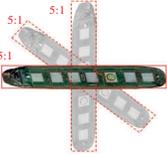
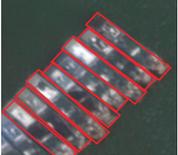
Disadvantages of the traditional bounding box		
1. The size and aspect ratios do not reflect the real shape of the target object.	2. Object and background pixels are not effectively separated.	3. Dense objects are difficult to be separated.
		
Advantages of the rotatable bounding box		
1. The width and height of RBox reflect the physical size of the object, which is helpful for customized designing of the prior boxes.	2. RBox contains less background pixels than BBox does, so classification between object and background is easier.	3. RBox can efficiently separate dense objects with no overlapped areas between nearby targets.
		

Table 1: Comparison of traditional bounding box and rotatable bounding box. Examples on ship locating are used to support the conclusions.

fine its location, size and orientation. Compared with BBox, RBox surrounds the outline of the target object more tightly, therefore overcomes all the disadvantages listed in the table. A detail comparison of RBox and BBox is demonstrated in Table 1. So we suggest that RBox is a better choice on detection of rotated objects.

Given two boxes, it is important for a detection algorithm to evaluate their distance, which is used to select positive samples during training, and suppress repeated predictions in detection. The common used criterion for BBox is the Intersection-over-Union (IoU), which can also be used by RBox. The IoU between two RBoxes A and B defines as following:

$$\text{IoU}(A, B) = \frac{\text{area}(A \cap B)}{\text{area}(A \cup B)}, \quad (1)$$

where \cap and \cup are Boolean operations between two RBoxes. The Boolean calculation between RBoxes is more complex than BBox because the intersection of two RBoxes can be any polygon with no more than eight sides.

Another criterion for RBox is angle-related IoU (ArIoU),

which is defined as following:

$$\text{ArIoU}(A, B) = \frac{\text{area}(\hat{A} \cap B)}{\text{area}(\hat{A} \cup B)} \cos(\theta_A - \theta_B), \quad (2)$$

or

$$\text{ArIoU}_{180}(A, B) = \frac{\text{area}(\hat{A} \cap B)}{\text{area}(\hat{A} \cup B)} |\cos(\theta_A - \theta_B)|, \quad (3)$$

where θ_A and θ_B are angles of RBox A and B , \hat{A} is an RBox which keeps the same parameters with RBox A except that the angle parameter is θ_B , not θ_A . ArIoU takes angle difference into consideration so that the ArIoU between RBox A and B decreases monotonically when their angle difference changes from 0 degree to 90 degrees. The two definitions differ in the behavior when $(\theta_A - \theta_B)$ is near 180 degrees. ArIoU₁₈₀ ignores the head and tail direction of the objects when distinguish them is too difficult.

IoU and ArIoU are used in different way. ArIoU is used for training so it can enforce the detector to learn the right angle, while IoU is used for non-maximum suppression (NMS) so the predictions with inaccurate angle can be effectively removed.

3. Rotation invariant detection

In this section, we apply RBox into detection, which means that the detector must learn not only the locations and sizes, but also the angles of the target objects. Once this purpose is achieved, the network can realize the existence of orientation difference between objects, rather than being confused by rotation. The performance of the detector becomes rotation invariant as the result.

3.1. Model

NETWORK STRUCTURE: DRBox uses a convolutional structure for detection, as shown in Figure 1. The input image goes through multi-layer convolution networks to generate detection results. The last convolution layer is for prediction and the other convolution layers are for feature extraction. The prediction layer includes K groups of channels where K is the number of prior RBoxes in each position. Prior RBoxes is a series of predefined RBoxes. For each prior RBox, the prediction layer output a confidence prediction vector indicating whether it is a target object or background, and a 5 dimensional vector which is the offset of the parameters between the predicted RBox and the corresponding predefined prior RBox. A decoding process is necessary to transform the offsets to the exact predicted RBoxes. At last, the predicted RBoxes are sorted with their confidence and passed through NMS to remove repeated predictions.

