

Vehicle Detection using Partial Least Squares

Aniruddha Kembhavi, *Student Member, IEEE*, David Harwood, *Member, IEEE*, and Larry S. Davis, *Fellow, IEEE*

Abstract—Detecting vehicles in aerial images has a wide range of applications from urban planning to visual surveillance. We describe a vehicle detector that improves upon previous approaches by incorporating a very large and rich set of image descriptors. A new feature set called *Color Probability Maps* is used to capture the color statistics of vehicles and their surroundings, along with the Histograms of Oriented Gradients feature and a simple yet powerful image descriptor that captures the structural characteristics of objects, named *Pairs of Pixels*. The combination of these features leads to an extremely high dimensional feature set (approximately 70,000 elements). Partial Least Squares is first used to project the data onto a much lower dimensional subspace. Then, a powerful feature selection analysis is employed to improve performance, while vastly reducing the number of features that must be calculated. We compare our system to previous approaches on two challenging datasets and show superior performance.

Index Terms—Vehicle detection, partial least squares, feature selection

I. INTRODUCTION

SEVERAL commercial earth observation satellites, such as IKONOS, GeoEye and QuickBird, provide publicly available imagery at a ground sampling distance (GSD) of 1 meter. High resolution images of a small number of locations are publicly available via Google Earth at an astonishing GSD of 0.15 meter. We consider the problem of detecting vehicles from such high resolution aerial and satellite imagery; this problem has a number of applications. Images of road networks, along with the distribution of vehicles in different regions, can provide information for urban planning and traffic monitoring. Detecting and tracking vehicles in aerial videos is also an important component in visual surveillance systems. In spite of the increasing availability of high resolution aerial and satellite images, vehicle detection still remains a challenging problem. In urban settings especially, the presence of a large number of rectilinear structures, such as trash bins, electrical units and air conditioning units on the tops of buildings can cause many false alarms. Figure 1(a) shows an example of an image patch for which previously published vehicle detectors produce a large number of false alarms.

Object detection systems have typically used image descriptors such as Scale Invariant Feature Transform (SIFT) [17] and Geometric Blur [2], calculated at a number of interest points within the image. These image descriptors are then combined using various aggregating approaches such as Bags-Of-Words [32] and Spatial Pyramids [16] to provide a rich description of the object. However, such approaches have not been extensively used for the problem of vehicle detection, due to the low resolution of traditional aerial images and the need of many commonly used image descriptors for a sufficiently large support region.

Vehicle detection has been previously treated as a template matching problem, and several algorithms have been proposed

that construct templates in 2D as well as 3D. Moon et al. [20] propose an approach to optimally detect two dimensional shapes. They derive an optimal one dimensional step edge detector which minimizes the noise power and mean squared error between the input and filter output. This turns out to be the derivative of the double exponential (DODE) function. The DODE filter is then extended along the shape's boundary contour to obtain a shape detector. The problem of vehicle detection is then equivalent to *detecting parallelograms*. They show impressive results on images of vehicle parking lots. A comprehensive analysis of this vehicle detector under a wide range of operating environments was carried out in [21]. A mathematical analysis was provided to quantify the degradation of the vehicle detector with decreasing illumination and acquisition angle.

Choi et al. [5] use a mean shift based clustering algorithm to extract candidate blobs that exhibit symmetry properties of typical vehicles. Each candidate is then classified using geometric and radiometric characteristics of the blob. Eikvil et al. [7] propose a similar two stage strategy for vehicle detection in satellite images. The first stage consists of segmenting regions into potential vehicles, roads, vegetation, etc. They also leverage multi-spectral information to identify regions of vegetation and Geographical Information System (GIS) data to obtain the road network. The second stage then consists of a region classification algorithm using geometrical properties such as area, moments, etc. Zheng et al. [34] obtain vehicle candidates using a morphological pre-processing stage, which are then classified using a neural network. Such two stage approaches typically suffer from errors obtained in the segmentation stage of the system. Furthermore, geometric properties of blobs are not powerful enough to detect vehicles with high accuracy in urban settings, where the presence of a large number of rectilinear structures cause many false alarms.

Zhao et al. [33] pose the vehicle detection problem as a 3D object recognition problem. They use human knowledge to model the geometry of a typical vehicle. Psychological tests revealed that human subjects most often used cues such as the rectangular shape of the car, layout of windshields and presence of shadows to detect cars. Such cues are then integrated using a Bayesian network. They also make use of camera calibration and illumination information to predict shadow cues. While effective, their algorithm cannot be easily extended to build other object detection systems, due to the large amount of human modeling that is required. Similar 3D models have also been used to model vehicle geometries for the purpose of car detection and counting in aerial images by Hinz [15] and Schlosser et al. [25].

Grabner et al. [12] propose an online version of boosting to efficiently train their vehicle detector. Their algorithm avoids building large pre-labeled training datasets by using an active learning framework. They use three classes of features - Haar wavelets, Histograms of Oriented Gradients and Local Binary Patterns, all of which can be very efficiently calculated using Integral Images and Integral Histograms. Their detection results are further improved by segmenting the image into streets, build-

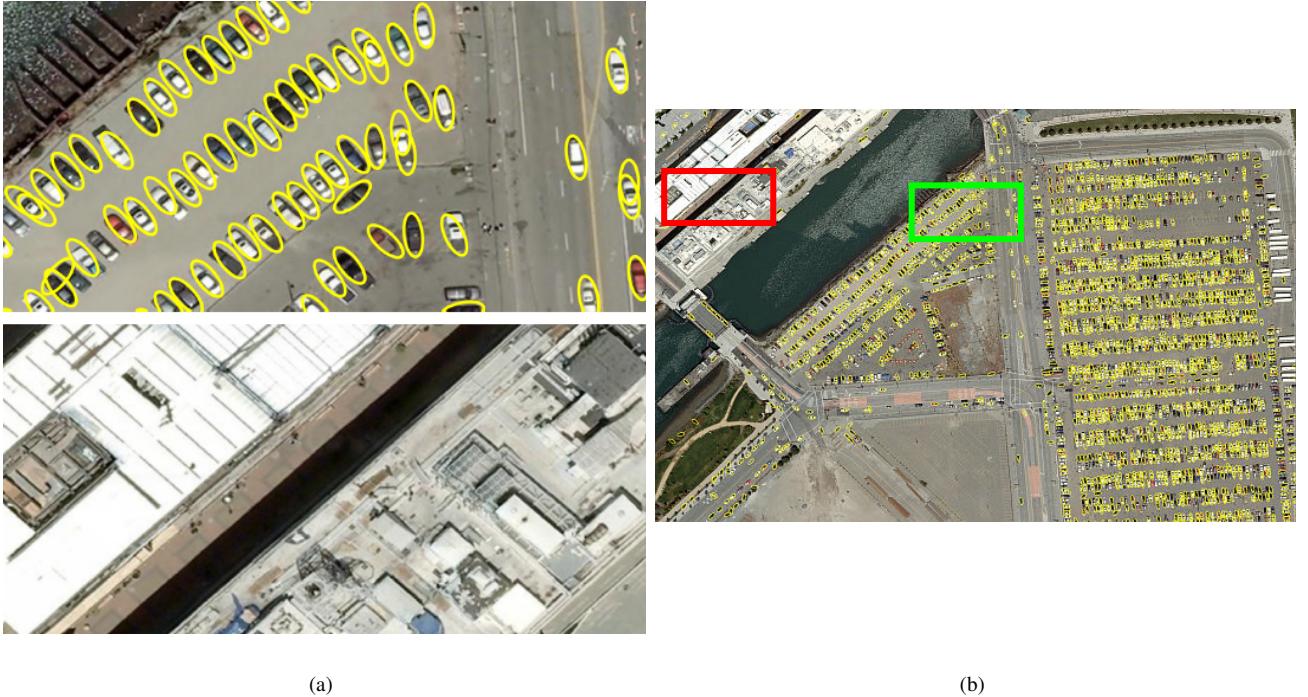


Fig. 1. (a) Two image patches showing the performance of our system. Both image patches were extracted from a much larger image (5007×7776 pixels) which is displayed in its entirety in (b) and at a higher resolution in Figure 23. The top figure (region marked in green) shows the high detection accuracy of our system in the presence of a large number of vehicles. The bottom image (region marked in red) shows the low false alarm rate of our system in a region that has many rectilinear structures. Typical vehicle detection systems often produce false alarms in the presence of such structures. © 2009 Google

ings, trees, etc. and then discarding vehicle detections that are not present on the streets.

