

# Histograms of Oriented Gradients

Edgar A. Margffoy-Tuay  
Universidad de los Andes  
201412566

ea.margffoy10@uniandes.edu.co

## Abstract

*Histograms of Oriented Gradients (HOG) are image features that describe and capture local orientation information on an image. Reminiscent of SIFT features, HOG binarizes the gradient responses on a local window cell, which allow to describe local objects based on their gradient filter responses, based on this features, it is possible to distinguish different objects and eventually, detect them. The objective of the present report consists on designing and testing a multiscale object detector based on HOG features, similar to the baseline approach proposed by Dalal and Triggs adapted to face detection.*

## 1. Introduction

HOG features were proposed by Dalal and Triggs on [1] to characterize and generalize pedestrian outlines and contour gradient orientations suitable to detect them on a street scene, initially, HOG descriptors were only employed on object categories that presented a similar geometry, point of view and aspect ratio, which favored the generalization of a single HOG activation template that could be applied onto a test image via convolution. Due to the similar geometry, aspect ratio and orientation present on pedestrian pictures, HOG representation was suitable to solve this binary classification problem, as Dalal and Triggs demonstrated.

However, if simple HOG features were to be applied to solve a complex multiclass object that presents several of the initial HOG limitations such as viewpoint translation, geometric obfuscation and intra-class instance differences, the expected classification accuracy rate should be low, due to the introduction of several false positive detections due to the similarity between the HOG representations of different objects that should be represented differently.

To solve those limitations, several HOG improvements were proposed, such as Deformable Part Models (DPM) [2], on which not only a single local object HOG representation is used but also the relation with nearby HOG features. This approach allows to represent an object as a graph of HOG features, giving more information about geometry, location and viewpoint of an image. Another approach consists on training a single SVM classifier per each HOG object instance, this approach, denominated Exemplar SVMs [4], allows not only to classify each object instance accurately, but it also allows to obtain geometrical information suitable to more advanced tasks, such as 3D reconstruction and pose estimation, however, this method is too expensive due to the formulation and minimization of a single classification model per each example present on the dataset, which means that this approach is not scalable, however it is highly parallel, due to the independence of each exemplar classification model.

To explore the limitations of a plain HOG detector on a complex detection task, such as face detection, a simple Dalal-Triggs inspired framework is proposed to solve this task over the Wider Face detection dataset [6], a modern state-of-the-art face detection benchmark.

## 2. Materials and Methods

HOG feature calculation is similar to SIFT [3] representation, the main difference between both descriptors is related to gradient histogram normalization across a segmentation of a window on a set of blocks, this procedure allows to decorrelate gradient histograms on different blocks members of a single HOG window, which allows to introduce more redundancy onto the model, increasing the accuracy of the final classification model.

To detect multiple object instances on each image, a normalized sliding window of size  $136 \times 100$  approach was used, this decision was taken taking in account all the possible aspect ratios of the bounding boxes present on the dataset, some of them are elongated, and other have a similar proportion. This window was also applied over different scales, to ensure that smaller background or near faces were also subject to examination. Finally, the detection pipeline was implemented without any addition nor any improvement such as DPMs or Exemplar SVMs. To calculate HOG histograms across each of the dataset images, a cell of size 8 was defined.

The model training was done using 5 iterations of hard negative mining, by selecting initial random regions present on all images. All the implementation was based upon Andrea's Vedaldi vlfeat and matconvnet libraries [5], as shown on the Oxford's VGG tutorial on object detection<sup>1</sup>.

<sup>1</sup>Available at: <http://www.robots.ox.ac.uk/~vgg/practicals/category-detection/>

## About the Datasets

- **Wider Face:** This dataset contains 32203 images, which in turn contains 393703 annotated faces (Bounding boxes) grouped on 60 event categories. All detection evaluations are based on the Intersection-over-Union (IoU) metric.