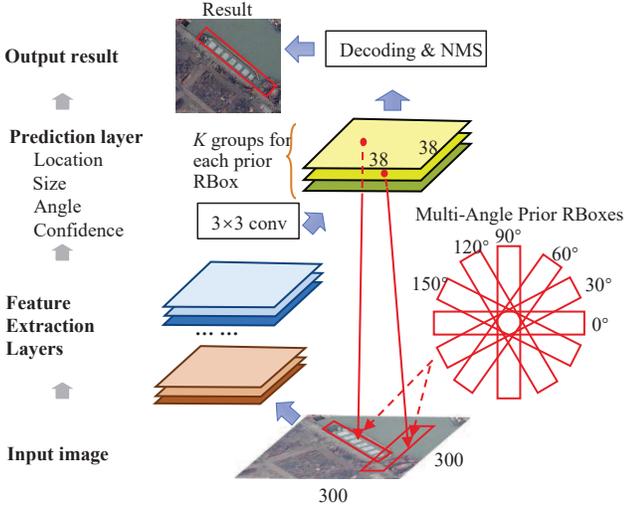


Figure 1: The networks structure of DRBox. The networks structure of DRBox is similar with other box based methods except for the use of multi-angle prior RBoxes. DRBox searches for objects using sliding and rotating prior RBoxes on input image and then output locations of objects besides with their orientation angles.

Multi-angle prior RBox plays an important role in DR-Box. The convolutional structure ensures that prior boxes can move over different locations to search for target objects. On each location, the prior RBoxes rotate at a series of angles to generate multi-angle predictions, which is the key difference between DRBox and other bounding box based methods. The aspect ratio used in detection is fixed according to the object type, which decreased the total number of prior boxes. By the multi-angle prior RBoxes strategy, the network is trained to treat the detection task as a series of sub-tasks. Each sub-task focuses on one narrow angle range, therefore decreases the difficulty caused by rotation of objects.

3.2. Training

The training of DRBox is extended from SSD training procedure [16] to involve angle estimation. During training, each ground truth RBox is assigned with several prior RBoxes according to their ArIoU. ArIoU is a non-commutative function, which means $\text{ArIoU}(A, B)$ is different from $\text{ArIoU}(B, A)$. A prior RBox P is assigned to a ground truth RBox G when $\text{ArIoU}(P, G) > 0.5$. After the assignment, the matched prior RBoxes are considered as positive samples and are responsible for generating the losses of location and angle regression. The use of ArIoU helps the training process to select prior RBox with proper angle as positive samples, so the angle information of the objects can be roughly learned during training. After the matching step, most of the prior RBoxes are negative. We

apply hard negative mining to decrease the number of the negative samples.

The objective loss function of DRBox is extended from SSD objective loss function by adding the angle related term. The overall objective loss function is as following:

$$L(x, c, l, g) = \frac{1}{N} (L_{\text{conf}}(c) + L_{\text{rbox}}(x, l, g)) \quad (4)$$

where N is the number of matched prior RBoxes. The confidence loss $L_{\text{conf}}(c)$ is a two class softmax loss over all selected positive and negative samples, where c is the two-dimension confidence vector. The RBox regression loss $L_{\text{rbox}}(x, l, g)$ is similar to SSD and Faster R-CNN, where we calculate the smooth L_1 loss between the predicted RBox l and the ground truth RBox g :

$$\begin{aligned} L_{\text{rbox}}(x, l, g) &= \sum_{i \in \text{Pos}} \sum_j \sum_{m \in \{cx, cy, w, h, a\}} x_{ij} \text{smooth}_{L1}(\hat{l}_i^m - \hat{g}_j^m) \end{aligned} \quad (5)$$

where $x_{ij} \in \{1, 0\}$ is an indicator for matching the i -th prior RBox to the j -th ground truth RBox. \hat{l} and \hat{g} are defined as following, which are the offsets of the parameters in l and g with their corresponding prior RBox p , respectively:

$$\hat{t}^{cx} = (t^{cx} - p^{cx})/p^w, \quad \hat{t}^{cy} = (t^{cy} - p^{cy})/p^h; \quad (6a)$$

$$\hat{t}^w = \log(t^w/p^w), \quad \hat{t}^h = \log(t^h/p^h); \quad (6b)$$

$$\hat{t}^a = \tan(t^a - p^a). \quad (6c)$$

Equations 6a , 6b and 6c are the location regression terms, the size regression terms and **the angle regression term**, respectively. The angle regression term applies tangent function to adapt to the periodicity of the angle parameter. The minimization of the angle regression term ensures that the correct angle is learned during training.

