

---

# Report - Autonomous Driving: Object Detection

---

**Zhen Li**

Department of Computer Science  
University of Toronto  
Toronto, ON, M5S 2E4  
zhen@cs.toronto.edu

**Zhicong Lu**

Department of Computer Science  
University of Toronto  
Toronto, ON, M5S 2E4  
luzhc@cs.toronto.edu

## 1 Introduction

Object detection has always been one of the core problems of Computer Vision and Machine Learning. Based on object detection, autonomous driving is the most challenging topic in this field which has a great impact on our life. The goal of this project is to detect the cars and pedestrians in the real road images, and locate the 2D bounding boxes of the objects correctly. Such detection techniques would make vision-based autonomous driving systems more robust and accurate, hence increasing the possibility of adopting them in the real life.

It is hard to apply classical machine learning techniques on this topic directly. Compared with the traditional object recognition database, for example, the MNIST [1] Database, the real road images have more occlusions, more complicated background textures, and as a result, objects from the real road images are harder to detect.

The basic idea of this project came from The KITTI Vision Benchmark Suite [2]. It provides a well labeled training set containing 7481 color images. Inspired by the ranking of different methods on the KITTI website, we would try to implement those with very good performance and without using other sensor data, for the reason that such methods are more accessible and have better potential to be implemented even on mobile devices. We would try to improve the training model based on the findings of the state-of-the-art to come up with our own method, and evaluate the performance of it. We expect that our method can reach a high accuracy on the test set with an optimized speed.

In this project, we focus on the performance of SVM, Regionlets [3], and CNN techniques as well as their differences. Our code is available on GitHub<sup>1</sup>.

## 2 Related Work

Object detection, especially cars and pedestrians detections for autonomous driving systems, has been a hot topic in computer vision for recent years. In many tasks, since the number of images and windows to evaluate is huge, we often rely on a weak classifier to get proposals for the more expensive classifier. Selective Search [4] is a successful algorithm, which emphasize recall to include all image fragments of potential relevance. It can be combined with many feature representation techniques, such as the histograms of oriented gradients (HOG) [5], and SIFT [6].

[TODO: introduce regionlets]

Convolutional Neural Network (CNN or ConvNet) [7] is able to learn the features of the object and handle variations such as poses, viewpoints, and lightings, with high accuracy and high efficiency. However, it doesn't perform well when occlusion occurs, which is often the case in pedestrian and cyclists detection.

---

<sup>1</sup><https://github.com/CommanderLee/ObjectDetection>

Recently, the Fully Convolutional Neural Network (FCN) based methods[5], with end-to-end approach of learning model parameters and image features, further improves the performance of object detection. DenseBox [8] is a unified end-to-end FCN that directly predicts bounding boxes and object class confidences through all locations and scales of an image with great accuracy and efficiency. It also incorporates with landmark localization during multi-task learning and further improves object detection accuracy. It has the best accuracy on car detection on KITTI by the time the proposal is finished. However, it has not been tested on the tasks of pedestrian or cyclists detection.

Felzenszwalb et al [9] combined a margin-sensitive approach for data-mining hard negative samples, which can be used by iteratively adding hard negative examples (false positive errors) [4].

