

A New Algorithm for Pedestrian Detection

Cheng Ke-yang^{1,2}, Bao Jun-xian²

¹ College of Information Science and Technology, Nanjing University of Aeronautics and Astronautics,
Nanjing Jiangsu, China

² College of Computer Science and Telecommunication Engineering, Jiangsu University, Zhenjiang Jiangsu,
China

Abstract. This article puts forward a novel framework for pedestrian detection tasks, which proposing a model with both sparse reconstruction and class discrimination components, jointly optimized during dictionary learning. We present an efficient pedestrian detection system using mixing sparse features of HOG, FOG and CSS to combine into a Kernel classifier. Results presented on our data set show competitive accuracy and robust performance of our system outperforms current state-of-the-art work

Keywords: Pedestrian detection, Sparse reconstruction, Class discrimination

1. Introduction

Pedestrian counting in public places plays a key role in many applications, such as evacuating from a dense region to a sparse one when an emergency happens, or optimizing the design of traffic infrastructures to provide better transportation services. Furthermore, social security and surveillance strongly depend on the effectiveness of pedestrian counting. A wide variety of pedestrian detection methods have been proposed [1-6].

Sparse representations have recently drawn much interest in signal, image, and video processing. Under the assumption that natural images admit a sparse decomposition in some redundant basis (or so-called dictionary), several such models have been proposed, e.g., curve lets, wedge lets, band lets and various sorts of wavelets [7]. Interestingly, while discrimination is the main goal of these papers, the optimization (dictionary design) is purely generative, based on a criterion which does not explicitly include the actual discrimination task, which is one of the key contributions of our work. In [8], a discriminative method is introduced for various classification tasks, learning one dictionary per class; the classification process itself is based on the corresponding reconstruction error, and does not exploit the actual decomposition coefficients. In [9], a generative model for documents is learned at the same time as the parameters of a deep network structure. In [10], multi-task learning is performed by learning features and tasks are selected using a sparsity criterion. The framework we present in this paper extends these approaches by learning simultaneously a single shared dictionary as well as models for different signal classes in a mixed generative and discriminative formulation (see also [11], where a different discriminative term is added to the classical reconstructive one). Similar joint generative/discriminative frameworks have started to appear in probabilistic approaches to learning, e.g., [12, 13, 14, 15, 16], and in neural networks [17], but not, to the best of our knowledge, in the sparse dictionary learning framework.

The remainder of this paper is organized as follows. In Section 2, we describe the procedure of feature extraction, and in Section 3, we present a formulation for learning a dictionary tuned for a classification task, which we call discriminative sparse learning. Experimental results are provided and analyzed in Section 4. Finally, Section 5 concludes this work.

2. Feature extraction

Obviously, the choice of features is the most critical decision when designing a detector, and finding good features is still largely an empirical process with few theoretical guidelines. We evaluate different combinations of features, and introduce a new feature based on the similarity of colors in different regions of the detector window, which significantly raises detection performance. The pedestrian region in our detection window is of size 48×96 pixels.

Histograms of oriented gradients (HOG) are a popular feature for object detection, first proposed in [18]. They collect gradient information in local cells into histograms using trilinear interpolation, and normalize overlapping blocks composed of neighboring cells. Interpolation, local normalization and histogram binning make the representation robust to changes in lighting conditions and small variations in pose. HOG was recently enriched by Local Binary Patterns (LBP), showing a visible improvement over standard HOG on the INRIA Person data set [24]. In our experiments we compute histograms with 9 bins on cells of 8×8 pixels. Block size is 2×2 cells overlapping by one cell size.

