

---

# Objects Recognition Based on HOG with SVM

## Final Report for CIS 519

---

Shangyi Cheng  
Yao Chu  
Chenyang Zhao

SHANGYI@SEAS.UPENN.EDU  
CHUYAO@SEAS.UPENN.EDU  
CHZHAO@SEAS.UPENN.EDU

### Abstract

We studied the question of object recognition using histograms of oriented gradients (HOG) and support vector machine (SVM) with Gaussian Kernel on ball detection as a test case. After reviewing several existing methods for feature extraction, our project verified that HOG performs well in ball recognition. The whole process is shown in this report, from choosing regions of interest manually, extracting features from these regions using HOG, using principal component analysis (PCA) to reduce dimension of the dataset and then applying SVM with Gaussian kernel to classify.

## 1. Introduction

The topic for this project is to detect the object of our interest on the given images and predict which class it belongs to. The dataset we use is from Caltech 256(Griffin, G. Holub, AD. Perona, P.), which contain images of four kinds of balls, including 98 images of golf in the folder “088.golf-ball”, 174 images of soccer ball in “193.soccer-ball”, 104 images of bowling-ball in “017.bowling-ball” and 98 images of tennis ball in “216.tennis-ball”. As shown in Figure 2, among images of the same kind (take soccer-ball as example), some of the images have a whole contour of soccer in the center and fill the image, while in other images(such as a photo of a soccer game), the soccer may cover a small portion in the corner, or only part of a soccer or a group of soccers appear on the image. We decide to use the image with a single target occupying most of the image as the training set. So the first thing we need to do before training is to create such standard images for training from the raw dataset. And our final goal is, given a new image, the model can detect whether this image contains any ball of our interest or just unrelated to balls and in the first condi-

tion, the model can give the label for the class the detected object belongs to.

## 2. Methodology

### 2.1. Overview of the Method

Inspired by the paper (Dalal & Triggs, 2005) presented by Navneet Dalal and Bill Triggs, we tried the HOG method to extract the features vector for each individual image and apply SVM learner to train the dataset. An overview of the whole process is shown in Figure 1.

For the training process, we first preprocess the images by selecting the region containing the ball manually, resize and restore such images as the dataset for feature extraction. Then, HOG is applied on those processes images and get a vector of length 900 as features for one instance. Thus, we reconstruct the training dataset by standardizing the dataset and apply PCA method to reduce the number of features. Then, SVM with Gaussian kernel is used to learn from the dataset. The metrics of accuracy, ROC and learning curve along with cross-validation are used to ensure that the model has learned to classify four kinds of balls with an acceptable accuracy without overfitting.

While detecting ball in an brand-new image, the model scan the image using mask with increasing size, doing the corresponding manipulation with the same parameters as those for training set, and then make prediction for each block no matter whether it contains any ball of our interest, what the portion of the ball covers the image and whether it appears in the center or in the corner, individually or in group. Among all predictions, the model will choose the ones with high confidence and decide whether they point to the same object. Finally, the target of our interest will be framed in a rectangle along with its label.

### 2.2. Region of Interest (ROI) Selection

Before we use HOG to extract features from labeled images, we need to select the ROI which contains an object we would like our model to recognize. We tried several ways to select such region, including Harris corner detec-

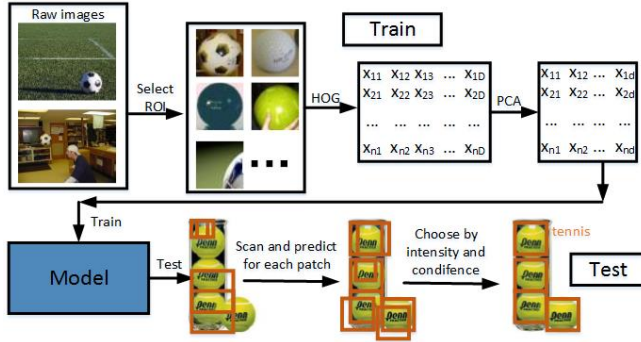


Figure 1. An illustration of HOG with SVM process.

tion, Hough transfer, (Brown et al., 2005), (Atherton T.J., 1999), (Victor A., 2006), (D'Orazio T., 2004).

To increase the number of instances in our training data, for each image, we firstly decide the a rectangle ROI region and then shift the area a little in eight directions respectively. Therefore, a raw image will create at most nine images and at least one image for training. Finally, the picked area is resized to a  $40 \times 40$  pixal square. Apart from the four kinds of balls, we also generate some squares with the label "Unknown" to represent the unrelated features to ball. The images in this class are mostly selected from the background, fragments of the ball and other unrelated area. Figure 2 shows the processed images after manual selection and resizing. First two rows are images for balls and the third row is for "Unknown".

