

SPARSE CODING WITH FAST IMAGE ALIGNMENT VIA LARGE DISPLACEMENT OPTICAL FLOW

Xiaoxia Sun[†], Nasser M. Nasrabadi[‡] and Trac D. Tran[†]

[†] Department of ECE, Johns Hopkins University, 3400 N. Charles Street, Baltimore, MD, 21218

[‡] Lane Department of CSEE, West Virginia University, 395 Evansdale Drive, Morgantown, WV, 26506

ABSTRACT

Sparse representation-based classifiers have shown outstanding accuracy and robustness in image classification tasks even with the presence of intense noise and occlusion. However, it has been discovered that the performance degrades significantly either when test image is not aligned with the dictionary atoms or the dictionary atoms themselves are not aligned with each other, in which cases the sparse linear representation assumption fails. In this paper, having both training and test images misaligned, we introduce a novel sparse coding framework that is able to efficiently adapt the dictionary atoms to the test image via large displacement optical flow. In the proposed algorithm, every dictionary atom is automatically aligned with the input image and the sparse code is then recovered using the adapted dictionary atoms. A corresponding supervised dictionary learning algorithm is also developed for the proposed framework. Experimental results on digit datasets recognition verify the efficacy and robustness of the proposed algorithm.

1. INTRODUCTION

Sparse coding has been successfully applied to numerous computer vision tasks, including face recognition [1], scene categorization [2] and object detection [3]. Application of sparse representation-based classifier (SRC) on face recognition [1] demonstrates a startling robustness over noise and occlusions, where the test subjects are still recognizable even when they wear sunglasses or scarf. However, SRC has been found to be highly sensitive to the misalignment of the image dataset: a small amount of image distortion due to translation, rotation, scaling and 3-dimensional pose variations can lead to a significant degradation on the classification performance [4].

One straightforward way to solve the misalignment problem is to register the test image with dictionary atoms before sparse recovery. By assuming the dictionary atoms are registered, Wagner *et al.* [4] parameterize the misalignment of the test image with an affine transformation. These parameters

are optimized using generalized Gauss-Newton methods after linearizing the affine transformation constraints. By minimizing the sparse registration error iteratively and sequentially for each class, their framework is able to deal with a large range of variations in translation, scaling, rotation and even 3D pose variations. Due to the adoption of holistic features, sparse coding is more robust and less likely to overfit.

In the case of local feature-based sparse coding, max pooling strategy [5] is often employed over the neighboring coefficients to produce local translation-invariant property. Based on spatial pyramid matching framework, Yang *et al.* [2] proposed a local sparse coding model with local SIFT features followed by multi-scale max pooling. The results on several large variance datasets achieved plausible performance that can hardly be pursued by simply applying holistic sparse coding. To improve the discriminability of the sparse codes, their dictionary was trained with supervised learning via back-propagation [6]. Classification performance of local feature-based sparse coding has also been evaluated on several large datasets in [7], demonstrating a state-of-art performance that is competitive with deep learning [8, 9]. Another interesting approach is the convolutional sparse coding [10], where the local features are reconstructed by convoluting the local sparse codes using local dictionary. Visualization of its dictionary shows that the dictionary atoms contain more complex features, therefore having more discriminative power.

In this paper, we present a novel sparse coding framework that is robust to image transformation. In the proposed model, each dictionary atom is constructed in the form of a tensor and is aligned with the test image using the large displacement optical flow concept [11]. We show experimentally that the proposed sparse coding framework outperforms most other sparsity-based methods. Specifically, our paper has the following novelties and contributions: (i) The proposed algorithm does not require the training dataset to be pre-aligned. (ii) Adapting the dictionary to the input test image is highly efficient: requiring only $\mathcal{O}(PT)$ operations for adapting each dictionary atom, where T is the number of pixels in a searching window and P is the total number of *subatoms* to be aligned. (iii) Supervised dictionary learning algorithm is developed for the proposed sparse coding framework.

The remainder of the paper is organized as follows: We first

This work has been partially supported by NSF under Grants NSF-CCF-1117545, NSF-CCF-1422995 and NSF-ECS-1443936.