3.3. Complement details

PYRAMID INPUT: DRBox applies pyramid input strategy that the original image is rescaled into different resolutions, and separated into overlapped 300×300 sub-images. The DRBox network is applied to each sub-image and the network only detects targets with proper size. Non-maximum suppression is applied on the detection results of the whole image, which suppresses repeated predictions not only within a sub-image, but also crossing overlap areas of different sub-images. The pyramid input strategy helps the detection network to share features between large and small objects. Besides, the satellite image used in this article is often very large, so the division process and non-maximum suppression across sub-images helps to detect objects in very large images.

CONVOLUTION ARCHITECTURE: DRBox uses truncated VGG-net for detection. All the full-connective layers, convolution layers and pooling layers after layer conv4_3 are removed. Then, a 3×3 convolution layer is added after layer conv4_3. The receptive field of DRBox is 108 pixels 108 pixels, so any targets larger than this scope cannot be detected. Besides, the feature map of layer conv4_3 is 38×38 , so the targets closer than 8 pixels may be missed.

PRIOR RBOX SETTINGS: In this article, three DRBox networks are trained separately for vehicle detection, ship detection and airplane detection, respectively. The scale and input resolution settings jointly ensure that the areas of the prior RBoxes cover the sizes of the objects sufficiently, thus the objects of different sizes can be effectively captured. In ship detection, it is hard to distinguish the head and tail of the target. In this case, the angles of the ground truth RBoxes and multi-angle prior RBoxes varies from 0 degree to 180 degrees. In detail, ship objects are detected with prior RBoxes of 20×8 , 40×14 , 60×17 , 80×20 , 100×25 pixels in size and 0:30:150 degrees in angle; vehicle objects are detected with prior RBoxes of 25×9 pixels in size and 0:30:330 degrees in angle; airplane objects are detected with prior RBoxes of 50×50 , 70×70 pixels in size and 0:30:330 degrees in angle. The total number of prior boxes per image is 43320, 17328 and 34656 for ship, vehicle and airplane detection, respectively.

DRBox reaches 70-80 fps on NVIDIA GTX 1080Ti and Intel Core i7. The input pyramid strategy produce no more than $4/3$ times time cost. Considering $1/3$ overlapped between sub-images, DRBox reaches processing a speed of 1600×1600 pixels² per second. The speed of SSD and Faster R-CNN are 70 fps and 20 fps on our dataset using the same convolution network architecture.

4. Experiments and results

4.1. Dataset

We apply our method on object detection of satellite images. We have not found any open source dataset on this problem, so we build one using the GoogleEarth images. The dataset includes three categories of objects: vehicles, ships and airplanes. The vehicles are collected from urban area in Beijing, China. The ships are collected near the wharfs and ports besides the Changjiang River, the Zhujiang River and the East China Sea. The airplanes are collected from the images of 15 airports in China and America. The dataset recently includes about 12000 vehicles, 3000 ships and 2000 airplanes and is still under expansion. About 2000 vehicles, 1000 ships and 500 airplanes are taken out as testing dataset and others are used for training.

Each object in the images are marked with a RBox, which indicates not only the location and size, but also the

angle of the object. A Matlab tool is developed to label the data with RBoxes.

4.2. Benchmark

In this section, we compare the performance of DRBox with the detectors that use BBox. SSD and Faster R-CNN are used as benchmarks of BBox based methods. All the detectors use the same convolution architecture and the same training data argumentation strategy. The prior boxes used in SSD and the anchor boxes used in Faster R-CNN are optimized for the datasets. All other hyper parameters are optimized for the dataset, too.

4.3. Detection results

Figure 2 show ships, vehicles and airplanes detected by using RBox (DRBox) and BBox (SSD). The predicted bounding boxes that matches the ground truth bounding boxes with $\text{IoU} > 0.5$ are plotted in green color, while the false positive predictions and false negative predictions are plotted in yellow and red color, respectively. Our method is better in the given scenes. DRBox successfully detects most of the ships on both the port region and open water region, while SSD almost fails on detection of nearby ships. The detection results for vehicles and airplanes also show that SSD generates more false alarms and false dismissals than does DRBox.