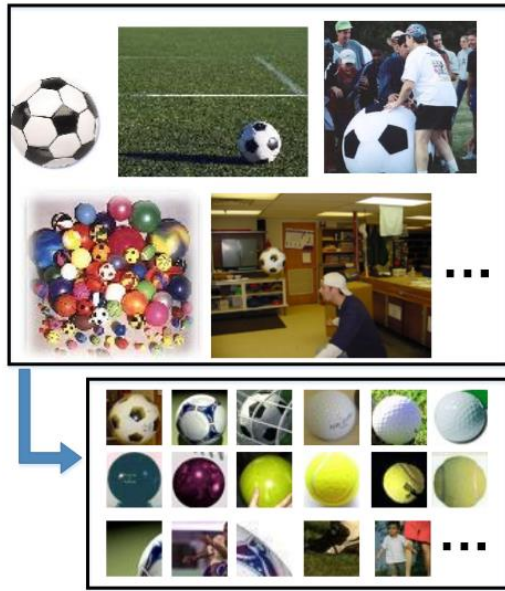


Figure 2. Resized Images Before & After ROI Selection

### 2.3. Histogram of Oriented Gradients Descriptors

As described in ROI selection section, the patches cropped out were resized to  $40 \times 40$  pixels. Each patch works as an dataset instance, however, there are two primary ways to extract features from these patches. First, each pixel value (grayscale, RGB or LAB color spaces) can be extracted directly as a feature. Second, apply HOG or SIFT to extract features. Considering the HOG/SIFT representation can capture edges and gradient structure which is characteristic of local shape as well as easily maintain the invariance to local geometric and photometric transformation, we finally chose second method. Also, HOG method can be implemented easily with fewer features compared to SIFT, so we chose to use HOG descriptors.

Implementing HOG to extract features from patches followed procedure shown in figure 3. For balls detection and classifier, following coefficient were chosen: All patches were converted into grayscale color space to normalize gamma and color; 1-D point derivative (centered  $[-1, 0, 1]$ ) mask with no smoothing was applied to patches; 9 orientation bins evenly spaced over  $0-180$  degrees (unsigned gradient) was chosen; Cell size  $6 \times 6$  and block size  $2 \times 2$  were chosen for R-HOG.

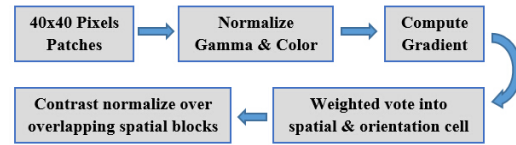


Figure 3. HOG Procedure

HOG features is shown in figure 4. After extract HOG features from a patch, reshape the feature matrix into a vector represent the patch. Eventually we generated a  $3660 \times 900$  datasets with 3660 instances and 900 features.

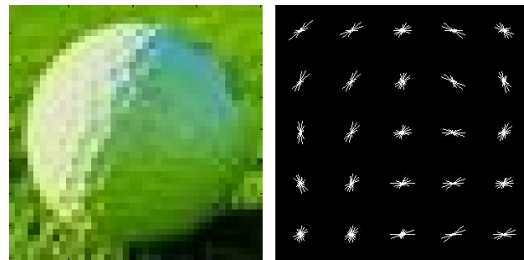


Figure 4. HOG Feature Extraction Result

### 2.4. PCA

The number of features we get using HOG is 900. Since the number of instances we use is about 4000, which is not large enough for 900 features and easy to result in over-

fit, it's necessary to reduce the dimension first. On the other hand, applying PCA would improve efficiency and decrease the time for training. On the training dataset, we firstly standardize the dataset  $X$  and record covariance matrix  $\Sigma$  and mean values  $\mu$ . Next, we choose the  $d$  most important PCA basis vectors by calculating the eigenvalues and eigenvectors of  $X^T X$ . Then we reconstruct the new training dataset by multiply the PCA basis vectors to the origin dataset  $X$ .

## 2.5. One-Class SVM

Support vector machine(SVM) is one of basic learning algorithm that most widely used to solve object recognition problems (Dalal & Triggs, 2005), (Massimiliano P., 1998). In our case, we applied one-class SVM first and extend it to multi-class to train the dataset with several different labels. Given a set of instances which belongs to either of two classes, a SVM classifier finds the optimal hyperplane leaving the largest possible fraction of instances of the same class on the same side, while maximize the distance of either class from hyperplane. The objective function of SV learner is

$$\arg\min_{\theta} \frac{1}{2} \sum_{i=1}^d \theta_i^2 \quad (1)$$

s.t.  $y_j(\theta \mathbf{x}_j) \geq 1 \forall j$ .