HOF Histograms of flow were initially also proposed by Dalal et al. [19]. We have shown that using them (e.g. in [19]’s IMHwd scheme) complementary to HOG can give substantial improvements on realistic datasets with significant ego motion. Here, we introduce a lower-dimensional variant of HOF, IMHd2, which encodes motion differences within 2×2 blocks with 4 histograms per block, while matching the performance of IMHwd (3×3 blocks with 9 histograms). Fig. 2(d) schematically illustrates the new coding scheme: the 4 squares display the encoding for one histogram each. For the first histogram, the optical flow corresponding to the pixel at the i th row and j th column of the upper left cell is subtracted from the one at the corresponding position of the lower left cell, and the resulting vector votes into a histogram as in the original HOF scheme. IMHd2 provides a dimensionality reduction of 44% (2520 instead of 4536 values per window), without changing performance significantly. We used the publicly available flow implementation of [20]. In this work we show that HOF continues to provide a substantial improvement even for flow fields computed on JPEG images with strong block artifacts (and hence degraded flow fields).

Several authors have reported improvements by combining multiple types of low-level features [21, 22, 23]. Still, it is largely unclear which cues should best be used in addition to the now established combination of gradients and optic flow. Intuitively, additional features should be complementary to the ones already used, capturing a different part of the image statistics. Color information is such a feature enjoying popularity in image classification [24] but is nevertheless rarely used in detection. Furthermore, second order image statistics, especially co-occurrence histograms, are gaining popularity, pushing feature spaces to extremely high dimensions [25, 22].

We propose to combine these ideas and use second order statistics of colors as additional feature. Color by itself is of limited use, because colors vary across the entire spectrum both for people (respectively their clothing) and for the background, and because of the essentially unsolved color constancy problem. However, people do exhibit some structure, in that colors are locally similar—for example (see Fig. 1) the skin color of a specific person is similar on their two arms and face, and the same is true for most people’s clothing. Therefore, we encode color self similarities within the descriptor window, i.e. similarities between colors in different sub-regions. To leverage the robustness of local histograms, we compute D local color histograms over 8×8 pixel blocks, using trilinear interpolation as in HOG to minimize aliasing. We experimented with different color spaces, including $3 \times 3 \times 3$ histograms in RGB, HSV, HLS and CIE Luv space, and 4×4 histograms in normalized rg, HS and uv, discarding the intensity and only keeping the chrominance. Among these, HSV worked best, and is used in the following.

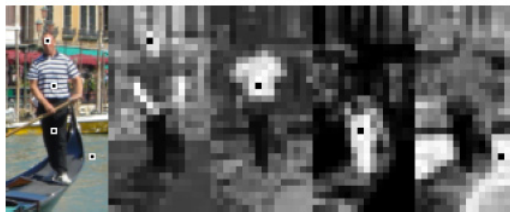


Fig.1. Self-similarity encodes relevant parts

3. Supervised dictionary learning

We present in this section the core of the proposed model. In classical sparse coding tasks, one considers a signal x in \mathbb{R}^n and a fixed dictionary $D = [d_1, \dots, d_k]$ in $\mathbb{R}^{n \times k}$ (allowing $k > n$, making the dictionary over complete). In this setting, sparse coding with an ℓ_1 regularization amounts to computing

$$R^*(x, D) = \min_{\alpha \in \mathbb{R}^k} \|x - D\alpha\|_2^2 + \lambda_1 \|\alpha\|_1 \quad (1)$$

It is well known in the statistics, optimization, and compressed sensing communities that the ℓ_1 penalty yields a sparse solution, very few non-zero coefficients in α , although there is no explicit analytic link between the value of λ_1 and the effective sparsity that this model yields. Other sparsity penalties using the ℓ_0 regularization can be used as well. Since it uses a proper norm, the ℓ_1 formulation of sparse coding is a convex problem, which makes the optimization tractable with algorithms such as those introduced in [26, 27], and has proven in practice to be more stable than its ℓ_0 counterpart, in the sense that the resulting decompositions are less sensitive to small perturbations of the input signal x . Note that sparse coding with an ℓ_0 penalty is an NP-hard problem and is often approximated using greedy algorithms.

In this paper, we consider a setting, where the signal may belong to any of p different classes. We first consider the case of $p = 2$ classes and later discuss the multiclass extension. We consider a training set of m labeled signals $(x_i)_{i=1}^m$ in \mathbb{R}^n , associated with binary labels $(y_i \in \{-1, +1\})_{i=1}^m$. Our goal is to learn jointly a single dictionary D adapted to the classification task and a function f which should be positive for any signal in class $+1$ and negative otherwise. We consider in this paper two different models to use the sparse code α for the classification task:

(i) linear in α : $f(x, \alpha, \theta) = w^T \alpha + b$, where $\theta = \{w \in \mathbb{R}^k, b \in \mathbb{R}\}$ parametrizes the model.