More results of DRBox are shown in 3. Ship detection in port region is more challengeable than in open water region, whereas DRBox works well on both situations. The vehicles are very difficult to detect due to small size and complex backgrounds. In our dataset, each car is around 20 pixels in length and 9 pixels in width. Fortunately, we find that DRBox successfully finds vehicles that hid in the shadows of tall buildings, or parked very close to each other. We also find that DRBox can not only output the locates of the cars, but also predict the head direction of each car, which is even a challenge task for human beings. The estimated direction of cars on road matches the prior knowledge that traffic always keeps to the right of the road. Airplanes with different sizes are successfully detected, including an airplane that is under repairing.

Figure 4 shows precision-recall (P-R) curves of DRBox, SSD and Faster R-CNN. The recall ratio evaluates the ability of finding more targets in the image, while the precision evaluates the quality of predicting only the right object rather than containing many false alarms. The P-R curve of SSD and Faster R-CNN are always below the P-R curve of DRBox. DRBox has the best performance in this test. We further show BEP (Break-Even Point), AP (Average Precision) and mAP (mean Average Precision) of each method in Table 2. BEP is the point on P-R curve where precision equals recall, AP is the area below P-R curve, and mAP is the mean value of the APs on all object detection tasks.

Method	Dataset	BEP(%)	AP(%)	mAP(%)
Faster R-CNN	Ship	79.20	82.29	85.63
	Vehicle	71.60	75.55	
	Airplane	98.07	99.06	
SSD	Ship	82.72	82.89	89.68
	Vehicle	83.13	87.59	
	Airplane	97.74	98.56	
DRBox	Ship	94.62	94.06	94.13
	Vehicle	86.14	89.07	
	Airplane	98.62	99.28	

Table 2: BEP, AP and mAP of Faster R-CNN, SSD and DRBox. DRBox outperforms the other two methods on all indexes listed in this table.

DRBox is always the best on all the indexes compared with SSD and Faster R-CNN.

4.4. Comparison of the robustness against rotation

The robustness against rotation of a method involves two aspects: robustness against rotation of input images and robustness against rotation of objects. Robustness against rotation of input images means that a detection method should output the same results when the input image rotates arbitrarily. We define STD_AP to quantify this ability. STD_AP is the standard deviation of AP when the angle of the input image changes. STD_AP is estimated by rotating the test images each 10 degrees, then calculate the AP values of the detection results, respectively, and the standard deviation of the APs. Robustness against rotation of objects means that the same object should always be successfully detected when its orientation changes. We define STD_AS to quantify this ability. STD_AS is the standard deviation of the average score for objects in different orientation angles, which is estimated by dividing the objects in testing dataset into different groups according to their angles, then calculate average score (softmax threshold) of each group, respectively, and the standard deviation over all groups. The two robustness evaluation methods are interrelated, except that STD_AP has a bias towards the robustness on rotation of backgrounds while STD_AS has a bias towards the robustness on rotation of objects. The smaller STD_AP and STD_AS value indicate that the detection method is more robust against rotation.

STD_AP and STD_AS values of DRBox, SSD and Faster R-CNN are shown in Table 3. All the three methods are robust against rotation on airplane detection. However, Faster R-CNN is relatively not robust against rotation on ship and vehicle detection, SSD is not robust against rotation on ship detection. DRBox remains good scores in all tasks.

Further comparison between DRBox and SSD are demonstrated in Figure 5 and Figure 6. Figure 5 shows P-R recalls of the same input image but rotated to different

Method	Dataset	STD_AP (%)	STD_AS (%)
Faster R-CNN	Ship	5.51	13.76
	Vehicle	7.21	7.82
	Airplane	0.28	0.17
SSD	Ship	5.97	13.14
	Vehicle	0.90	3.20
	Airplane	0.80	0.51
DRBox	Ship	0.88	1.56
	Vehicle	0.72	2.06
	Airplane	0.51	0.41

Table 3: Quantitative evaluation of robustness against rotation for Faster R-CNN, SSD and DRBox. STD_AP evaluates robustness against rotation of input images. STD_R evaluates robustness against rotation of objects. DRBox outperforms the other two methods except that Faster R-CNN is slightly more robust to rotation of airplane targets.

angles, where the results of DRBox and SSD are plotted in shallow red and shallow blue colors, respectively. The curves generated by DRBox are more concentrated, which indicates that DRBox is more robust against rotation of input images compared with SSD.