## 3. Results

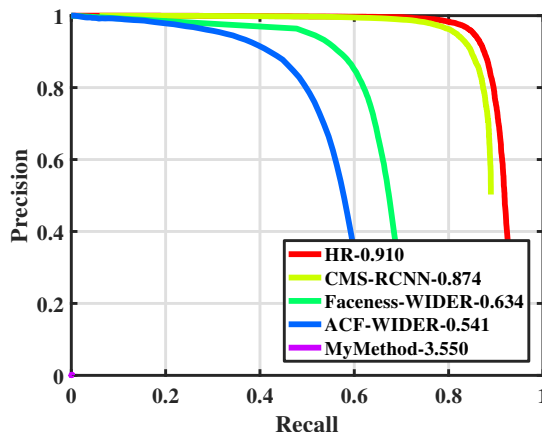


Figure 1: Precision-Recall curve results of the HOG proposed model

## 4. Results

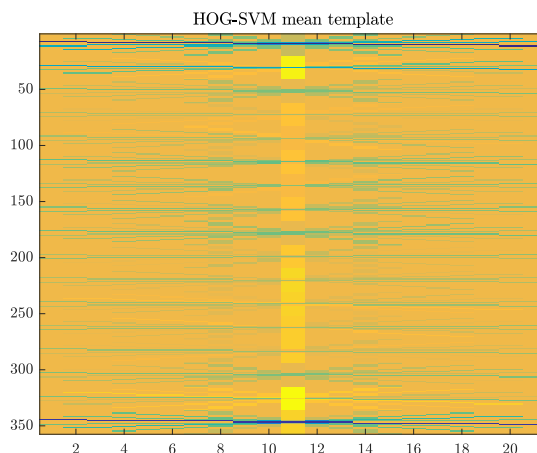


Figure 2: Mean HOG-SVM template result after evaluating the dataset using Hard-Negative Mining

As it can be shown on Figure 1, the proposed HOG detection model only outputs random boundingboxes, and it is no more different from a random coordinate sampler over all the 3302 validation set images. This implies that the multiscale HOG approach taken to solve the detection problem is not adapting well to the geometric, photometric and occlusion features present on the image, this phenomena can be confirmed after visualizing the model HOG templated that was generalized after evaluating a linear SVM based classification model using hard-negative mining.

As it can be observed on Figure 2, the mean hog model does not adapt nor represents the average expect face detector HOG template, which should contain defined image oriented features, such as round eyes and well defined image boundaries. Instead the model generalizes and adapts to distinguish a single vertical region, which is related to the limited detection contours computed as part of the HOG model, the width of all the detections was never larger than 4px, which confirms and allows to explain the low accuracy of the proposed approach.

As it was discussed previously, all the different occlusions, blur and geometric transformations present on the Wider dataset allow to propose several different HOG representations, which can not be generalized onto a single HOG template that can be used as a filter over all faces present on an image. With respect to the INRIA pedestrian dataset, there are images that present several faces with different anomalies such as image cluttering and occlusion, To solve this limitations, it should be possible to propose a model based on Exemplar SVMs, which should not only be able to detect faces present on an image, but it also should be able to recognize an specific individual and estimate its face pose.

## 5. Conclusions

HOG features are simple and very expressive local image features that allows to describe an object viewpoint in terms of its oriented gradients,

which in turn, are very similar to SIFT features. Throughout the present results it was possible to confirm and demonstrate the overall limitations of using a simple multiscale HOG detector, which is prone to errors due to its dependency on constant geometrical features such as scale, viewpoint rotation, image cluttering or blurring, which causes HOG descriptors to be variant with respect to these variables.

To account and improve the accuracy of a HOG-based object detector, it is possible to add more information related to the composition of each object, such as DPMs. Which allow to break with several of the limitations noted previously. However, even with a part-based detection model, the HOG feature representation will still present several drawbacks and limitations in terms of the representation power due to the information loss caused by this transformation.

## References

- [1] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, volume 1, pages 886–893 vol. 1, June 2005.
- [2] 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.
- [3] D. G. Lowe. Object recognition from local scale-invariant features. In *Proceedings of the International Conference on Computer Vision-Volume 2 - Volume 2, ICCV '99*, pages 1150–, Washington, DC, USA, 1999. IEEE Computer Society.
- [4] T. Malisiewicz, A. Gupta, and A. A. Efros. Ensemble of exemplar-svms for object detection and beyond. In *Computer Vision (ICCV), 2011 IEEE International Conference on*, pages 89–96. IEEE, 2011.
- [5] A. Vedaldi and B. Fulkerson. VLFeat: An open and portable library of computer vision algorithms. <http://www.vlfeat.org/>, 2008.
- [6] S. Yang, P. Luo, C.-C. Loy, and X. Tang. Wider face: A face detection benchmark. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5525–5533, 2016.