More recently, vehicle analysis has been extended from single images to video sequences. Yue et al. [31] propose a system for vehicle verification in airborne video sequences. The vehicle of interest may leave the field of view for a while or may be obscured. When a new vehicle is observed, verification is needed to confirm whether it was the previously detected vehicle. A homography based view synthesis method is used to generate novel views of the exemplars that are provided during training. This enables the system to be robust to large aspect angle variations of the test sequence. The synthesized novel view and testing object are then compared using a weighted combination of a rotationally invariant color matcher and a spatial feature matcher.

Our proposed vehicle detector improves upon previous systems by incorporating a much larger and richer feature set than previous approaches. First, a novel set of image descriptors are proposed that capture the color properties of an object and its surrounding, called *Color Probability Maps (CPM)*. Then, the commonly used Histograms of Oriented Gradients (HOG) feature is incorporated to capture the spatial distribution of edge orientations. Finally a very simple yet powerful image descriptor, named *Pairs Of Pixels (PoP)*, is proposed to capture the structural properties of objects. The concatenation of these three classes of features leads to a very high dimensional feature vector (approximately 70,000 elements). In contrast, the number of samples that we have to train our vehicle detector is much smaller, rendering many popular machine learning techniques unusable. Furthermore, our features are extracted from neighboring pixels within a detection window, which greatly increases their multi-collinearity. We take advantage

of the nature of our problem by employing a classical statistical regression analysis technique called Partial Least Squares (PLS). Our PLS analysis extracts a low dimensional subspace within which we can use a simple quadratic discriminant classifier to classify image patches into vehicles and background.

Using such a large number of features greatly increases the computational cost of the vehicle detector. A common approach to speed up detectors using a large number of features is to use a boosting algorithm along with a rejection cascade as in [28]. We reduce the computational cost of our system using a dual feature selection approach. First, a recently proposed feature selection method called Ordered Predictors Selection, which combines a number of informative vectors that rank features based on their predictive performance, is used. This is coupled with multi-stage, multi-resolution image analysis, where a large number of image windows are discarded at an early stage of processing (processing at lower resolutions) and only a small fraction of image patches are analyzed at the highest image resolution (second stage). Our feature selection approach not only increases the speed of our system but also its performance.

We demonstrate our proposed vehicle detector on two datasets. The first consists of color images collected from a satellite and obtained via Google Earth. This is a set of 40 high resolution images over the city of San Francisco. The second dataset is the publicly available Overhead Imagery Research Data Set which consists of a large number of aerial images, annotated with vehicles [26]. We compare our approach to several previously proposed object detection approaches including the Histograms of Oriented Gradients approach of Dalal et al. [6], the Spatial Pyramidal Matching algorithm [16] incorporating SIFT features, the recently proposed Intersection Kernel Support Vector Machines

using HOG features [19] and a vehicle detector proposed in [20]. Our solution obtains favorable results as compared to previous approaches.

II. FEATURE EXTRACTION

We use three classes of features in our proposed solution: Color Probability Maps, Histograms of Oriented Gradients and Pairs-of-Pixels (PoP) features.

A. Color Probability Maps

Vehicles often lie on homogeneously colored backgrounds such as asphalt, cement, dirt roads, etc. Sometimes, the immediate neighborhood of a vehicle might be composed of multiple surfaces (a vehicle parked on the side of the road has an asphalt surface on three sides and a cement sidewalk on the fourth side). Thus, an image patch containing typical vehicle colors towards the center and colors representing typical backgrounds towards the periphery is likely to contain a vehicle. Color Probability Maps capture such color statistics of objects and their immediate environment.

We begin by identifying colors that are typically present in the background category. First, pixels are sampled from the entire set of background image patches in the training set of images. Each pixel is represented in a 3 dimensional space (r, g, s) , where r and g are chromaticity variables representing the fraction of the red and green components respectively: $r = \frac{R}{R+G+B}$ and $g = \frac{G}{R+G+B}$. s represents the brightness component: $s = \frac{R+G+B}{3}$. These pixels (obtained from all the background training images) are clustered to determine the dominant colors present in the background. Each cluster of pixels is used to build a color density model in RGB space using Kernel Density Estimation.

$$\begin{aligned} p^c(r, g, s) &= \frac{1}{N_{pts}^c} \sum_{i=1}^{N_{pts}^c} K_{\sigma_r}(r - r_i^c) K_{\sigma_g}(g - g_i^c) K_{\sigma_s}(s - s_i^c) \\ K_{\sigma}(x) &= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{x^2}{2\sigma^2}} \end{aligned} \quad (1)$$

where $p^c(r, g, s)$ refers to the probability of the (r, g, s) triple for the c^{th} cluster. σ_r , σ_g and σ_s represent the channel bandwidths. r_i^c, g_i^c and s_i^c refer to the chromaticity and brightness components of the i^{th} pixel in the c^{th} cluster. N_{pts}^c refers to the number of points in the c^{th} cluster. Given an image, one can obtain a color probability map for every color model. These probability maps are concatenated to form a feature vector representing the color statistics of the image patch, F_{cmap} . In order to limit the number of probability maps that must be computed for each image patch, only the most discriminating clusters are used. First, clusters that contain a very small number of points are rejected. Then the remaining clusters are ranked based on their discriminative capability. For a given cluster, we generate the corresponding probability map for all positive and negative samples in the training set and calculate the average misclassification error using a 5-fold cross validation procedure. The top N_{cmap} clusters are then chosen.

The Improved Fast Gauss Transform (IFGT) [22] is used for efficient computation of the probabilities. The bandwidths for each channel are estimated independently using a bandwidth selection criterion given in [23]. Given a test image patch, the computation of its Color Probability Maps is further speeded up

by *a priori* constructing look-up tables that directly map entries in the $R - G - B$ color space to a probability value. A look-up table is constructed for each of the N_{cmap} color clusters.

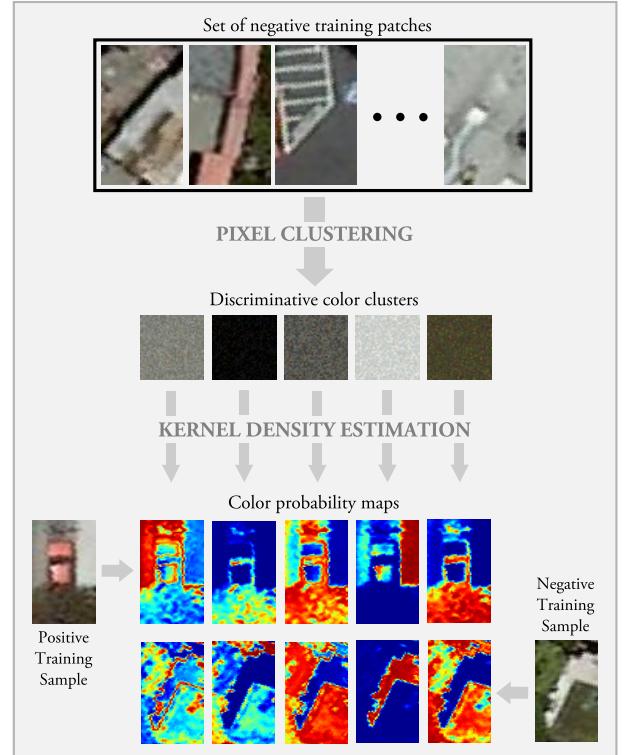


Fig. 2. **Color Probability Maps.** Pixels extracted from negative training image patches are clustered to obtain models of typical colors observed in the background. Given a new image patch, kernel density estimation is then used to obtain a probability map corresponding to each color cluster.

