
Cyclist detection in aerial images

Anonymous Author(s)

Affiliation

Address

email

Abstract

In this contribution, I propose an algorithm for detection of cyclists in aerial images. This task is somewhat different from the standard object recognition problems, since the bird's eye perspective prohibits direct transfer learning from datasets as eg. COCO or ImageNet. Nevertheless, a transfer learning solution is proposed but the base model is chosen with special care. The training and evaluation is performed on a dataset containing annotated pictures from Giro d'Italia 2015.

1 Object detection in aerial pictures

Object detection is one of the most important branches of computer vision. Since in any realistic environment, the machine-learning agent receives visual input containing multiple objects, it is vital to be able to detect and correctly classify environmental objects to give meaningful responses to external stimuli.

With the rise of unmanned aerial vehicles (UAVs or 'drones') the amount of visual data captured from a bird's eye perspective is fastly growing. It is the anticipation that one day, autonomous UAVs will considerably contribute to many businesses eg. delivery or agriculture.¹

These trends call for increased attention and effort to computer vision for aerial applications. This domain is different from conventional CV applications, because the data is unlike widely adopted standard datasets (COCO, ImageNet) since everything is captured from the height. Henceforth, re-training conventional object detection architectures is one promising way to go with preserving the basic principles of modern object recognition but tailoring weights to our special needs.

2 Literature survey, existing solutions

In this section, I first give an short overview of the fundamentals of modern object recognition. In the following, I introduce the particular results of earlier investigations on which the project is based. For a very concise overview, I direct the reader to a series of blog posts that fantastically capture the main points and put the field in context.⁽¹⁾ Thank you Lilian.

2.1 Object detection algorithms

The coming parts are rather standard introduction and laying the foundations for later discussion but by no means a comprehensive review of the enormous field of object detection.

¹And sadly, military applications and surveillance.

31 2.1.1 Preliminaries and basic notions

32 The task of object detection algorithms is twofold: (i) detect objects on a given input image, (ii)
33 classify the objects it has identified.

34 The solution to this task can be broken down to two parts: first, finding proposals for possible objects
35 and then using a conventional convolutional architecture to classify proposals.
36

37 It is not trivial how to measure the soundness of a proposal for a Region of Interest (ROI) or actually
38 how to represent objects inside images. In this work, we stick with the *bounding box* representation,
39 which is assigning a rectangular box to every object that is possibly the smallest box fully containing
40 the object. Our language is obviously loose here, but at the end of the day, boxes are determined by
41 humans annotating the images so no precise definition is possible. More advanced approaches use
42 segmentation masks that respect the shape of objects, but here we cannot go into details.
43

44 For bounding boxes, a very common metric to assess the soundness of the ROI is the *Intersection*
45 *Over Union* (IOU) which has a pretty straightforward meaning for a ROI proposal and a ground
46 truth box: measure the intersection's area and divide by the union's area. How ROIs are assigned to
47 ground truth (GT) boxes can vary from implementation to implementation. A very common choice is
48 to assign ROIs to that GT box for which the maximal IOU is attained.
49

50 For a given proposed region, classification is done with a CNN architecture. Unfortunately, due to
51 scarcity of room, I cannot elaborate more on the details of how ROIs are divided into foreground and
52 background and similar, more advanced topics. The interested reader can look it up in any competent
53 GitHub repository.

54 2.2 Architectures

55 In this section, I only introduce the *Faster-RCNN* architecture in detail since this is the one used in
56 the project.
57

58 2.2.1 History

59 Surprisingly from the two fundamental tasks of object recognition systems, the problem of ROI
60 proposals has been addressed earlier in a competent way. The famous *Felzenszwalb algorithm* uses a
61 graph representation of the images to basically perform clustering² on the pixels.⁽⁹⁾

62 The other task has been first satisfyingly tackled in the 2012 landmark AlexNet network. Blending
63 together the two parts followed fast.^{(2),(3)}

64 2.2.2 Faster-RCNN

65 This method was novel in the sense that surpassing the traditional segmentational methods, it also
66 used a neural network architecture to generate ROI proposals. The basic sturcture of a Faster-RCNN
67 network is depicted in figure 1.
68

69 The three main parts:

- 70 • **backbone**: this part is in fact a convolutional neural network with top classification layers
71 'peeled off'. Its purpose is feature extraction for later stages.
- 72 • **RPN**: aka Region Proposal Network. This is the part that produces ROI proposals from
73 the output of the backbone. Again, we cannot dive into the details but you shall look up
74 *anchor boxes* that are basically some fixed 'basic box proposals' from which the actual
75 proposition's coordinates are determined via regression.