(ii) bilinear in x and α : $f(x, \alpha, \theta) = x^T w \alpha + b$, where $\theta = \{W \in \mathbb{R}^{n \times k}, b \in \mathbb{R}\}$. In this case, the model is bilinear and f acts on both x and its sparse code α .

The number of parameters in (ii) is greater than in (i), which allows for richer models. Note that one can interpret w as a linear filter encoding the input signal x into a model for the coefficients, which has a role similar to the encoder in [27] but for a discriminative task. A classical approach to obtain for (i) or (ii) is to first adapt D to the data, solving

$$\min_{D, \alpha} \sum_{i=1}^m \|x_i - D\alpha_i\|_2^2 + \lambda_1 \|\alpha_i\|_1 \quad (2)$$

Note also that since the reconstruction errors $\|x_i - D\alpha_i\|_2^2$ are invariant to scaling simultaneously D by a scalar and α_i by its inverse, we need to constrain the ℓ_2 norm of the columns of D . Such a constraint is classical in sparse coding [28]. This reconstructive approach provides sparse codes α_i for each signal x_i , which can be used a posteriori in a regular classifier such as logistic regression, which would require to solve

$$\min_{\theta} \sum_{i=1}^m c(y_i f(x_i, \alpha_i, \theta)) + \lambda_2 \|\theta\|_2^2 \quad (3)$$

where C is the logistic loss function ($C(x) = \log(1 + e^{-x})$), which enjoys properties similar to that of the hinge loss from the SVM literature, while being differentiable, and λ_2 is a regularization parameter, which prevents over fitting. This is the approach chosen in [29] (with SVMs). However, our goal is to learn jointly D and the model parameters θ . To that effect, we propose the formulation

$$\min_{D, \theta, \alpha} \left(\sum_{i=1}^m c(y_i f(x_i, \alpha_i, \theta)) + \lambda_0 \sum_{i=1}^m \|x_i - D\alpha_i\|_2^2 + \lambda_1 \|\alpha_i\|_1 \right) + \lambda_2 \|\theta\|_2^2 \quad (4)$$

where λ_0 controls the importance of the reconstruction term, and the loss for a pair (x_i, y_i) is

$$S^*(x_i, D, \theta, y_i) = \min_{\alpha} S(\alpha, x_i, D, \theta, y_i) \quad (5)$$

Where $S(\alpha, x_i, D, \theta, y_i) = c(y_i f(x_i, \alpha_i, \theta)) + \lambda_0 \|x_i - D\alpha_i\|_2^2 + \lambda_1 \|\alpha_i\|_1$. In this setting, the classification procedure of a new signal x with an unknown label y , given a learned dictionary D and parameters θ , involves supervised sparse coding:

$$\min_{x \in \{-1, +1\}} S^*(x, D, \theta, y) \quad (6)$$

The learning procedure of Eq. (4) minimizes the sum of the costs for the pairs $(x_i, y_i)_{i=1}^m$ and corresponds to a generative model. We will refer later to this model as SDL-G (supervised dictionary learning, generative). Note the explicit incorporation of the reconstructive and discriminative component into sparse coding, in addition to the classical reconstructive term (see [29] for a different classification component).

However, since the classification procedure from Eq. (6) compares the different costs $S^*(x, D, \theta, y)$ of a given signal for each class $y = -1, +1$, a more discriminative approach is to not only make the costs $S^*(x_i, D, \theta, -y_i)$ small, as in (4), but also make the value of $S^*(x_i, D, \theta, -y_i)$ greater than $S^*(x_i, D, \theta, y_i)$, which is the purpose of the logistic loss function C. This leads to:

$$\min_{D, \theta} \left(\sum_{i=1}^m c(S^*(x_i, D, \theta, -y_i) - S^*(x_i, D, \theta, y_i)) \right) + \lambda_2 \|\theta\|_2^2 \quad (7)$$