### 3 Methods

#### 3.1 Object Detection System

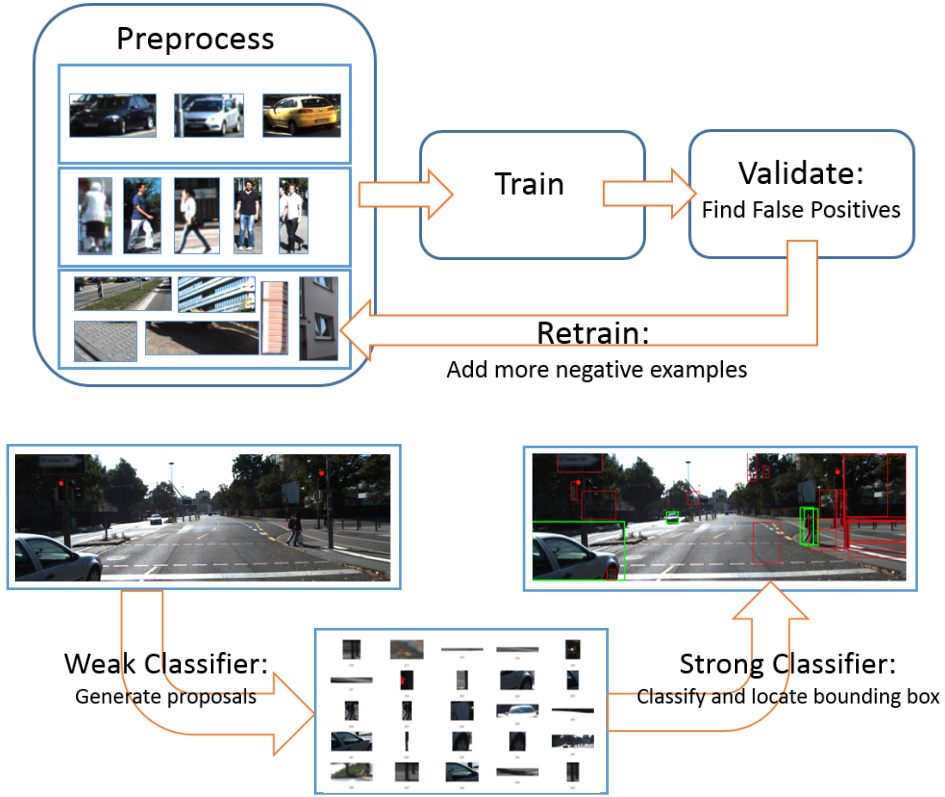


Figure 1: Top: The pipeline for training procedure. Bottom: The pipeline for testing procedure.

The whole architecture of our detection system is illustrated in Figure 1. In the training procedure, cropped images with ground truth labels are utilized to train the initial model. After validation, false positive errors are added to the negative samples, which iteratively focus on harder examples.

In the testing procedure, there are mainly two steps in our detection procedure:

Step 1 Weak classifier: Generate proposals and save the bounding box.

Step 2 Strong classifier: Find out cars and pedestrians, and save the results.

We use the Selective Search [4] as the common weak classifier. For the strong classifier after it, we use SVM, Regionlets, and CNN.

### 3.2 Ground Truth Generation

In training, image segments of cars and pedestrians are cropped and resized to  $64 \times 64$ , according to the corresponding labels. Then we adjust the original image by adding a random bias to the histogram map (from  $-20\%$  to  $+20\%$ ) repeatedly. We also reverse the original image horizontally to make full use of the training set. This adjustment will not only balance the different light conditions in the original images, but also make full use of the hidden information.

### 3.3 Negative Samples Generation

We use the Selective Search [4] to generate the negative samples. Selective Search is a segmentation algorithm that focus on recall rate and ensure its results do not stem from parameter tuning. It starts from initial regions generated by [10], and use a greedy algorithm to repeatedly merge similar segments. The similarity function is defined to combine the size similarity, by calculating the fraction of joint area, and the texture similarity, by calculating the histogram intersection.

We randomly select images from the training set, run Selective Search, and find out background segments, which are defined to have a 0% to 30% overlap with a positive sample (car or pedestrian). We generate 20000 negative samples for training.

### 3.4 Retrain Procedure

Generating the negative samples is a important part of road object detection, since the background textures are too complicated to cluster easily. Too enhance the model's ability to solve hard problems, we add retrain procedure [9] to iteratively adding false positive errors to the negative sample set.

## 4 Experiments

### 4.1 Data set

We use the KITTI data set for Object Detection [2], which contains 7481 color images with ground truth bounding box labels, including 28782 cars and 4487 pedestrians. We partition the data set to a training set contains 5237 images(70%), and a testing set contains 2244 images(30%). Then we run different algorithms on the data set and analyze the results.

### 4.2 Regionlets

TODO: modify format, add subsection title, add reference(I have add the ref.bib file, so we can use the auto-ref now), and change math font

Regionlets with Adaboost for generic object detection

To try some different features as well as different classifiers, we implemented regionlets for generic object detection [] to compare with the SVM classifier. We also use selective search to get candidate bounding boxes which may contain target object, and classify the candidate boxes with the learned classifier.

Regionlets

Regionlet defination Regionlets are defined as sub-parts of a region in the image. They are introduced to act as basic units to extract appearance features, and are organized into small groups to describe features with different degrees of deformation. In this way, features extracted from small regions and a big region can work together, which can provide a good localization ability as well as tolerate more variations. Figure [] shows the definition of regionlets. The outer black rectangle is a candidate bounding box. The blue rectangle inside the bounding box is a feature extraction region denoted as R, which will contribute a weak classifier to the boosting classifier. Within the region R, some small sub-regions are selected and defined as a group of regionlets, as shown in Figure []. The features of these regionlets will be aggregated to a single feature for a region R during training, and a bounding box will be represented by a number of such regions. By introducing regionlets, it is

straightforward to think that it would improve object recognition especially for occluded situations, which are common in automated driving, because the discerning features are extracted while irrelevant appearance from background is largely discarded. Besides, more regionlets in a single region R will increase the capacity to model deformations.