The problem of minimizing the cost function could be simplified as maximizing the  $J(\alpha)$  in 2.

$$J(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle \quad (2)$$

s.t.  $\alpha_i \leq 0 \forall i$ ,  $\sum_i \alpha_i y_i = 0$ , where  $\alpha_i$  is the constraints weight scaler, and  $\langle x_i, x_j \rangle$  is a scaler given by the kernel function.

Take linear kernel as an example, the parameter of the hyperplane follows that

$$\theta = \sum_{i=1}^N \alpha_i y_i x_i \quad (3)$$

while  $b$  can be determined from  $\alpha$  solution of the dual problem, and from the Khn-Tucker conditions

$$\alpha_i (y_i (\theta x_i + b) - 1) = 0, i = 1, 2, \dots, N \quad (4)$$

The problem of classifying a new data instance could be simply solved by looking at the sign of the score function

$$P(\mathbf{x}) = \theta \mathbf{x} + \mathbf{b} \quad (5)$$

While applying SVM algorithm, choosing proper kernel is essential to train the classifier. We used two different kernels to train the dataset, which are Gaussian kernel, linear kernel, and drew the Receiver Operating Characteristic(ROC) curve respectively. The result are shown in the

figure (5,6). As shown in the figure, the two kernels give similar result while Gaussian kernel provides a faster speed to converge.

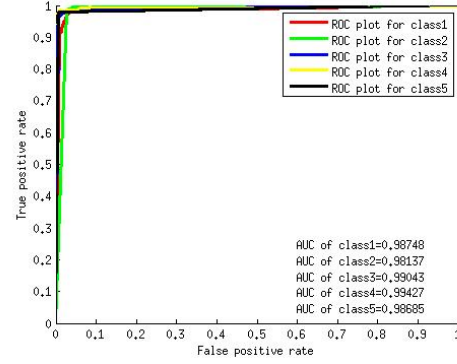


Figure 5. ROC curve of SVM with Gaussian Kernel

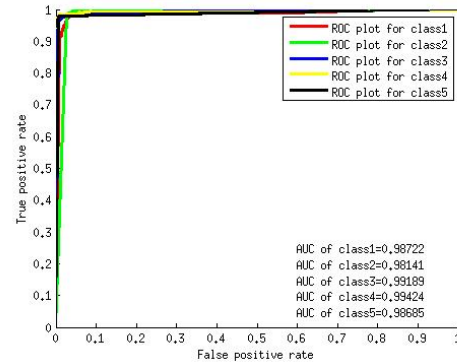


Figure 6. ROC curve of SVM with Linear Kernel

## 2.6. Multi-Class SVM

Traditional SVM algorithm only solves 2-classes problems. However, detecting ball object in a image should be able to tell whether there is a ball and figure out the type of the ball object, which is a multi-class classification problem. Therefore, one-vs-rest is introduced for solving multi-class problem.

Given a dataset with  $k$  different labels,  $k$  different one-class SVM classifiers is trained for each class. When training for individual class, the instances that belongs to this class are labeled as 1, and all the others are labeled as 0. A new instance should be predicted as in the class which gives the largest value from score function 5.

## 2.7. Detect Balls from New Image

Once classifier was trained, balls could be detected from the new input image based on the predict results.

In order to detect the balls, the algorithm should search every possible detect window in the image. So square detect window with changeable size move around in the image to extract test patches. Window size changes from minimum value (80 pixels or 10 percent of maximum side) to maximum value (small one between image width and height). And overlapping ratio between two consecutive detect windows is 85 percent. Every patch is resized to 40x40 pixels, extracted HOG features, and normalized in the same way with training data. Following this, apply classifier to predict labels for test instances. Based on the predict results, overlap test instances can be eliminated by adaptive non-maximum suppression. A neighbour positive detection threshold is set to improve the detection result.

## 3. Results and Performance Study

### 3.1. Learning Curve

To evaluate our SVM classifier, we run our classifier through 10-fold cross validation and drew the learning curve in figure 7 to show the performance versus the number of training examples. As shown in the figure, the testing accuracy increases as the set of training instances grows larger, and the variance of accuracies decreases.

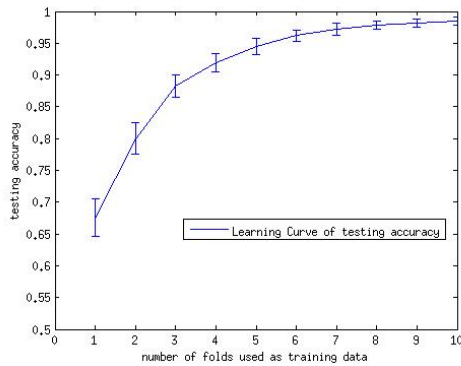


Figure 7. Learning Curve of testing accuracy

### 3.2. Precision & Recall Metric

We evaluated our learning model with the precision and recall value. The result are shown in the 1. As shown in the table, both precision and recall value could reach 0.9, which shows that our classifier model works fine in our problem.