B. Histograms of Oriented Gradients

The second class of features are the Histograms of Oriented Gradients (HOG), which have been used in many object detection algorithms [6]. These features capture the spatial distribution of gradients that are typically observed in image patches that contain vehicles. Since histograms are computed over regions, they are fairly robust to some variability in the location of the parts of the object. Moreover, the HOG descriptor is also invariant to rotations smaller than the size of the histogram orientation bin.

Each detection window is divided into square cells and a 9-bin HOG feature is calculated for each cell. Grids of 2×2 cells are grouped into a block, resulting in a 36 dimensional feature vector per block. Each block feature vector is normalized to an L2 unit length. Dalal et al. [6] used blocks defined at a single scale. In their approach, for a detection window of size 64×128 pixels, 105 blocks were used, each having a size of 16×16 pixels. Zhu et al. [35] extended this approach by using blocks at varying scales and varying aspect ratios (1:1, 1:2 and 2:1). We incorporate this multi-scale approach employed in [35].

C. Pairs Of Pixels

Properties of pixel pairs within an image patch can provide structural information about the object present in that patch. Consider the image patch shown in Figure 3. The structure of a car shown in the image can be described using the relationships

between the regions highlighted in red, green and blue. The three regions highlighted in red represent the body of the vehicle and typically have the same color. Similarly, regions highlighted in green represent the windows of the vehicle which usually appear dark in color. Regions that are not present on the vehicle might have colors that differ from the color of the vehicle but they might be similar to one another. Such relative color statistics can capture the structure and relationships amongst different parts of a given object. Using relative properties of pixel pairs, as opposed to pixel properties themselves, is robust to illumination changes in the scene, changes in the background, as well as the color of the object itself.

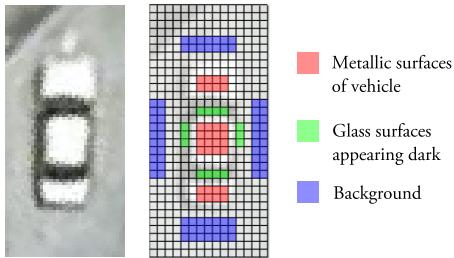


Fig. 3. The structure of a typical car and its surrounding regions can be described using pairwise relationships between the highlighted regions. Regions that are marked with the same color typically have the same color and texture properties. This information is captured by the PoP feature.

In principle, one can encode many different relative properties of regions, such as the difference between their colors, textural properties, gradient magnitudes, etc. Here, we restrict these regions to be single pixel locations and the relative property to be the Euclidean distance between their color values. The feature vector, F_{pop} , encoding this relative property can be obtained by concatenating the distances between the colors of a large number of pixel pairs.

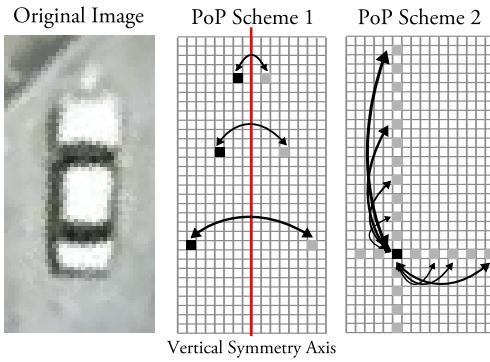


Fig. 4. POP Schemes. Scheme 1 captures differences between pairs of pixels that lie symmetrically across the central vertical axis of the image patch. Scheme 2 captures differences between pairs of pixels that lie in the same row and column.

If we were to consider all pixel pairs in an image patch, the feature vector would be of length M^2N^2 where M and N are the length and width of an image patch. For a small image patch of size 100×100 , this would result in a feature vector of length 100 million. In order to restrict this dimensionality, we propose two alternatives. These alternatives are designed to capture a large portion of the structural information, while taking advantage of the symmetry exhibited by vehicles about a central vertical axis.

Scheme 1: Consider all pairs of pixels that exhibit a symmetry in location about the central vertical axis, as shown in Figure 4. This reduces the dimensionality of the feature vector to $M \times \frac{N}{2}$. However, these pixel pairs only capture horizontal differences.

Scheme 2: Consider a pixel p at location (x, y) . Consider its differences with all pixels that lie in the same row and column as pixel p at intervals of distance d . This provides an advantage over Scheme 1 by capturing structural properties in both the horizontal and vertical directions while restricting the dimensionality of the resulting feature vector. Using Scheme 2 results in a feature vector of length $M \times N \times (1 - \frac{M+N}{d})$.

We use Scheme 2, since we are able to accommodate the length of the resulting PoP feature vector.

The three classes of features are finally combined to form the resulting feature vector describing an image patch.

$$F = [F_{cmap} \ F_{hog} \ F_{pop}] \quad (2)$$

III. PARTIAL LEAST SQUARES

The combination of the three feature classes results in an extremely high dimensional feature space (approximately 70,000 dimensions). In contrast, the number of samples in our training dataset is much smaller (about 200 in the positive and 1500 in the negative class). Furthermore, our features are extracted from neighboring pixels within a detection window, which tremendously increases the correlation between them, rendering traditional Ordinary Least Squares (OLS) regression estimates unreliable. This phenomenon is also known as the multicollinearity of the feature set. The nature of our proposed feature set makes an ideal setting for a statistical technique known as Partial Least Squares (PLS) regression [30].

The PLS method was first developed by Herman Wold in the 1960s and 1970s to address problems in econometric path-modeling [29], and was subsequently adapted in the 1980s to problems in chemometric and spectrometric modeling. In the late 1980's and 1990's, PLS attracted the attention of statisticians [10][14], due to its ability to deal with a small number of examples and a large number of variables.

We present a brief introduction to PLS. For a more detailed analysis, see [3]. Consider a set of p predictor variables, X_1, X_2, \dots, X_p , which are used to predict q response variables, Y_1, Y_2, \dots, Y_q . Let n equal the number of observation pairs denoted as (x_i, y_i) where $\{i = 1, 2, \dots, n\}$. The data samples are assumed to be mean-centered. They are concatenated to form the matrices $X_{(n \times p)}$ and $Y_{(n \times q)}$. When $n < p$, classical regression tools cannot be applied since the covariance matrix $X^T X_{(p \times p)}$ is singular.

PLS regression is based on the following latent component decomposition

$$X = TP^T + E \quad (3)$$

$$Y = UQ^T + F \quad (4)$$

where, T and U give the latent components (known as the *scores* matrices), P and Q provide the coefficients (known as the *loadings* matrices), and E and F are the error matrices. Note that a decomposition similar to Equation 3 is obtained by Principal Components Analysis.

The latent components given by T are obtained by a linear transformation of X as follows,

$$T_{n \times d} = X_{n \times p} W_{p \times d} \quad (5)$$

where d is the dimensionality of the latent space. The latent components T , are used for prediction in place of the original data vectors X . There are many variants of the basic PLS algorithm. They can be broadly classified based on their ability to deal with univariate response variables versus multivariate response variables. Multivariate response PLS has two popular implementations. The first variant leads to the NIPALS algorithm, whereas the second variant leads to the SIMPLS algorithm. These two methods differ in the matrix deflation process within the PLS algorithm. In our analysis, we have used the NIPALS algorithm. The NIPALS algorithm is essentially one of many methods that exist for finding the eigenvectors of a matrix. It was originally developed for Principal Components Analysis, but was subsequently used to iteratively extract factors for Partial Least Squares. Algorithm 1 provides a brief outline of the NIPALS algorithm. For more details we refer the reader to [11].