**Region feature extraction** According to [], feature extraction from R takes 2 steps. First, we extract appearance features of HOG descriptors from each regionlets respectively. Second, we generate the representation of R based on regionlets features. We apply max-pooling over regionlet features for the feature of region R. Denote  $T(R)$  as the feature representation for region R,  $T(r_j)$  as the feature extracted from the  $j$ th regionlet  $r_j$  in R, then the operation is defined as follows:  $T(R) = \max T(r_j)$ . The max-pooling happens for each feature dimension independently. For each regionlet, we first extract HOG feature for the regionlet. Then we pick a 1D feature from the same dimension of HOG feature in each regionlet and apply max-pooling to get the feature for region R. We have millions of such 1D features in a detection window and the most discriminative ones are determined through a boosting learning process.

**Learning the object detection model** We use the ensemble method of boosting to learn the discriminative regionlet groups and their configurations from a huge pool of candidate regions and regionlets.

**Regionlets pool construction** To deal with deformation at different scales, we first build a largely over-complete pool for regions and regionlets with various positions, aspect ratios and sizes. we denote the 1D feature of a region relative to a bounding box as  $R = (l, t, r, b, k)$ , where  $k$  denotes the  $k$ th element of the low-level feature vector of the region. The region pool is spanned by  $X \times Y \times W \times H \times F$ , where  $X$  and  $Y$  are respectively the space of horizontal and vertical anchor position of R in the detection window,  $W$  and  $H$  are the width and height of the feature extraction region R, and  $F$  is the space of HOG feature. Enumerating all possible regions is impractical and not necessary. We employ a sampling process to reduce the pool size. The algorithm is the same as Algorithm 1 in [].

After getting the region pool, we propose a set of regionlets with random positions inside each region. Although the sizes of regionlets in a region could be arbitrary in general, we restrict regionlets in a group to have the identical size because our regionlets are designed to capture the same appearance in different possible locations due to deformation. The sizes of regionlets in different groups could be different. A region may contain up to 5 regionlets in our implementation.

The final feature space used as the feature pool for boosting is spanned by  $R \times C$ , where  $R$  is the region feature prototype space and  $C$  is the configuration space of regionlets.

**Training with boosting regionlet features**

We use Gentle Adaboost[] to train classifiers for our object detector. One boosting classifier consists of a set of selected weak classifiers. We define the weak classifier as a decision tree. Gentle Adaboost puts less weight on outlier data points, which would make the classifiers more robust.

### 4.3 SVM

#### 4.3.1 Training

We noticed that there are far more cars than pedestrians in the KITTI data set. To balance the difference, the multiplication trick in Section 3.2 can be used to increase the number of pedestrians, which make better use of the training set. In addition, since the KITTI benchmark only evaluate objects larger than 25 pixels (height), we ignore these small objects in the training set. After all of these augmenting and filtering strategies, we have 22576 cars and 20960 pedestrians as positive samples for training.

The SVM classifier is trained using the cropped positive samples with ground truth labels, and negative background segments generated by the Selective Search. Since SVM is a binary classifier, we trained 3 classifiers: (1) Car vs All, (2) Pedestrian vs All, and (3) Car vs Pedestrian. We use the `fitcsvm` function from the MATLAB toolbox to train out model.

#### 4.3.2 Validation

We use the  $k$ -fold ( $k = 7$ ) cross validation to evaluate our model, as required by the project guide (60% training, 10% validation, and 30% testing). According to our description in Section 3.2, we

multiply the data set by changing the contrast and brightness of the images, as well as reversing the images horizontally. Then we can compare the cross loss before and after the multiplication, using the `crossval` function from the MATLAB toolbox.

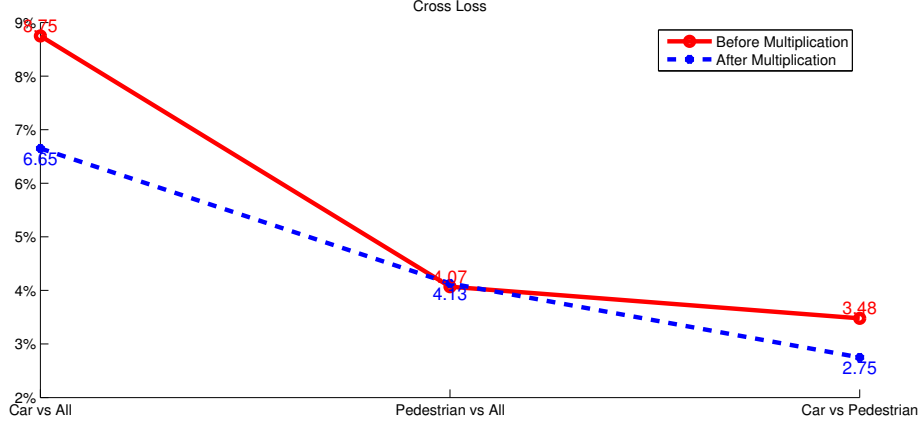


Figure 2: Cross validation ( $k = 7$ ) loss before and after the image multiplication.