²Our language is loose again: this may rather be the *selective search* component of the segmentation solution

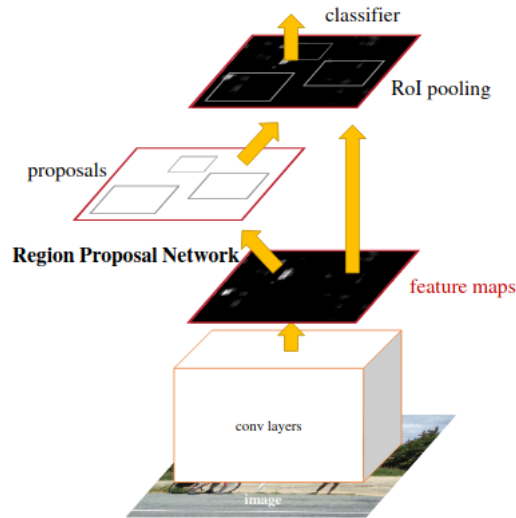


Figure 1: The three main parts of the Faster-RCNN architecture ⁽⁴⁾

- **ROI head:** in this part, we basically do the classification of the proposed regions. I repeatedly stress that the above description is only a vague one, you shall look up literature.

In summary, the output of a forward pass on the network is a list of proposed boxes with the predicted labels and confidency scores for the labels. I also used a Faster-RCNN net in the project. In particular, the used model goes by the name of `e2e_faster_rcnn_X_101_32x8d_FPN_1x`, you may look it up at the respective FAIR model zoo. (For more details, it is the best to refer to the config file from which the model is actually built.)

2.2.3 More advanced networks

In spite of being a great improvement over its predecessors, Faster-RCNN is not the most advanced arch actually.

In terms of speed, there are much more competent solutions.^{(6),(7),(8)}

Another qualitative improvement is implemented in the *Mask-RCNN* familiy where not only boxes but the very segmentation of the picture is also computed.⁽⁵⁾

2.3 Leveraged earlier solutions

The final solution has been built upon an earlier solution for the **VisDrone** dataset, introduced in the next section. As we had access only to limited GPU time and as being newcomers to the field, we realized that this was the way to go. The solution is due to *Oytun Ulutan* at UC Santa Barbara, and can be found under:

https://github.com/oulutan/Drone_FasterRCNN

Fortunately, the weights of the model are also published so the transfer scenario could be implemented. We had to amend the number of classes to our specific application, so practically trimmed the the last layers of the network and replaced them with the ones in agreement with our class numbers.

Oytun's work is actually a fork of FAIR's `maskrcnn-benchmark` repo,⁽¹²⁾ under

<https://github.com/facebookresearch/maskrcnn-benchmark>

3 Data sources

As we have decided to do an UAV project, we started looking up accessible databases. Two have been found:

- **VisDrone DET-2019 dataset:**^{(10),(11)} This dataset consists of images taken by different UAVs under various conditions, eg. in urban and rural areas, under varying lightning and weather conditions, from different heights, with few or many objects in one picture. The dataset was collected by the AISKEYE team at Lab of Machine Learning and Data Mining, Tianjin University, China. Although the repo is unclear about licensing, we believe that the data is free to be used for research purposes.
- **MultiDrone** The public MultiDrone Dataset⁽¹³⁾ has been assembled using both pre-existing audiovisual material and newly filmed UAV shots. The dataset is curated by the Multidrone Consortium. A large subset of these data has been annotated for facilitating scientific research, in tasks such as visual detection and tracking of bicycles, football players, human crowds, etc. A dataset for bicycle detection/tracking was assembled, consisting of 3 Youtube videos³ (resolution: 1920 x 1080) at 25 frames per second. Annotations are not exhaustive, i.e., there may be unannotated objects in the given video frames. The license agreement allows to use only for scientific research, testing and development purposes and we aren't authorized to put any part of the dataset on the publicly accessible Internet.