Algorithm 1 NIPALS Algorithm

```

1: for  $i = 1$  to  $d$  do
2:   Matrix Projection
    $W_i = X^T Y / \|X^T Y\|$ 
    $T_i = X W_i / \|X W_i\|$ 
3:   Matrix Deflation
    $X = X - T_i T_i^T X$ 
    $Y = Y - T_i T_i^T Y$ 
4: end for
```

W_i and T_i represent the i^{th} columns of the matrices W and T respectively. The regression model is given by,

$$Y = XB + F \quad (6)$$

where, $B_{p \times q}$ is the matrix of regression coefficients. Algebraic manipulations yield,

$$B = WQ^T = W(T^T T)^{-1} T^T Y \quad (7)$$

A new observation x_{new} thus yields a response value given by,

$$y_{new} = \frac{1}{n} \sum_{i=1}^n y_i + B^T (x_{new} - \frac{1}{n} \sum_{i=1}^n x_i) \quad (8)$$

The data we are interested in, falls into two classes - vehicles and background. We use the PLS regression algorithm as a *class aware* dimensionality reduction tool by setting the class label of a sample in Y to 1 or -1 . Thus, for our purpose, $q = 1$. Note that the matrix of regression coefficients B , is now a single vector ($B_{p \times 1}$), with a single coefficient for every feature. In practice, we do not project a new observation onto B . Instead, we project it onto the first k columns of matrix W , and then apply a classifier on that subspace. This method allows us to apply any classifier within this subspace (linear or non-linear) and has been shown to provide improved performance. The number of PLS factors k is obtained using cross validation.

Dimensionality reduction techniques can be broadly classified in two ways: linear vs non-linear methods and supervised vs unsupervised methods. Non-linear methods are generally computationally intensive and are not very suitable for extremely

high dimensional feature spaces, such as ours. For the purposes of classification, supervised methods hold an advantage over unsupervised methods, owing to their use of the class information within the training dataset.

Principal Component Analysis (PCA) is a classical linear unsupervised method. While PCA creates orthogonal latent vectors by maximizing the covariance between the data vectors x_i , PLS (a supervised approach) also considers the class labels. Figure 5 demonstrates the advantage of PLS over PCA. A subset of the training dataset (80% of the data points with all 69,552 variables) is provided to both PCA and PLS, and then projected onto the first two factors given by each dimensionality reduction method. The first row shows the training points projected onto the subspaces. The second row shows the test points (the remaining 20% of the data) projected onto the subspaces. The first two factors given by PLS are clearly more discriminating than PCA.

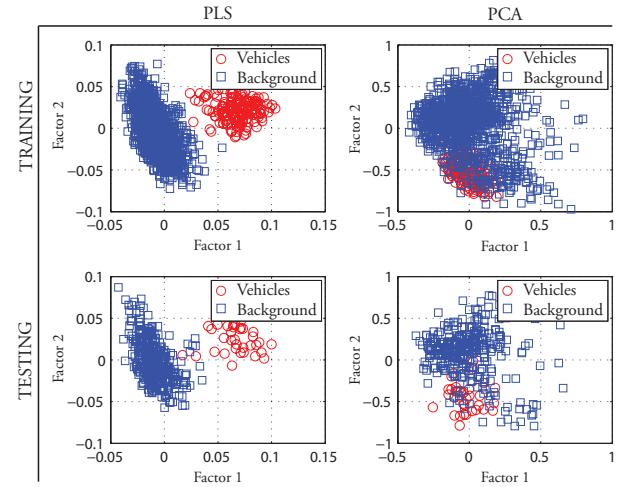


Fig. 5. Projection of data points from a 69,552 dimensional feature space onto a 2 dimensional subspace. In this illustration 80% of the training dataset is used to obtain the subspaces, and the remaining 20% of the data from each class are used as the testing samples. The left column shows the subspace extracted using PLS. The right column shows the subspace extracted using PCA. Clearly, PLS extracts a subspace that is more discriminating than PCA.

Fisher Discriminant Analysis (FDA) is in this way similar to PLS. It is a linear supervised method. However, FDA suffers from the *small sample size* problem. When the number of features exceeds the number of samples ($n < p$), as in our case, the covariance estimates are not full rank, which are required to obtain the projection vectors. A number of extensions to LDA have been proposed to deal with this problem. Belhumeur et al. [1] first projected points onto a lower dimensional subspace using PCA (which does not suffer from the small sample size problem), and then applied FDA on the reduced subspace. Chen et al. [4] used a modified version of Fisher's criterion and proposed an efficient and stable algorithm to calculate the discriminant subspace. However, FDA has a further limitation, in that it retains only $l - 1$ meaningful latent vectors, where l is the number of classes being considered. For our 2 class problem, a 1 dimensional subspace might not be sufficiently discriminative.

A. Visualization of PLS factors

It is useful to be able to visualize data points in the reduced PLS latent space. We propose a technique to visualize the separation

between the data points in the two classes as the dimensionality of the PLS latent space increases.

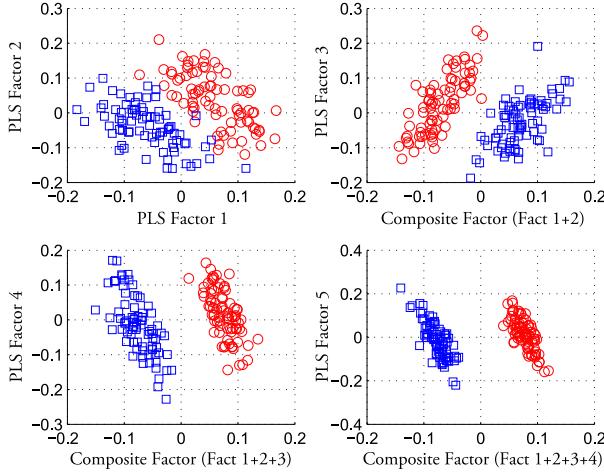


Fig. 6. Visualization of multiple PLS factors. The effect of the n^{th} PLS factor can be visualized by plotting it against the composite factor obtained by combining the first $n - 1$ PLS factors.

Figure 6 shows the visualization of data in a 5 dimensional PLS factor space. The top left plot shows the data plotted in the space spanned by the first two PLS factors. One can observe that the two classes are not completely separated in this 2 dimensional space. In order to visualize the effect of the third PLS factor, the data points are plotted in the space spanned by the third factor and the composite PLS dimension obtained from the first two PLS factors. The composite factor is essentially the vector of regression coefficients (B) given in equation 7. The data projected on B represents the output of PLS regression obtained using a 2 dimensional latent space. The top right plot shows the data projected on B as against the third PLS factor. In this way, data projected on factor i is plotted against the data projected on the composite factor obtained from the PLS factors 1 to $i - 1$.

One can observe the increase in the separation between the two classes as the number of PLS factors is increased. This visualization is a useful tool to observe the structure of the data in the two classes in a high dimensional PLS factor space.

Note that it cannot be used to determine the optimum number of PLS factors for a given problem, which is obtained by a cross validation procedure on the training dataset.

IV. FEATURE SELECTION

The large number of features greatly increases the computational cost of the vehicle detector. Calculating the features requires a substantial amount of time. At the same time, projecting a data point from an p dimensional space down to a d dimensional PLS factor space requires $p \times d$ multiplications and $(p-1) \times d$ additions. A typical image which must be scanned at multiple orientations and scales has hundreds of thousands of windows to be evaluated. Clearly, this would be a very time consuming process. This can be greatly speeded up, without a significant loss in performance, by using an effective feature selection process.

Furthermore, feature selection can often improve the performance of a PLS based classification system. In our set of thousands of features, one can expect many of them to be very

noisy and redundant. Feature selection can help discard a large fraction of such variables.

We use two methods to perform feature selection: Ordered Predictors Selection and a Multi-Stage Multi-Resolution Analysis.

A. Ordered Predictors Selection (OPS)

Variable Importance on Projection (VIP) is a widely used technique for PLS based feature selection. VIP provides a score for every variable, that ranks them according to their predictive power. A cross validation scheme can then be used to select the number of variables required to obtain a desired level of classification accuracy. In general any informative vector which provides a measure of the predictive power of the variables, can be used for feature selection.