We can observe from Figure 2 that the cross loss gets lower when we multiply the data set from 13908 samples (11288 cars and 2620 pedestrians) to 43536 samples (22576 cars and 20960 pedestrians). With this trick, we can make full use of the limited data set when we lack certain classes of objects. On the other hand, this trick also helps to balance the different light conditions among the images by adding random brightness bias. The final cross loss range from 2.75% to 6.65%, which is satisfying.

In addition to the randomly selected negative samples ( $N_1 = 20000$ ), we add retrain step to enhance the training procedure. Though the number of positive samples is limited by the training set, we can crop much more negative samples from the existing images. But we need to balance that with the training cost. To be more effective, we save the false positive image segments to files during the validation, and add these images to the training set, labeled as the background image. These are considered to be the hard samples ( $N_2 = 10500$ ), so finally we have 30500 negative samples.

### 4.3.3 Testing

First we can test on the cropped testing set to evaluate the object recognition performance. Similarly, we can compare the different accuracy, precision, and recall rate we got before the image multiplication and after the image multiplication. We use the `predict` function from the MATLAB toolbox to classify with our trained SVM models.

From Figure 3, we can conclude that multiplying images will also increase the accuracy, precision, and recall rate. Using 43536 positive samples and 20000 negative samples, we obtain a precision rate range from 92.96% to 98.26% and a recall rate range from 88.72% to 96.56%. It is reasonable to consider that if we collect more training data or multiply more images with existing data, we will be able to get better performance. But that will also leads to more expensive training and testing cost.

However, it doesn't perform well on the real test set, especially when the cars are small enough. We can observe from Figure 4 that the proposals are relatively big compared to the small cars. However, if we decrease the expected size of bounding box, we will miss the big car that close to the observer.

## 4.4 CNN

TODO: some description of test

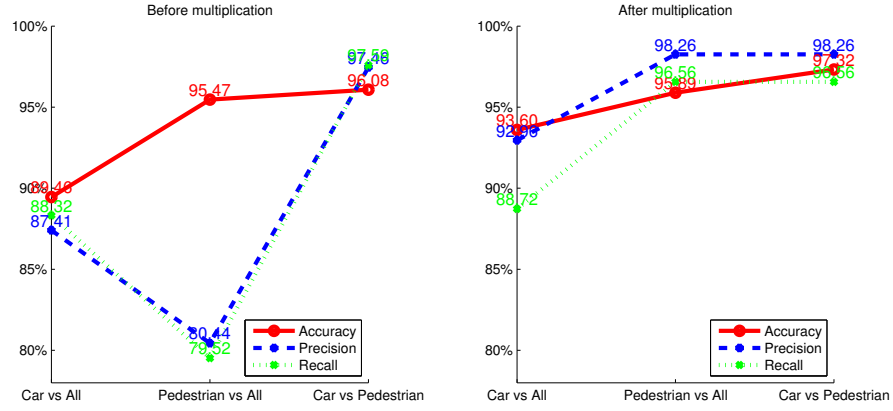


Figure 3: Accuracy, precision, and recall rate before and after the image multiplication.