We calculate the recall ratio of targets in different angle scopes, respectively. Figure 6 show the results. Each angle scope corresponds to one curve which shows the relationship between the softmax threshold and recall ratio. The performance of DRBox approximately remains the same for each angle scopes, whereas the performance of SSD shows strong instability when the angle of objects changes.

5. Conclusions

Robustness to rotation is very important on detection tasks of arbitrarily orientated objects. Existing detection algorithms uses bounding box to locate objects, which is a rotation variant structure. In this article, we replace the traditional bounding box with RBox and reconstruct deep CNN based detection frameworks with this new structure. The proposed detector, which is called DRBox, is rotation invariant due to its ability of estimating the orientation angles of objects. DRBox outperforms Faster R-CNN and SSD on object detection of satellite images.

DRBox is designed as a box-based method, whereas it is also possible to apply RBox into proposal based detection frameworks, e.g. R-FCN or Faster R-CNN. Training with RBox enforces the network to learn multi-scale local orientation information of the input image. We are looking forward of this interesting property to be used in other orientation sensitive tasks.



Figure 2: Detection results using RBox and traditional bounding box. The first row are results of DRBox; the second row are results of SSD. Columns 1-2, columns 3-4 and column 5 are results of ship detection, vehicle detection and airplane detection, respectively. DRBox performs better in the given examples.



Figure 3: Detection results of DRBox. Ship detection results are shown in the first row, where DRBox works well in both open water region and port region. Vehicle detection results are shown in the second row, where DRBox successfully finds cars hidden in complex backgrounds and correctly indicates their head directions. Airplane detection results are shown in the third row, including an airplane that is under repairing.

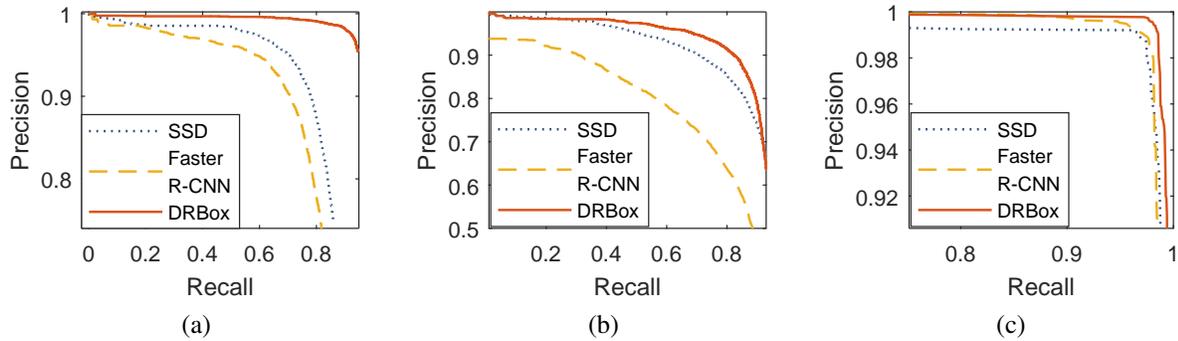


Figure 4: Precision-recall curves of (a) ship detection results, (b) vehicle detection results and (c) airplane detection results. The performance of DRBox is the best in each detection task.

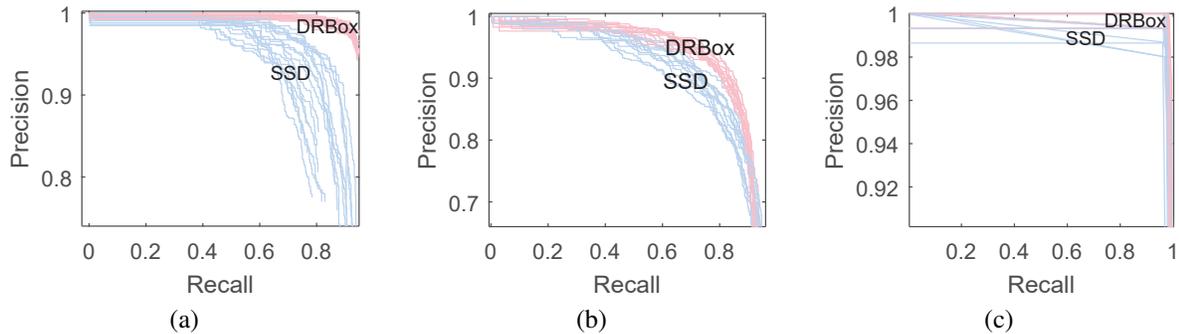


Figure 5: Precision-recall curves when input images are rotated on multiple angles in (a) ship detection task, (b) vehicle detection task and (c) airplane detection task. The curves generated by SSD are in blue color, and the curves generated by DRBox are in red color. The curves generated by DRBox are more concentrated, which indicates that DRBox is more robust to rotation of input images.