Teofilo et al. [27] used several informative vectors and their combinations to perform feature selection for regression problems. Their method is computationally efficient as compared to other variable selection methods, such as genetic algorithms, and can be completely automated. They used several different datasets in their analysis, some having a multiple number of dependent variables. Their analysis showed a rather surprising result. The number of PLS factors that were obtained by a cross validation procedure for the purpose of feature selection was often higher than the number of PLS factors that were optimal for the purpose of regression.

We perform a similar feature selection analysis, enabling us to reject a large fraction of noisy features. We only deal with a single dependent variable, which is set to the class label (+1/ -1). The following informative vectors are used in our analysis.

- 1) Variable Importance on Projection (VIP) - The VIP score for the j^{th} variable is a measure based on the weighted PLS coefficients. The higher the score, the more importance a variable presents. The average VIP score over all variables equals 1. The *thumb rule* used to select variables according to their VIP score is to retain only those variables whose VIP score is greater than 1.

$$\text{VIP}_j = \sqrt{p \sum_{k=1}^d B_k^2 W_{jk}^2 / \sum_{k=1}^d B_k^2} \quad (9)$$

- 2) Regression Coefficients (B) - The regression coefficients B defined in Equation 7 represent the expected change in the response, per unit change in the variable. The absolute value of the regression coefficients are thus used as informative scores.
- 3) Correlation (CORR) - The correlation informative vector contains the Pearson correlation coefficients between every predictor variable X_i and the response variable Y . A high correlation indicates that the predictor variable is very informative about the PLS classification model.

$$R = \frac{X^t Y}{n-1} \quad (10)$$

Since the Pearson correlation coefficient lies between -1 and $+1$, with 0 indicating an absence of any correlation, the absolute value of the correlation coefficients is used.

- 4) Covariance Procedures Vector (CVP) - The CVP score was proposed by Reinikainen et al. [24] as a measure of variable importance. The ranking of variables is based on

the covariance of the dependent and independent variables which is given by,

$$\text{CVP} = \text{diag}(X^T Y Y X) \quad (11)$$

Informative Vectors CORR and CVP do not depend on the dimensionality of the PLS factor space. However, VIP and BETA vary with the dimensionality of the factor space. Motivated by the OPS results obtained in [27], a nested cross validation procedure is used to determine the optimum dimensionality of the PLS factor space for classification and the optimum dimensionality of the PLS factor space to calculate the informative vector.

Algorithm 2 OPS Cross Validation Procedure

```

1: for  $\theta = 1$  to  $N_{\text{select}}$  do
2:   Build PLS model using  $\theta$  factors
3:   Calculate informative vector and sort variables by their
   informative score
4:   Choose top  $K$  variables based on feature selection criteria
5:   for  $\phi = 1$  to  $N_{\text{model}}$  do
6:     Build PLS model using  $\phi$  factors and top  $K$  variables
7:     Calculate classification accuracy on the validation set
8:   end for
9: end for
```

B. Multi-Stage Multi-Resolution Analysis

The nature of the three feature classes introduces a fair amount of redundancy in the feature set. For instance, the color probability maps capture color statistics of every pixel within an image patch. Clearly, neighboring pixels are highly correlated and introduce a lot of redundancy. Downsampling an image can help remove this redundancy somewhat, at the cost of performance.

We build 2 PLS models using features computed from training images at different resolutions. The first model is built from training images reduced to a width and height that equal 1/2 the original image dimensions, vastly reducing the dimensionality of the feature vector. The second model is computed from the original training images. The features that contribute towards each PLS model, undergo the OPS feature selection strategy that was outlined in Section IV-A. Thus we obtain two trimmed PLS models. In principle, one may add more downsampling stages to obtain a further speedup, possibly at the cost of accuracy.

Given a testing image, each and every image patch is classified using the first and fastest PLS model. A subset of these are then sent to the second stage which contains the second PLS model (lowest speed, highest performance).

V. EXPERIMENTS

A. Google Earth San Francisco Dataset

We test the performance of our vehicle detector on satellite images of the city of San Francisco taken from Google Earth. This dataset consists of 40 satellite images at a resolution of 979×1348 pixels and a color depth of 24 bits per pixel (RGB). Figure 7 shows one such image. Each image looks down on an urban scene with multiple cars present. The total number of vehicles in the dataset is 650. The average size of a vehicle is 48×16 pixels¹.

¹The highest resolution for current commercial satellite imagery is 0.5 meters, where as the images used in the San Francisco dataset correspond to a resolution of 0.1 meters. We have used images captured using the software: Google Earth. Images provided by Google Earth beyond this resolution are obtained by digitally enlarging the images.

The first 5 images in the dataset are used for training and the performance of our vehicle detector is tested on the remaining 35 images. 184 positive image patches (vehicles) are extracted from the training images and aligned vertically. Similarly 1500 negative image patches are randomly chosen from the training images. Each image patch has a size of 81×41 pixels. Figure 7 shows a few training image patches from the dataset. The test images are scanned using a sliding window approach. The size of the window is fixed to 81×41 pixels, since the training and test images have been captured at the same resolution and all vehicles are observed at a single scale. More generally however, one may need to scan the image at multiple scales. The horizontal and vertical step sizes are set to 5 pixels.

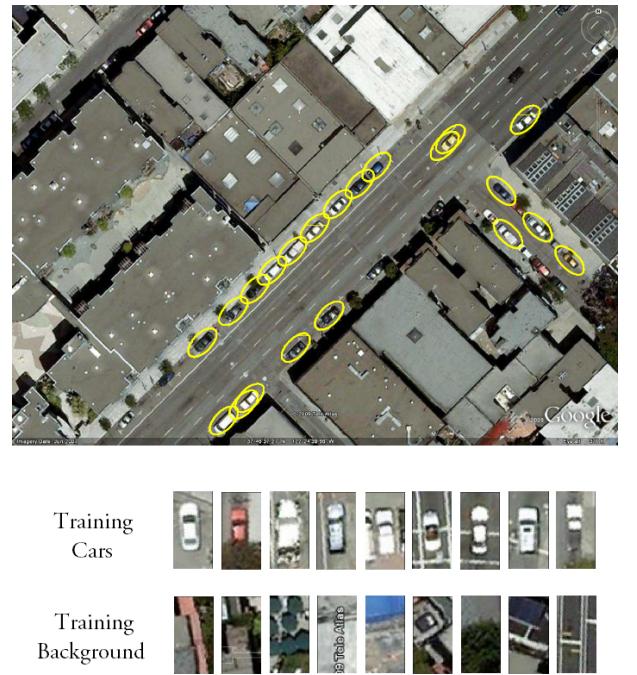


Fig. 7. Example images from the Google Earth San Francisco dataset. Top image shows a test image from the dataset, looking down on an urban scene, with overlaid detections using our vehicle detector. Bottom image shows a few training images patches from both classes. © 2009 Google

B. Feature Extraction and Evaluation

Each image patch has a size of 81×41 pixels. 6 color clusters are used to build our color probability maps. This results in a feature vector of length 19,926. The HOG feature is calculated for 661 blocks ranging from a size of 12×12 pixels to a size of 80×40 pixels. Each block yields a 36 dimensional feature vector resulting in a total length of 23,796. The PoP feature is calculated using images reduced to size 41×21 , which gives a feature of length 25,830. Thus, the length of the resulting feature vector obtained by combining the three feature classes is 69,552.

All the parameter selection in our system is performed by using 3 iterations of a 5-fold cross validation scheme. The analysis presented in Figures 8 to 15 was carried out using this cross validation procedure on the Google Earth San Francisco training dataset.

Figure 8 shows the mean classification error obtained using each of the three feature classes separately as well as their combination. As the number of PLS factors is increased, the classification error reduces. Beyond a certain value however,

addition of PLS factors does not yield an improved performance. This saturation point is generally regarded as the optimum number of PLS factors. In some cases, performance decreases as more PLS factors are introduced. The PoP feature class outperforms the other two feature classes, but the best performance is obtained when all three feature classes are combined. The number of factors chosen, when the entire set of features is used, is 5.