As detailed below, this problem is more difficult to solve than (4), and therefore we adopt instead a mixed formulation between the minimization of the generative Eq. (4) and its discriminative version (7), (see also [30])—that is,

$$\left(\sum_{i=1}^m \mu c(S^*(x_i, D, \theta, -y_i) - S^*(x_i, D, \theta, y_i)) + (1 - \mu) S^*(x_i, D, \theta, y_i) \right) + \lambda_2 \|\theta\|_2^2 \quad (8)$$

where μ controls the trade-off between the reconstruction from Eq. (4) and the discrimination from Eq. (7). This is the proposed generative/discriminative model for sparse signal representation and classification from learned dictionary D and model θ . We will refer to this mixed model as SDL-D, (supervised dictionary learning, discriminative). Note also that, again, we constrain the norm of the columns of D to be less than or equal to one.

All of these formulations admit a straightforward multiclass extension, using softmax discriminative cost functions $c_i(x_1, \dots, x_p) = \log\left(\sum_{j=1}^p e^{x_j - x_i}\right)$, which are multiclass versions of the logistic function, and learning one model θ_i per class. Other possible approaches such as one-vs-all or one-vs-one are of course possible, and the question of choosing the best approach among these possibilities is still open. Compared with earlier work using one dictionary per class [32], our model has the advantage of letting multiple classes share some features, and uses the coefficients of the sparse representations as part of the classification procedure, thereby following the works from [31, 32, 29], but with learned representations optimized for the classification task similar to [30, 33].

Our bilinear model with $f(x, \alpha, \theta) = x^T w \alpha + b$ does not admit a straightforward probabilistic interpretation. On the other hand, it can easily be interpreted in terms of kernels: Given two signals x_1 and x_2 , with coefficients α_1 and α_2 , using the kernel $K(x_1, x_2) = \alpha_1^T \alpha_2 x_1^T x_2$ in a logistic regression classifier amounts to finding a decision function of the same form as f . It is a product of two linear kernels, one on the α 's and one on the input signals x . Interestingly, Raina et al. [29] learn a dictionary adapted to reconstruction on a training set, then train an SVM a posteriori on the decomposition coefficients. They derive and use a Fisher kernel, which can be written as $K'(x_1, x_2) = \alpha_1^T \alpha_2 r_1^T r_2$ in this setting, where the r 's are the residuals of the decompositions. In simple experiments, which are not reported in this paper, we have observed that the kernel K , where the signals x replace the residuals r , generally yields a level of performance similar to K' and often actually does better when the number of training samples is small or the data are noisy.

4. Experiments

To evaluate the performance of the proposed algorithm, we carry out a series of experiments on a dataset extracted 500 images of size 48*96 from a video. If the image is contain a pedestrian, the label of it will be 1, else will be -1. Fig. 2(a) shows several images with label 1. Fig. 2(b) shows several images with label -1. 100 images from the dataset are selected as the test examples. Different number images of the dataset are selected as the training examples to compare the accuracy rate.

Fig.3 shows the compare results of recognition between with HOG, HOF and Color features respectively and with the corresponding sparse features. Fig.4 shows the result of using mixing features to compare the two methods. As shown in the graph, our method performs better than the method directly using HOG, HOF

and Color features to recognition. In addition, with the increasing number of training samples, our method performs better.

Fig.5 shows the result of these two methods using shading images to test. Compared with the traditional method, our method has better recognition accuracy and shows good robustness

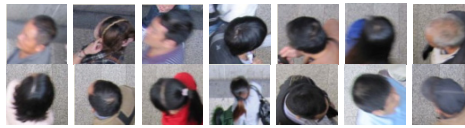


Fig.2 (a) images with label 1

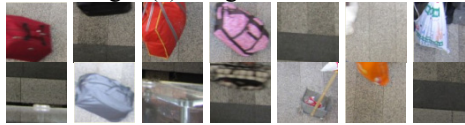


Fig.2 (b) images with label -1.