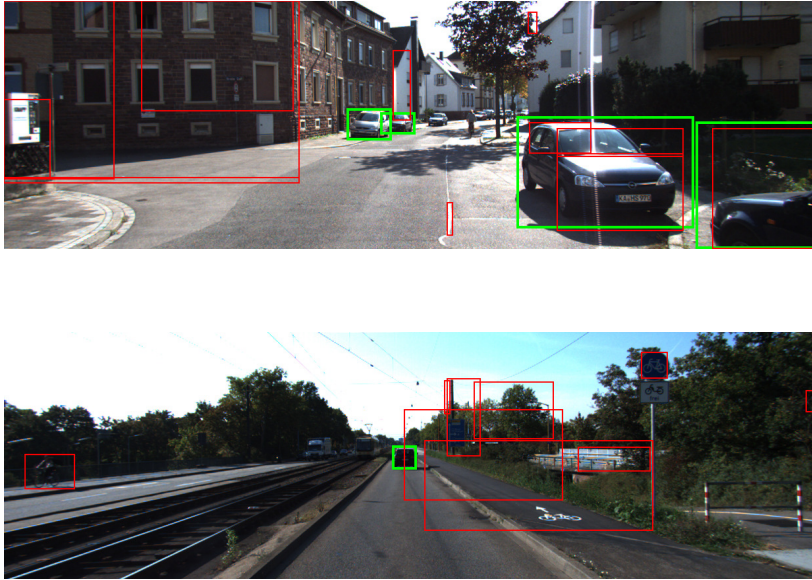


Figure 4: Red bounding box is output result of the classifier, and the green bounding box is the ground truth. Top: Found the two big cars on the right, but lost other two small cars in the middle. Bottom: Found the cyclist on the left, but lost the small car in the middle.

## 5 Conclusion

[old related work, maybe changed to future work :)]

[old]Besides KITTI, we would also use some other datasets to test our methods, especially on pedestrian detection, including Caltech Pedestrian Detection Benchmark[7] and Daimler Pedestrian Segmentation Benchmark Dataset[1].

[old]Inspired by the ranking of different methods on the KITTI website, we would try to implement those with very good performance and without using other sensor data, for the reason that such methods are more accessible and have better potential to be implemented even on mobile devices. We would try to improve the training model based on the findings of the state-of-the-art to come

up with our own method, and evaluate the performance of it. We expect that our method can reach a high accuracy on the test set with an optimized speed.

## References

- [1] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [2] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
- [3] X. Wang, M. Yang, S. Zhu, and Y. Lin, "Regionlets for Generic Object Detection," in *2013 IEEE International Conference on Computer Vision*. IEEE, dec 2013, pp. 17–24.
- [4] K. E. Van de Sande, J. R. Uijlings, T. Gevers, and A. W. Smeulders, "Segmentation as selective search for object recognition," in *Computer Vision (ICCV), 2011 IEEE International Conference on*. IEEE, 2011, pp. 1879–1886.
- [5] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, vol. 1. IEEE, 2005, pp. 886–893.
- [6] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International journal of computer vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [7] A. Krizhevsky, I. Sutskever, and G. Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in Neural Information Processing Systems 25*, 2012.
- [8] L. Huang, Y. Yang, Y. Deng, and Y. Yu, "DenseBox: Unifying Landmark Localization with End to End Object Detection," *arXiv*, sep 2015. [Online]. Available: <http://arxiv.org/abs/1509.04874>
- [9] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 32, no. 9, pp. 1627–1645, 2010.
- [10] P. F. Felzenszwalb and D. P. Huttenlocher, "Efficient graph-based image segmentation," *International Journal of Computer Vision*, vol. 59, no. 2, pp. 167–181, 2004.

## Old References

- [1] C. G. Keller, M. Enzweiler, and D. M. Gavrila, A new benchmark for stereo-based pedestrian detection, in 2011 IEEE Intelligent Vehicles Symposium (IV), 2011, pp. 691696.
- [2] A. Geiger, P. Lenz, and R. Urtasun, Are we ready for autonomous driving? The KITTI vision benchmark suite, in 2012 IEEE Conference on Computer Vision and Pattern Recognition, 2012, pp. 33543361.
- [3] Y. Tian, P. Luo, X. Wang, and X. Tang, Deep Learning Strong Parts for Pedestrian Detection, in 2015 IEEE International Conference on Computer Vision, 2015.
- [4] L. Huang, Y. Yang, Y. Deng, and Y. Yu, DenseBox: Unifying Landmark Localization with End to End Object Detection, *arXiv:1509.04874*, 2015.
- [5] J. Long, E. Shelhamer, and T. Darrell, Fully convolutional networks for semantic segmentation, *arXiv1411.4038*, 2014.
- [6] A. Krizhevsky, I. Sutskever, and G. Hinton, Imagenet classification with deep convolutional neural networks, *Adv. Neural Inf. Process. Syst.* 25, 2012.
- [7] P. Dollar, C. Wojek, B. Schiele, and P. Perona, Pedestrian detection: A benchmark, in 2009 IEEE Conference on Computer Vision and Pattern Recognition, 2009, pp. 304311.
- [8] K. Sande, J. Uijlings, T. Gevers, and A. Smeulders, Segmentation as Selective Search for Object Recognition, in 2011 IEEE International Conference on Computer Vision, pp. 1879-1886.