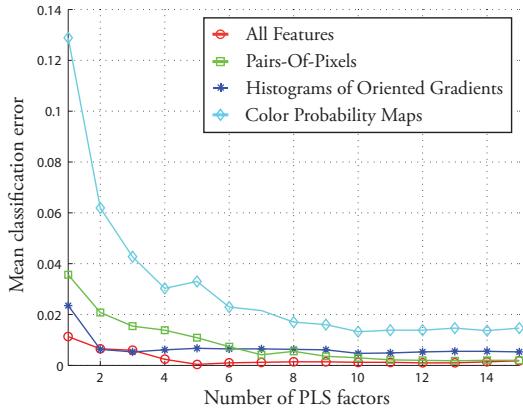


Fig. 8. Mean classification error obtained after 3 iterations of 5-fold cross validation on the Google Earth San Francisco training dataset. The error plot is shown for each class of features individually, as well as their combination.

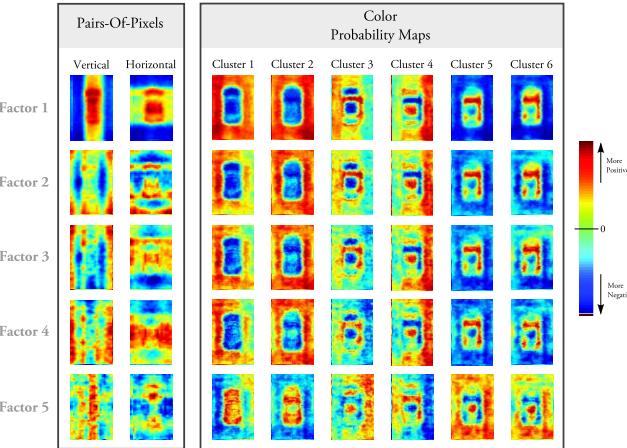


Fig. 9. The first 5 PLS latent vectors. The components of the latent vectors corresponding to the Color Probability Maps and Pairs of Pixels feature classes are displayed. The images shown for the color probability maps directly correspond to the latent vectors. The images displayed for the PoP feature correspond to accumulator matrices. (*best viewed in color*)

Figure 9 shows the components of the first 5 PLS latent vectors (W' 's) that correspond to the color probability maps and the PoP features. The latent vectors corresponding to the color probability maps are directly displayed as shown in the figure, since each coefficient corresponds to a single pixel location in the image patch. But each coefficient of the latent vector from the PoP features corresponds to two pixel locations in the image (a pair of pixels), and thus cannot be displayed directly. Instead, the images displayed in the figure for the PoP feature class correspond to accumulator matrices. We initialize the accumulator matrix as a matrix of zeros and increment a location with the corresponding coefficient in the latent vector, when that location forms one half of the corresponding pixel pair. Furthermore, the images corresponding to the PoP feature class are split into the

vertical and horizontal components. Each coefficient in the HOG feature vector corresponds to a bin of a histogram and captures information about a neighborhood of pixels. These neighborhoods also have varying sizes (multiple block sizes), and thus cannot be easily visualized as the other two feature classes.

A very positive (dark red) or very negative (dark blue) color indicates a high degree of importance afforded to that pixel location. It is noticeable that the feature classes complement each other well, by extracting information from different parts of the image patch. The first (and most important) factor for the color probability maps shows that information extracted from the neighborhood of the object is given a high weight. This indicates that these features are able to capture the immediate context within which vehicles are typically observed. The PoP feature is seen to primarily focus on pixels that lie on top of the vehicle, in order to capture its regular structure.

C. Feature Selection: OPS

Figure 10 shows the result of the commonly used feature selection criterion using VIP. The VIP informative vector has been calculated using 5 PLS factors (which was the optimal number when using the entire set of features - refer to Figure 8). Using the $VIP > 1$ thumb rule reduces the number of features to 21,322 which represents about 30% of the total feature set. The performance does not drop significantly; instead a comparable performance is obtained using a smaller number of PLS factors. Results using varying values of the VIP cut-off score are shown.

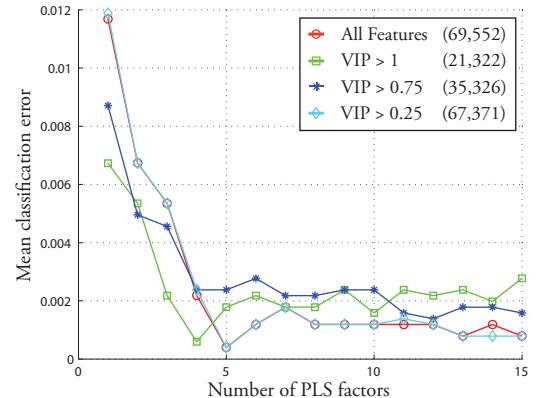


Fig. 10. Basic feature selection using VIP. The error plot is shown for the original set of features and for a subset of features using the VIP criterion. $VIP > 1$ is the *thumb rule* usually applied in PLS analysis.

Figure 11 shows the results of feature selection using the four informative vectors. The VIP and B informative vectors have been calculated using 5 PLS factors. The Correlation and Covariance Procedures Vector are calculated independent of the number of PLS factors. As a fair comparison, we use the same number of features with each informative vector. This number, 21,322, is the number of features retained using the $VIP > 1$ rule. The best performing informative vector is the regression vector B .

Figure 12 shows the results of feature selection using the B informative vector when the number of PLS factors (θ) used to compute this vector are varied. As was observed in [27], the optimum number of PLS factors used to calculate the informative vector (θ) is not the same as the optimum number of PLS factors used to build the model (ϕ), when all features are used (ϕ was

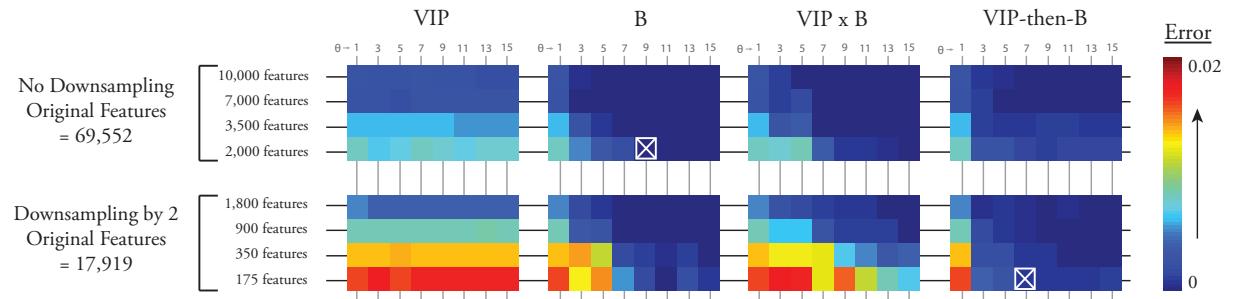


Fig. 13. Results of the OPS feature selection approach on the original training set and the downsampled set. The errors represent the mean misclassification error after cross validation. (Refer to the text for more details).

only 175 features in the fast first stage (almost 0.25% of the original number of features 69,552). This enables a fast rejection strategy and improves the performance of the system. Figure 15 shows the performance at a downsampling factor of 2 when the popular $VIP > 1$ feature selection strategy is used, as well as when the best feature selection criteria is used.

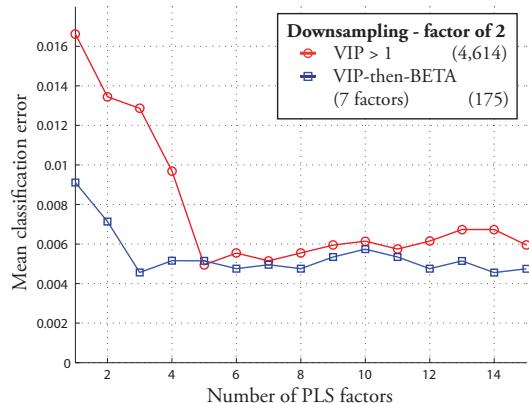


Fig. 15. Mean misclassification error at downsampling factor of 2. The error rates drop after the feature selection process. The informative vector combination VIP -then-B outperforms $VIP > 1$.