Since the input model for transfer learning was trained on the VisDrone data, I only used the Multidrone images in the transfer learning. Due to licensing issues, I had to invent a way not to make annotations publicly available but to preserve the reproducibility of the results. The final solution was setting up a secret GitHub Gist to that I will provide access upon reasonable request.

4 Data preparation

The frames were fetched from YouTube. What I realised were that

- The data has to be manually curated. In fact, the videos contained (annotated) frames that were captured from the ground so showed the competitors from a totally different angle from the usual aerial view. This was remedied by watching all the relevant parts and marking the segments in which real bird's eye perspective was presented. (Part of the reasons I worked with only three videos)
- The annotations have been prepared on 1920x1080 resolution but the images have been fetched in 1280x720. Unfortunately, one of my former teammates – who had written the script – failed to notice/mention this which would have saved me a lot of time and GPU time...

5 Utility engineering

As there is no rose without thorn, the VisDrone project came almost with no documentation whatsoever. After correcting and extending the Dockerfile, a final version has been uploaded to the GitHub repository. With this, any interested person should be able to set up the environment on his/her device. (After PM-ing me for the code of the annotations. For detailed instructions, refer to the repository.)

Handling the dataset with dataloaders was only possible after writing a new module for handling the Giro images, since this is by no means a standard dataset. Interestingly, the original code did not contain any type of on the fly validation for training nor callbacks.⁴ Even more interestingly, I am not sure if Pytorch contains *native* callbacks at all. (Ignore aside, since that was opted out for this project –

³In fact, there were seven, but due to scarcity of time (the videos had to be hand-curated, see later, only three were used

⁴Or they have been hidden deliberately :)

147 I am not sure if it would have required rewriting it all.)

148

149 This means that I had to write these utilities – with reloading custom weights for the network and
150 customizing trainable weights – almost from scratch. Nevertheless, it has been done since without
151 them, it would not have been possible to do any meaningful training.

152

153 Finally, as for 'visible' evaluation, I also put together (mostly from other parts of the library) a
154 pipeline that does the inference on single images and then draws the predicted boxes on top of them.

155 6 Training and evaluation

156 This is actually the section where all the mentioned sufferings paid off. After having a *stable*
157 train-validation pipeline, the training was quite straightforward⁵

- 158 1. Loading weights from the VisDrone project
- 159 2. Peeling off the classification layers and replacing them with the appropriate sized layers for
160 our problem (number of classes)
- 161 3. Updating state dict if applicable
- 162 4. Freezing custom layers. First, I only unfroze the newly added (so 'factory initialized')
163 layers. In the second run, I also unfroze the whole RPN and ROI head, to attain even better
164 convergence
- 165 5. Finally start training

166 6.1 Technical details

167 For every video, frames have been randomly split into train, test and validation sets with approximate
168 ratios 0.85 : 0.10 : 0.05. This is not exact since I had to remove frames that were inappropriate for
169 training: either bad perspective or containing unidentified objects (in the annotations, we had '-'
170 instead of some class).

171 So finally, the numbers:

video	giro1	giro4	giro8
train	2389	431	60
test	258	7	3

Table 1: Number of images in respective datasets

172 I haven't included the validation numbers. This is merely because validation was not really there to
173 assess performance precisely, but to see if the IOU values were acceptable or not (when I started
174 experimenting they were absolutely not). To save time, I only validated on few images.

175 This is kind of unfortunate practice, but I had to speed up everything because I was in acute need
176 of time after all my teammates abandoned me and skeletons unexpectedly started popping out of the
177 closet (weight loading, having received unmatching pics to annotations, need to do hand-curating,
178 etc.) Nevertheless, the convergence on the trainset was saturated by the end but the accuracy of
179 predictions on the validation set was in a raising trend even when train loss stagnated.