Table 1. Confusion Matrix: B for Bowling, G for Golf, S for Soccer, T for Tennis and U for Unknown

TRUE LABEL	B	G	S	T	U	PRECISION
PREDICTION						
BOWLING	234	3	4	2	0	0.96
GOLF	21	357	10	7	10	0.88
SOCCER	12	1	393	0	1	0.97
TENNIS	6	3	4	285	1	0.96
UNKNOWN	1	1	0	0	474	0.96
RECALL	0.86	0.98	0.96	0.97	0.97	
ACCURACY	0.95					

### 3.3. Results on Raw Images

Using the model after learning, we test on several raw images, including the ones where we selected ROI for training and new images which the model never see before. In Figure 8, we show some of the results that our model detect and label the target successfully. We use different colors of rectangle to frame the area of the detected target, red for soccer, cyan for tennis, purple for bowling-ball and green for golf. In theory, our model can detect all balls with the same label on an image, such as the upper middle and right ones with several tennis balls. The lower middle image with golf is a brand-new image from website and our model can detect and label it correctly.

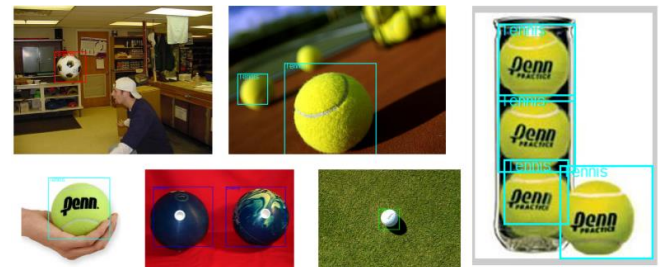


Figure 8. Successful Results for Detection and Label

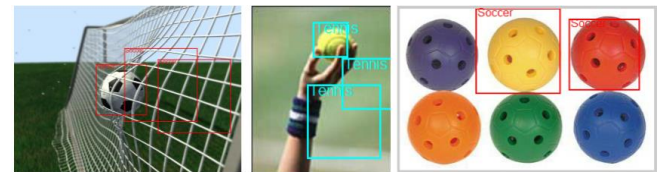


Figure 9. False Positive and True Negative Results

However, there are also some cases in which our model either fails to detect all the targets (as the right one with several bowling-balls in Figure 9), or gives false positive prediction (label some irrelevant areas as a specific kind of



ball, such as the left and middle ones). One of the possible reasons is that the target is sheltered by other objects have regular patterns, such as the soccer in the net shown Figure 9. Another reason is related to the size of the raw images. The image with a lifting hand holding a tennis has  $96 \times 96$  pixels. Since the background of it is also green, it's more likely to detect and label the background mistakenly as tennis.

#### 4. Further Work

Although our current Gaussian SVM detector is reasonably accurate for the cropped patches (test part of manually cropped ROI), there is still a problem when test thousands of test instances extracted from new input image (downloaded from other sources). False positive results were the main problem we encountered. The most likely explanation for this problem we believe is the insufficiency of the datasets, especially in the negative sets ('Unknown'). Also the original datasets for each kind ball is limited, the cropped patches generated based on the original datasets may be not representative. It would be useful to use a larger datasets to improve detector accuracy and implement cascade training method (Viola & Jones, 2001) to improve the training efficiency.

#### Acknowledgments

This project is supported by Uenn CIS 519. Thanks Eric Eaton and Xiaoxiang Hu for their instruction. Thanks MATLAB for providing so many useful toolboxes.

#### References

- Atherton T.J., Kerbyson D.J. Size invariant circle detection. pp. 795–803, 1999.
- Brown, Matthew, Szeliski, Richard, and Winder, Simon. Multi-image matching using multi-scale oriented patches. In *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, volume 1, pp. 510–517. IEEE, 2005.
- Dalal, Navneet and Triggs, Bill. Histograms of oriented gradients for human detection. In *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, volume 1, pp. 886–893. IEEE, 2005.
- D'Orazio T., Guaragnella C., Leo M. Distant A. A new algorithm for ball recognition using circle hough transform and neural classifier. pp. 393–408, 2004.
- Massimiliano P., Alessandro V. Support vector machines for 3d object recognition. pp. 637–646, 1998.

Victor A., Carlos H.G., Arturo P. Raul E.S. Circle detection on images using genetic algorithms. pp. 652–657, 2006.

Viola, Paul and Jones, Michael. Rapid object detection using a boosted cascade of simple features. In *Computer Vision and Pattern Recognition, 2001. CVPR 2001. Proceedings of the 2001 IEEE Computer Society Conference on*, volume 1, pp. 1–511. IEEE, 2001.