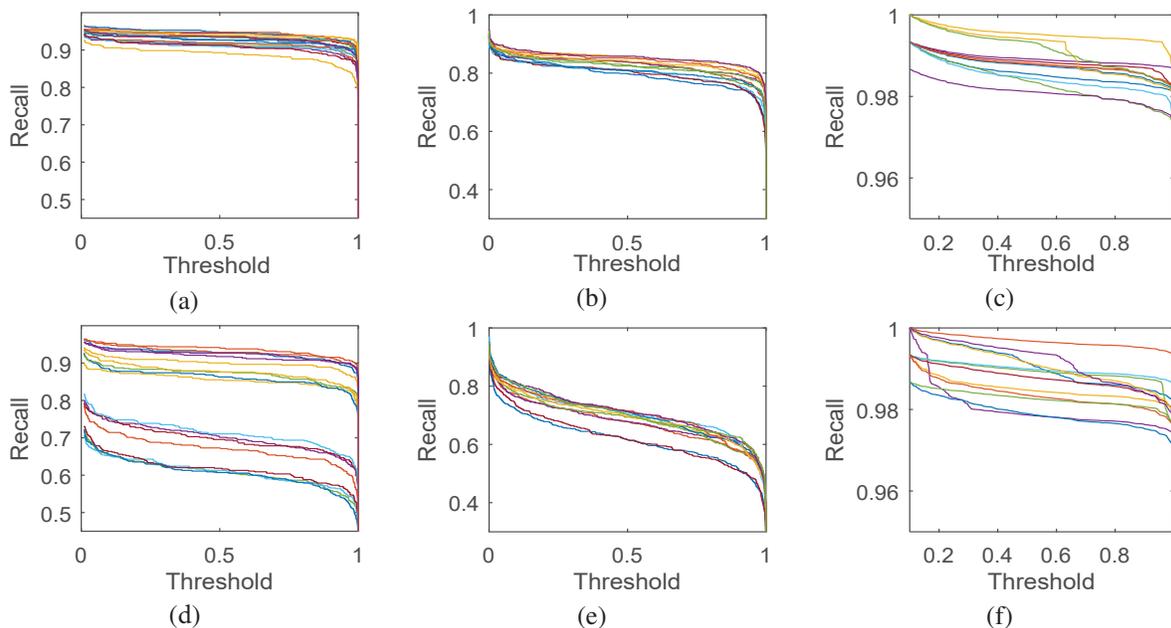


Figure 6: The curves of recall ratios on different softmax threshold values. In each sub-figure, different curves are for different object orientation angles. (a) (b) (c) are results of DRBox in ship detection task, vehicle detection task and airplane detection task, respectively; (d) (e) (f) are results of SSD in ship detection task, vehicle detection task and airplane detection task, respectively. The curves generated by DRBox are more concentrated, which indicates that DRBox is more robust to rotation of objects.