180 6.2 Evaluation

181 The evaluation is based on the transfer-trained model's performance on the testsets.

182

⁵Or it seems so. If you want to train your net, you must refer to the repo for the caveats. I somewhat wasted most of my time because the Visdrone weights have not been loaded properly. I only found out recently, when I haven't been able to overfit a minibatch. So for the record: it is done by the checkpoint, surprisingly.

183 I used two widely used metrics: accuracy of predictions on the ground truth boxes and the
 184 mean recall for IOU values which I will define in the following.

185 First, one defines threshold values for assessing IOU values. In our case, they were: [0.5000, 0.5500,
 186 0.6000, 0.6500, 0.7000, 0.7500, 0.8000, 0.8500, 0.9000, 0.9500]. After this, one calculates the
 187 proportion of IOU values that lie above the threshold: this gives the recall for that threshold. The
 188 mean of those values is the mean threshold. (If you are dissatisfied with this much aggregation, I
 189 totally understand, please refer to the logfiles for more details.)

190 The results are presented in table 2. What we can see that the performance is satisfying for sets 1 and
 191 8 but really bad for set 4. I hypothesize that it is due to the fact that different videos showed different
 192 conditions and the data was not of the same distribution.

193 By this I mean that the notorious data featured cyclists of label 'front view' whereas cyclists are
 194 typically shown in competitions in 'back view'. Unfortunately, I haven't had time to elaborate on this
 195 issue but the probable imbalance in class distributions may be remedied.

testset	giro1	giro4	giro8
number of GT boxes	1516	18	25
accuracy	0.81	0.0	0.80
mean recall	0.46	0.03	0.42

Table 2: Evaluation results

196 6.3 Demo

197 Here, I present some impressions what I acquired while looking at predictions drawn on figures.
 198

199 First and foremost, looking at figure 2 is really sobering: it is pretty much possible that the algorithm
 200 learned to recognize the TV graphics! I also found some inconsistent behaviour (seeming mismatch
 201 between annotations and images). This must be cleared in any further project.

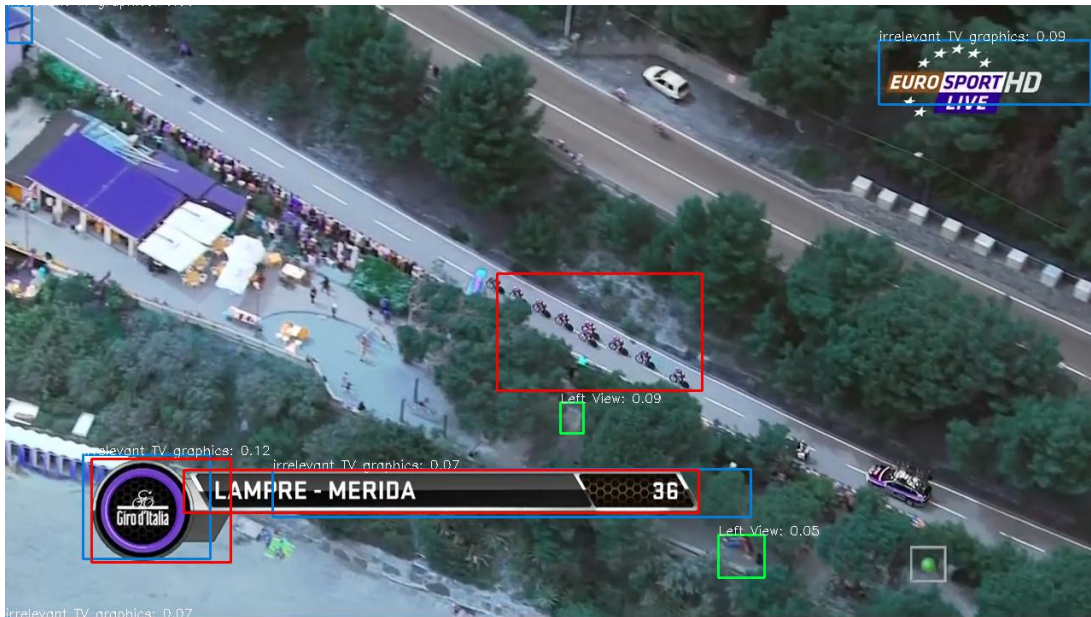


Figure 2: Predicted boxes vs. ground truth (red)

202 It is really hard to write down this, but the above suggests that our results are barely meaningful. In
 203 any future project, class balance (TV graphics!) shall be investigated in detail. Also, this means that
 204 the work I introduced is only making work the machinery but not yet getting acceptable results.