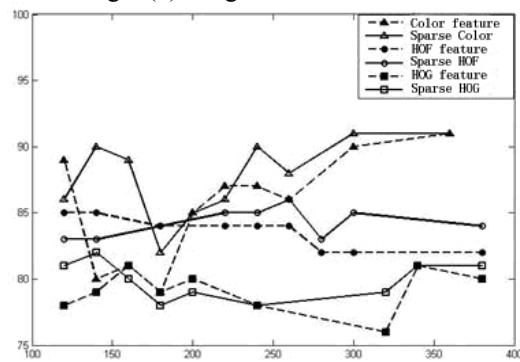


Fig.3 The compare results of recognition between with HOG, HOF and Color features respectively and with the corresponding sparse features

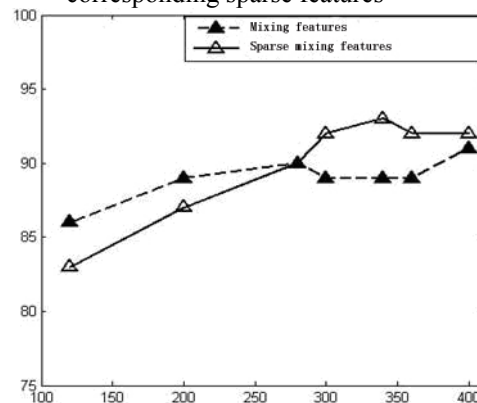


Fig.4 The result of using mixing features to compare the two methods

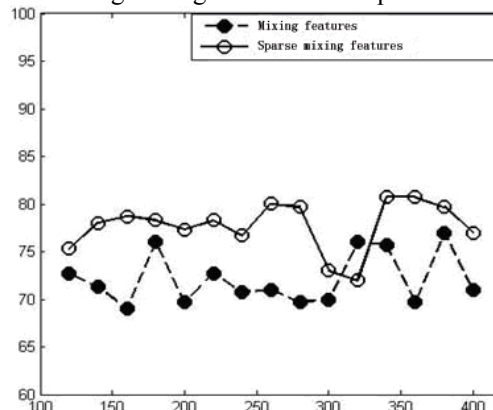


Fig.5 The result of these two methods using shading images to test

5. Conclusion

We proposed a system for pedestrian detection with very good accuracy. To achieve good classification performance, we put forward a novel framework for pedestrian detection tasks, which proposing a model with both sparse reconstruction and class discrimination components, jointly optimized during dictionary learning. We present an efficient pedestrian detection system using mixing sparse features of HOG ,FOG and CSS to combine into this a Kernel classifier. Results presented on our data set show competitive accuracy and robust performance of our system outperforms current state-of-the-art work. Although we use the system for the detection of pedestrians, the general idea can be applied to the detection of other object classes as well.

6. Acknowledgments.

This research is supported by the national science foundation of China (NFSC) No.61170126,61003183, provincial universities natural science foundation of Jiangsu province(11KJD520004) , funding of Jiangsu innovation program for graduate education (CXZZ11_0216) and international communication foundation of Nanjing University of Aeronautics and Astronautics.