Thus we obtain a 2 stage PLS model, whose parameters are given in Table I.

E. Performance: Google Earth San Francisco dataset

We test the performance of our system on the test images of the Google Earth San Francisco dataset. All test images are fully ground truthed to show the presence of vehicles along with their orientation. Vehicle detections that overlap the ground truth locations by an area equal to 33% of the size of the bounding box, are considered true detections. As per the evaluation criteria commonly used in the PASCAL VOC challenge [8], if multiple detections overlap with a ground truthed location, only one of them is considered a true positive detection. The remaining detections are considered to be false alarms. Since the vehicles in the entire dataset are roughly the same size, we only scan the images at a single scale. If the size of the vehicles was unknown, the images would have to be scanned at multiple scales. Since the orientation of the vehicles is unknown, the image must be scanned at multiple rotations. In order to decrease the number of image patches that must be scanned, we employ a coarse to fine rotation strategy. First the image is rotated in increments of 30° and scanned completely. Candidate windows are selected and are

then finely rotated in increments of 5° in a neighborhood of 25° from the original detection². In an urban setting, roads typically form a rectangular grid. Determining this grid orientation can help to initialize the scan angles and further speed up the system.

We compare our system to a number of other approaches. The first approach is a traditional object detection approach³. SIFT features are calculated on a dense grid of points in the training images. The SIFT descriptors from the training set undergo vector quantization to form visual words. An image patch can then be described using a histogram of the visual words present in the patch. In order to enforce some spatial constraints on the location of these features, we use Spatial Pyramidal Matching [16]. This involves repeatedly subdividing the image and computing histograms of SIFT features at increasingly finer resolutions. The distance measure used is histogram intersection and the classifier used is the Support Vector Machine.

The second approach we compare to, is the vehicle detector proposed by Moon et al. [20]. They derive an optimal one dimensional step edge detector to be the derivative of the double exponential (DODE) function. This is extended along the shape's boundary contour to obtain the shape detector⁴. This results in detecting shapes in the image that resemble parallelograms. Their vehicle detector does not require a training phase.

The third approach we compare to is the popularly used HOG based approach used to detect objects, proposed by Dalal et al. [6]. HOG features are calculated for a large number of blocks within an image window and concatenated together to form the feature vector. Blocks of only a single size are used in this approach. The feature vectors thus obtained are used to train a linear Support Vector Machine (SVM). We used libLinear [9]⁵, an efficient linear SVM, for this purpose. Dalal et al. note in their work that a kernel SVM can provide improved performance at the cost of a significant reduction in the efficiency of the system. This is because the standard approach to evaluate the kernel for a test vector involves a comparison with each of the support vectors.

Recently, Maji et al. [19] proposed a method to significantly improve the efficiency of kernel SVMs for a class of kernels such as the histogram intersection kernel and the chi squared kernel. Their approximate SVM has constant runtime and space requirements, independent of the number of support vectors, as opposed to the typical linear dependency. Furthermore, the loss in classification accuracy using this approximation is negligible. As a

²The identical scanning strategy is used for the proposed vehicle detector and the approaches we compare to.

³Code obtained from <http://www.cs.unc.edu/~lazebnik/>

⁴Code obtained from the authors

⁵Package available at <http://www.csie.ntu.edu.tw/~cjlin/liblinear/>

TABLE I
THE 2 STAGE VEHICLE DETECTOR

	STAGE 1	STAGE 2
Image size	41 × 21 pixels	81 × 41 pixels
Original features	17,919	69,552
Feature selection	VIP-then-B ($\theta = 7, \phi = 7$)	B ($\theta = 9, \phi = 6$)
Retained features	175	2,000
Percentage Windows Processed	100%	0.48%

fourth comparison, we applied this improved kernel SVM⁶ known as the approximate intersection kernel support vector machine (approx IKSVM) to the HOG features proposed by Dalal and Triggs. While [19] demonstrated that the IKSVM evaluates nearly as fast as a linear classifier, it did not address the problem of efficiently training such a classifier. Subsequently Maji et al. [18] proposed very efficient training algorithms for additive classifiers in a max-margin framework.

Figure 16 shows the Precision-Recall curves for the vehicle detectors. The dataset we test on, contains images in an urban setting. The presence of a large number of rectilinear structures in such images leads to false alarms with all approaches. Objects present on top of buildings, such as air conditioning units are seen to cause errors. The DODE approach is able to correctly find many vehicles but also produces a very large number of false alarms on rectangular structures. The HOG features when applied to a linear SVM outperform the SIFT based approach. The use of the histogram intersection kernel greatly improves performance over the linear kernel. Our proposed solution outperforms all the other approaches.

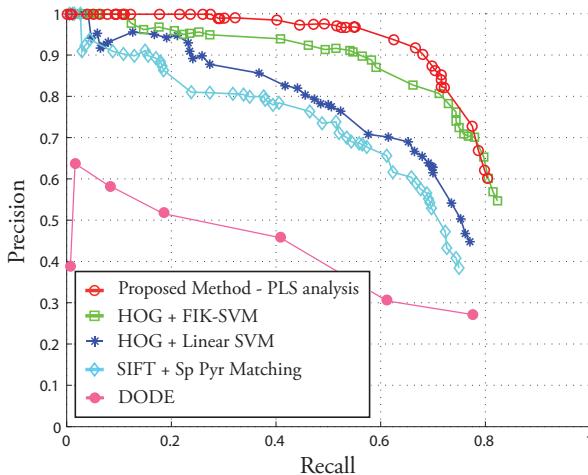


Fig. 16. Performance of the five vehicle detectors on the Google Earth San Francisco dataset. Our vehicle detector outperforms the other ones.

Figure 17 compares our proposed approach to directly applying an SVM to our proposed set of 69,552 features. The PLS based approach comfortably outperforms the approach using a linear SVM. We used the LibLinear package for this purpose. Using the approximate and fast IKSVM method improves results over using a linear kernel. However, applying the IKSVM to the original 69,552 features is quite time-consuming. Table II compares the speeds of the 3 classification approaches shown in Figure 17. Note that the numbers provided in Table II account for the

⁶Code obtained from <http://www.cs.berkeley.edu/~smaji/projects/fiksvm/>

TABLE II
COMPARISON OF CLASSIFICATION SPEEDS (WINDOWS/SECOND) WHILE
USING 69,552 FEATURES

	PLS method	LibLinear	Fast IKSVM
Detection speeds (win/sec)	8565	308	95

classification time only, and not the time required to calculate the features. Projection onto the subspace followed by classification by a quadratic classifier is very fast compared to the other two approaches. Applying the fast approximate intersection kernel is slower than the other two methods.

Overall, our two stage vehicle detection system is able to process approximately 5600 detection windows per second. Our system is implemented in MATLAB and all our experiments are run on a Intel Xeon 2.8 GHz processor. While the machine we use has multiple processing cores, our program currently makes use of only a single core. The number of detection windows processed per second can be improved by taking advantage of multi-core architectures. All three stages of our detection system - feature calculation, projection onto a subspace as well as classification - can be parallelized to yield a faster detection system.

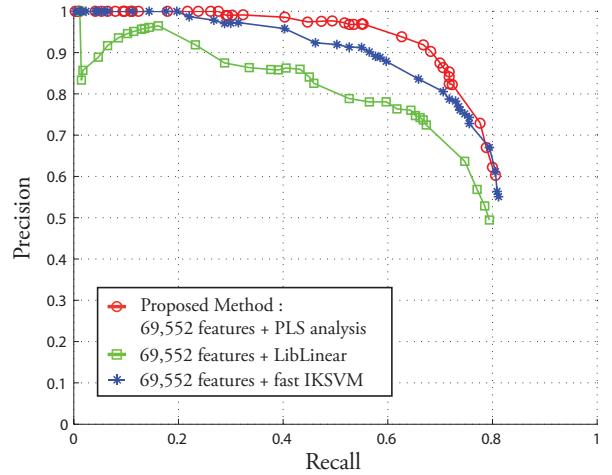


Fig. 17. Performance of the proposed vehicle detector compared to detectors composed of the proposed features and SVMs as classifiers.