References

- [1] X. Chen, R.-X. Gong, L.-L. Xie, S. Xiang, C.-L. Liu, and C.-H. Pan. Building regional covariance descriptors for vehicle detection. *IEEE Geoscience and Remote Sensing Letters*, 14(4):524–528, 2017.
- [2] X. Chen, S. Xiang, C.-L. Liu, and C.-H. Pan. Vehicle detection in satellite images by hybrid deep convolutional neural networks. *IEEE Geoscience and remote sensing letters*, 11(10):1797–1801, 2014.
- [3] Z. Chen, C. Wang, H. Luo, H. Wang, Y. Chen, C. Wen, Y. Yu, L. Cao, and J. Li. Vehicle detection in high-resolution aerial images based on fast sparse representation classification and multiorder feature. *IEEE Transactions on Intelligent Transportation Systems*, 17(8):2296–2309, 2016.
- [4] Z. Chen, C. Wang, C. Wen, X. Teng, Y. Chen, H. Guan, H. Luo, L. Cao, and J. Li. Vehicle detection in high-resolution aerial images via sparse representation and superpixels. *IEEE Transactions on Geoscience and Remote Sensing*, 54(1):103–116, 2016.
- [5] G. Cheng, P. Zhou, and J. Han. Learning rotation-invariant convolutional neural networks for object detection in vhr optical remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing*, 54(12):7405–7415, 2016.
- [6] J. Dai, Y. Li, K. He, and J. Sun. R-fcn: Object detection via region-based fully convolutional networks. In *Advances in neural information processing systems*, pages 379–387, 2016.
- [7] W. Diao, X. Sun, X. Zheng, F. Dou, H. Wang, and K. Fu. Efficient saliency-based object detection in remote sensing images using deep belief networks. *IEEE Geoscience and Remote Sensing Letters*, 13(2):137–141, 2016.
- [8] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part-based models. *IEEE transactions on pattern analysis and machine intelligence*, 32(9):1627–1645, 2010.
- [9] R. Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015.
- [10] R. Girshick, J. Donahue, T. Darrell, and J. 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.
- [11] J. Han, D. Zhang, G. Cheng, L. Guo, and J. Ren. Object detection in optical remote sensing images based on weakly supervised learning and high-level feature learning. *IEEE Transactions on Geoscience and Remote Sensing*, 53(6):3325–3337, 2015.
- [12] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. *arXiv preprint arXiv:1703.06870*, 2017.
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. In *European Conference on Computer Vision*, pages 346–361. Springer, 2014.
- [14] Q. Jiang, L. Cao, M. Cheng, C. Wang, and J. Li. Deep neural networks-based vehicle detection in satellite images. In *Bioelectronics and Bioinformatics (ISBB), 2015 International Symposium on*, pages 184–187. IEEE, 2015.
- [15] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie. Feature pyramid networks for object detection. *arXiv preprint arXiv:1612.03144*, 2016.
- [16] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016.
- [17] Y. Long, Y. Gong, Z. Xiao, and Q. Liu. Accurate object localization in remote sensing images based on convolutional neural networks. *IEEE Transactions on Geoscience and Remote Sensing*, 55(5):2486–2498, 2017.
- [18] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 779–788, 2016.
- [19] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.
- [20] I. Ševo and A. Avramović. Convolutional neural network based automatic object detection on aerial images. *IEEE Geoscience and Remote Sensing Letters*, 13(5):740–744, 2016.
- [21] L. W. Sommer, T. Schuchert, and J. Beyerer. Fast deep vehicle detection in aerial images. In *Applications of Computer Vision (WACV), 2017 IEEE Winter Conference on*, pages 311–319. IEEE, 2017.
- [22] H. Sun, X. Sun, H. Wang, Y. Li, and X. Li. Automatic target detection in high-resolution remote sensing images using spatial sparse coding bag-of-words model. *IEEE Geoscience and Remote Sensing Letters*, 9(1):109–113, 2012.
- [23] S. Tuermer, F. Kurz, P. Reinartz, and U. Stilla. Airborne vehicle detection in dense urban areas using hog features and disparity maps. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 6(6):2327–2337, 2013.
- [24] J. R. Uijlings, K. E. Van De Sande, T. Gevers, and A. W. Smeulders. Selective search for object recognition. *International journal of computer vision*, 104(2):154–171, 2013.
- [25] L. Wan, L. Zheng, H. Huo, and T. Fang. Affine invariant description and large-margin dimensionality reduction for target detection in optical remote sensing images. *IEEE Geoscience and Remote Sensing Letters*, 2017.
- [26] F. Zhang, B. Du, L. Zhang, and M. Xu. Weakly supervised learning based on coupled convolutional neural networks for aircraft detection. *IEEE Transactions on Geoscience and Remote Sensing*, 54(9):5553–5563, 2016.
- [27] Y. Zhou, Q. Ye, Q. Qiu, and J. Jiao. Oriented response networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, July 2017.
- [28] C. L. Zitnick and P. Dollár. Edge boxes: Locating object proposals from edges. In *European Conference on Computer Vision*, pages 391–405. Springer, 2014.