7 Conclusion, plans

To sum up, I must say that to achieve acceptable results, more training and deeper analysis of the data is indispensable.

It is however a fact calling for respect that I had to do the project almost totally alone while my teammates left before doing any meaningful work.

- First and foremost, it shall be investigated if the annotations really match the frames. (Some results on set 'giro4' hint otherwise, but I noticed this so late that I couldn't go into much detail)
- Then, accuracy and IOU values shall be broken down to classes not to have dominant TV graphics or suppressed classes eg. front view
- It must be re-considered what kinds of metrics to use to assess performance we really care about

References

- [1] Lilian Weng,
<https://lilianweng.github.io/lil-log/2017/10/29/object-recognition-for-dummies-part-1.html>
- [2] Ross Girshick et al.
Rich feature hierarchies for accurate object detection and semantic segmentation
<https://arxiv.org/abs/1311.2524>
- [3] Ross Girshick,
Fast R-CNN
<https://arxiv.org/pdf/1504.08083.pdf>
- [4] Shaoqing Ren et al,
Faster R-CNN: Towards Real-Time ObjectDetection with Region Proposal Networks
<https://arxiv.org/pdf/1506.01497.pdf>
- [5] Kaiming He et al,
Mask R-CNN
<https://arxiv.org/pdf/1703.06870.pdf>
- [6] Joseph Redmon et al
You Only Look Once:Unified, Real-Time Object Detection
https://www.cv-foundation.org/openaccess/content_cvpr_2016/papers/Redmon_You_Only_Look_CVPR_2016_paper.pdf
- [7] Wei Liu et al,
SSD: Single Shot MultiBox Detector
<https://arxiv.org/abs/1512.02325>
- [8] Tsung-Yi Lin et al,
Focal Loss for Dense Object Detection
<https://arxiv.org/abs/1708.02002>
- [9] Pedro F. Felzenszwalb & Daniel P. Huttenlocher
Efficient Graph-Based Image Segmentation
International Journal of Computer Vision 59, 167–181 (2004).
<https://doi.org/10.1023/B:VISI.0000022288.19776.77>
- [10] Vision meets drones: A challenge, Zhu, Pengfei and Wen, Longyin and Bian, Xiao and Ling, Haibin and Hu, Qinghua, arXiv preprint arXiv:1804.07437, 2018
- [11] Vision Meets Drones: Past, Present and Future, Zhu, Pengfei and Wen, Longyin and Du, Dawei and Bian, Xiao and Hu, Qinghua and Ling, Haibin, arXiv preprint arXiv:2001.06303, 2020

- 252 [12] Massa, Francisco and Girshick, Ross, maskrcnn-benchmark: Fast, modular reference imple-
253 mentation of Instance Segmentation and Object Detection algorithms in PyTorch, 2018
254 <https://github.com/facebookresearch/maskrcnn-benchmark>
- 255 [13] I. Mademlis, V. Mygdalis, N.Nikolaidis, M. Montagnuolo, F. Negro, A. Messina and I.Pitas,
256 “High-Level Multiple-UAV Cinematography Tools for Covering Outdoor Events”, IEEE Trans-
257 actions on Broadcasting, vol. 65, no. 3, pp. 627-635, 2019. I. Mademlis, N.Nikolaidis, A.Tefas,
258 I.Pitas, T. Wagner and A. Messina, “Autonomous UAV Cinematography: A Tutorial and a For-
259 malized Shot-Type Taxonomy”, ACM Computing Surveys, vol. 52, issue 5, pp. 105:1-105:33,
260 2019