7. References

- [1] Ess, Andreas Schindler, Konrad; Leibe, Bastian; Van Gool, Luc Object detection and tracking for autonomous navigation in dynamic environments International Journal of Robotics Research, v 29, n 14, p 1707-1725, December 2010
- [2] Davies, Anthony C.; Yin, Jia Hong; Velastin, Sergio A. Crowd monitoring using image processing Electronics and Communication Engineering Journal, v 7, n 1, p 37-47, Feb 1995
- [3] Kim, Chan-Young Sin, Bong-Kee Human gait analysis using Self Organizing Map.Proceedings of the 2009 Chinese Conference on Pattern Recognition, CCPR 2009, and the 1st CJK Joint Workshop on Pattern Recognition, CJKPR, p 888-891
- [4] Ma, Guanglin Park, Su-Birm; Ioffe, Alexander; Müller-Schneiders, Stefan; Kummert, Anton.A real time object detection approach applied to reliable pedestrian detection. IEEE Intelligent Vehicles Symposium, Proceedings, p 755-760, 2007
- [5] Marana, A.N. Cavenaghi, M.A.; Ulson, R.S.; Drumond, F.L. Real-time crowd density estimation using images. Lecture Notes in Computer Science, v 3804 LNCS, p 355-362
- [6] Kong, Chunyu ;Yang, Jikuang; Nie, Jin. A study on pedestrian detection models based on the analysis on real accident scenarios. Qiche Gongcheng/Automotive Engineering, v 32, n 11, p 977-983, November 2010
- [7] S. Mallat. A wavelet tour of signal processing, second edition. Academic Press, New York, 1999.
- [8] . Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman. Learning discriminative dictionaries for local image analysis. In CVPR, 2008.
- [9] M. Ranzato and M. Szummer. Semi-supervised learning of compact document representations with deep networks. In ICML, 2008.
- [10] A. Argyriou and T. Evgeniou and M. Pontil Multi-Task Feature Learning. In NIPS, 2006.
- [11] F. Rodriguez and G. Sapiro. Sparse representations for image classification: Learning discriminative and reconstructive non-parametric dictionaries. IMA Preprint 2213, 2007.
- [12] D. Blei and J. McAuliffe. Supervised topic models. In NIPS, 2007.
- [13] A. Holub and P. Perona. A discriminative framework for modeling object classes. In CVPR, 2005.
- [14] J.A. Lasserre, C.M. Bishop, and T.P. Minka. Principled hybrids of generative and discriminative models. In CVPR, 2006.
- [15] R. Raina, Y. Shen, A. Y. Ng, and A. McCallum. Classification with hybrid generative/discriminative models. In NIPS, 2004.
- [16] R. R. Salakhutdinov and G. E. Hinton. Learning a non-linear embedding by preserving class neighbourhood structure. In AI and Statistics, 2007.

- [17] H. Larochelle, and Y. Bengio. Classification using discriminative restricted boltzmann machines. in ICML, 2008.
- [18] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In CVPR, 2005.
- [19] N. Dalal, B. Triggs, and C. Schmid. Human detection using oriented histograms of flow and appearance. In ECCV, 2006.
- [20] M. Werlberger, W. Trobin, T. Pock, A. Wedel, D. Cremers, and H. Bischof. Anisotropic Huber-L1 optical flow. In BMVC, 2009.
- [21] P. Dollár, Z. Tu, P. Perona, and S. Belongie. Integral channel features. In BMVC, 2009.
- [22] W. Schwartz, A. Kembhavi, D. Harwood, and L. Davis. Human detection using partial least squares analysis. In ICCV, 2009.
- [23] C. Wojek, S. Walk, and B. Schiele. Multi-cue onboard pedestrian detection. In CVPR, 2009.
- [24] K. E. A. van de Sande, T. Gevers, and C. G. M. Snoek. Evaluation of color descriptors for object and scene recognition. In PAMI, 2009.
- [25] T. Watanabe, S. Ito, and K. Yokoi. Co-occurrence histograms of oriented gradients for pedestrian detection. In PSIVT, 2009.
- [26] B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *Ann. Stat.*, 32(2), 2004.
- [27] M. Ranzato, C. Poultney, S. Chopra, and Y. LeCun. Efficient learning of sparse representations with an energy-based model. In NIPS, 2006.
- [28] M. Elad and M. Aharon. Image denoising via sparse and redundant representations over learned dictionaries. *IEEE Trans. IP*, 54(12), 2006.
- [29] F. Rodriguez and G. Sapiro. Sparse representations for image classification: Learning discriminative and reconstructive non-parametric dictionaries. IMA Preprint 2213, 2007.
- [30] R. Raina, Y. Shen, A. Y. Ng, and A. McCallum. Classification with hybrid generative/discriminative models. In NIPS, 2004.
- [31] K. Huang and S. Aviyente. Sparse representation for signal classification. In NIPS, 2006.
- [32] J. Wright, A. Y. Yang, A. Ganesh, S. Sastry, and Y. Ma. Robust face recognition via sparse representation. In PAMI, 2008. to appear.
- [33] D. Blei and J. McAuliffe. Supervised topic models. In NIPS, 2007.