F. Speedup using Feature Selection

Figure 19 compares the performance of the 2 stage vehicle detector with a detector when only the high resolution stage (Stage 2) is used. The Precision-Recall curves of the two detectors are quite comparable. This enables us to obtain the speedup given by the 2 stage approach, without a significant loss of performance.



Fig. 18. Sample vehicle detection results from the Google Earth San Francisco dataset using our proposed approach.

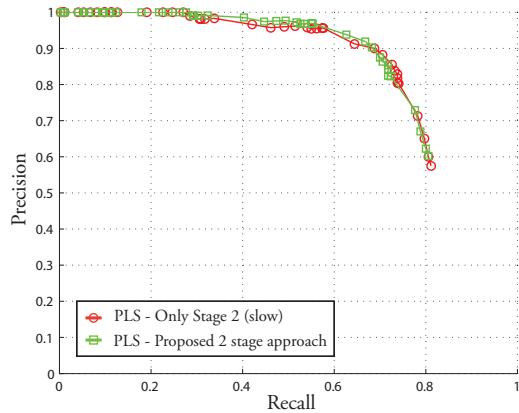


Fig. 19. Performance of the 2 stage vehicle detector compared to the performance of the detector when only the the high resolution stage (Stage 2) is used.

Table I shows the percentage of windows that are processed in each stage, over the entire Google Earth San Francisco dataset. All windows are processed by the fast Stage 1, but only a very small number are passed on to the second stage. Note that the image patches that obtain a high detection probability in stage 2, undergo further processing by rotating them at finer angular intervals. This processing is done at the full resolution.

G. Overhead Imagery Research Data Set

We also test the performance of our system on the publicly available Overhead Imagery Research Data Set (OIRDS)⁷. The

OIRDS is a large collection of almost 900 overhead images, captured using aircraft mounted cameras. The total number of vehicles annotated in the dataset is around 1800.

The OIRDS dataset is a very challenging dataset. The images in this set have varying levels of zoom and a wide degree of difficulty. The images have been captured primarily in suburban settings. The presence of a large number of trees in these images often causes vehicles to be partially occluded. We divide this dataset into 3 parts (with roughly 300 images in each part). We label these parts OIRDS 1, OIRDS 2 and OIRDS 3⁸. OIRDS 1 contains images that have the best picture quality and vehicles that are clearly visible. These images are similar in quality to the images in the Google Earth San Francisco dataset. OIRDS 2 contains images that are of a poorer quality and the vehicles are also harder to find. Some of these vehicles are partially occluded. OIRDS 3 contains images in which vehicles are very difficult to find. Many of these images have vehicles that were almost fully occluded. Figure 20 shows sample images from these three sets.

We compare the performance of the five vehicle detectors on OIRDS 1 and OIRDS 2. No new training was carried out for this dataset. The detectors trained on the Google Earth San Francisco training dataset were directly used. Figures 21 and 21 shows the Precision-Recall curves for all three detectors. We outperform all detectors on OIRDS1 and obtain a comparable performance to the HOG + kernel SVM approach on OIRDS 2. As expected, the performance of all the detectors drops for OIRDS 2.

Finally, Figure 23 shows the performance of our vehicle detector on a large panoramic image (5007 × 7776 pixels). This

⁷Downloaded from <http://sourceforge.net/apps/mediawiki/oirds>

⁸More details available at:
www.umiacs.umd.edu/~ani/vehicleDetection.html

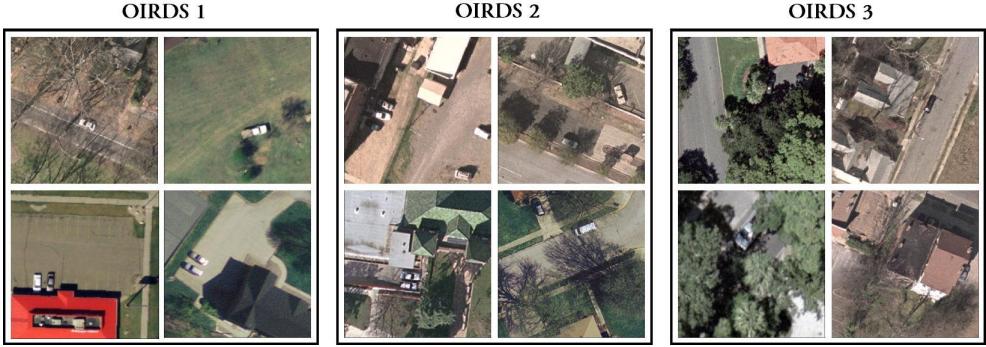


Fig. 20. Sample images taken from the OIRDS dataset. We divide this dataset into three parts based on the degree of difficulty.



Fig. 23. Performance of our vehicle detector on a large panoramic image overlooking the parking lot of the San Francisco Giants stadium. © 2009 Google

was obtained by stitching a large number of images obtained from Google Earth. The image overlooks the parking lot and adjacent areas of the San Francisco Giants stadium in San Francisco city. Our vehicle detector is able to accurately locate a large number of vehicles in this image.

H. Misclassifications

Figures 18 and 23 show the presence of a few missed detections as well as false alarms. The presence of stark and long shadows often cause corruption of the measured feature vectors and lead to missed detections. Similarly, saturated pixels on the body of vehicles, caused by the reflection of the sun, also give rise to missed detections. Some dark cars are also missed by the vehicle detector, due to the difficulty in capturing the structure present within the body of the vehicle (such as windows). Since the algorithm densely scans the image at multiple locations and orientations, we employ non-maximal suppression to only retain detections that are dominant within a given radius. However, this also leads to the removal of some true detections. False alarms

are typically caused by rectangular structures on top of buildings such as air-conditioning units. Markings on the road such as the ones to denote parking spots may also lead to some false alarms.

Vehicles typically follow a small set of patterns when they are stationary as well as moving. For example, they are typically present on roads or parking lots; parked cars usually align to lie parallel to each other; cars on the road and in the same lane usually move in the same direction. Using such contextual cues (as demonstrated in [13]) could further reduce the number of false detections in a given image, as well as improve the orientation estimation of the detected vehicles.

VI. CONCLUSION

We propose a vehicle detection system that incorporates a large set of rich features capturing color, gradient and structural properties of vehicles and their surroundings. A Partial Least Squares analysis enables us to project points from a very high dimensional feature space on to a low dimensional subspace. We experiment with a number of informative vectors which



Aniruddha Kembhavi Aniruddha Kembhavi received the bachelors degree in electronics and telecommunications engineering from the Government College of Engineering, Pune, India, in 2004. He is currently working toward the PhD degree in the Computer Vision Laboratory at the University of Maryland, College Park. His current research focuses on the problem of variable and feature selection for object classification. His research interests include human detection, object classification, and scene understanding. He is a member of the IEEE.



David Harwood David A. Harwood is a graduate of the University of Texas and M.I.T. His research, with many publications, is in the fields of computer image and video analysis and AI systems for computer vision. He is a longtime member of the research staff of the Computer Vision Laboratory of the Institute for Advanced Computer Studies.



Larry S. Davis Larry S. Davis received the BA degree from Colgate University in 1970 and the MS and PhD degrees in computer science from the University of Maryland in 1974 and 1976, respectively. From 1977 to 1981, he was an assistant professor in the Department of Computer Science at the University of Texas, Austin. He returned to the University of Maryland as an associate professor in 1981. From 1985 to 1994, he was the director of the University of Maryland Institute for Advanced Computer Studies, where he is currently a professor, and is also a professor and the chair of the Computer Science Department. He was named a fellow of the IEEE in 1997.