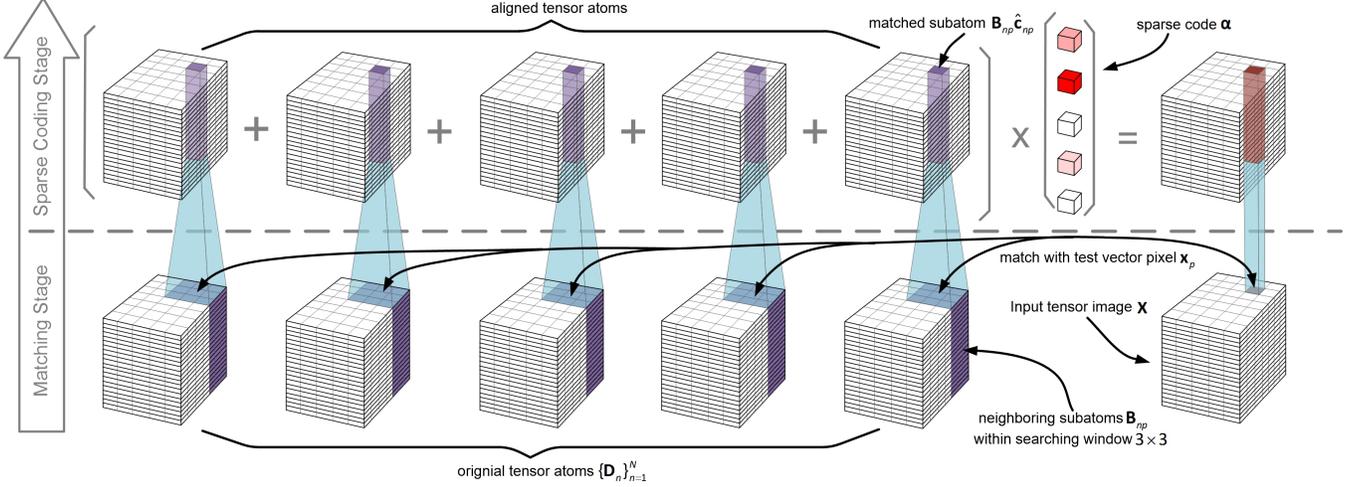


Fig. 1: Proposed sparse coding framework: Dictionary tensor atoms $\{\mathbf{D}_n\}_{n=1}^N$ and the test tensor image \mathbf{X} are shown in the lower part of the figure. Searching window of size $T = 3 \times 3$ within each tensor atom is colored with purple. Each group of neighboring T subatoms \mathbf{B}_{np} is matched with the corresponding vector pixel \mathbf{x}_p of the test tensor image, resulting in an aligned subatom $\mathbf{B}_{np}\hat{\mathbf{c}}_{np}$. After the matching process, the sparse code for \mathbf{x}_p is recovered using all the aligned subatoms $\{\mathbf{B}_{np}\hat{\mathbf{c}}_{np}\}_{n=1}^N$. For illustration purposes, only five dictionary tensor atoms are shown in the figure and the magnitude of the sparse codes are displayed with various intensities in red.

introduce the proposed sparse coding framework for dealing with dataset misalignment in Section 2. Next, in Section 3, we show how to train the dictionary in a supervised manner by solving a bilevel optimization problem. Finally, in Section 4, experimental results demonstrate that the proposed framework has a state-of-art performance, which is more promising over most existing sparsity-based methods.

2. SPARSE CODING WITH IMAGE ALIGNMENT VIA LARGE DISPLACEMENT OPTICAL FLOW

In this section, we first introduce how to construct the dictionary atoms and input images in the form of tensors. We then illustrate how to eliminate the misalignment by dynamically adapt the tensor dictionary atoms to the input tensor image.

In the proposed sparse coding model, as shown in Fig. 1, both dictionary atom and input image are represented by image tensors. Each pixel in the tensor image is a vectorized version of a local patch in the original image, referred to as a vector pixel. Denote the n^{th} tensor atom as $\mathbf{D}_n = [\mathbf{d}_{n1}, \dots, \mathbf{d}_{nP}] \in \mathbb{R}^{M \times P}$ and a given test tensor image as $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_P] \in \mathbb{R}^{M \times P}$, where $\mathbf{d}_{np} \in \mathbb{R}^M$ is the p^{th} subatom of the n^{th} tensor atom and $\mathbf{x}_p \in \mathbb{R}^M$ is the p^{th} vector pixel of the input image. M is the dimension of vector pixel, n is the dictionary atom index and P is the total number of subatoms in the tensor atom, which is the same number of vector pixels in the test tensor image. The dictionary is denoted as $\mathbf{D} = [\mathbf{D}_1, \dots, \mathbf{D}_N] \in \mathbb{R}^{M \times NP}$. Given a dictionary with N tensor atoms, a typical sparse recovery problem [1] is

formulated as:

$$\hat{\alpha} = \arg \min_{\alpha} \frac{1}{2} \sum_{p=1}^P \left\| \sum_{n=1}^N \alpha_n \mathbf{d}_{np} - \mathbf{x}_p \right\|_2^2 + \lambda \|\alpha\|_1, \quad (1)$$

where $\alpha = [\alpha_1, \dots, \alpha_N]^T \in \mathbb{R}^N$ is the sparse coefficient and $\lambda > 0$ is the regularization parameter. Problem (1) is a standard form of ℓ_1 -sparse recovery problem that can be efficiently solved using alternating direction method of multipliers (ADMM) [12].

When images in both the training and test datasets are misaligned, sparse coefficients recovered by solving the problem (1) become unreliable, thus resulting in poor classification performance. To alleviate the misalignment problem, we propose to register each tensor atom with the input test image via large displacement optical flow [11]. The notion of optical flow field is used here to describe the displacements of vector pixels within each tensor atom, and the sparse recovery is then performed by using only the best matching subatoms selected from the tensor atoms. The proposed framework is illustrated in Fig. 1. Denote $\mathbf{B}_{np} \in \mathbb{R}^{M \times T}$ as the T subatoms within the searching window centered at the location p of the n^{th} tensor atom. The recovery of the optical flow and sparse codes can be formally described as follows:

$$\begin{aligned} (\hat{\alpha}, \{\hat{\mathbf{c}}_{np}\}) &= \arg \min_{\alpha, \{\mathbf{c}_{np}\}} \frac{1}{2} \sum_{p=1}^P \left\| \sum_{n=1}^N \alpha_n \mathbf{B}_{np} \mathbf{c}_{np} - \mathbf{x}_p \right\|_2^2 + \lambda \|\alpha\|_1, \\ \text{s.t. } & \|\mathbf{c}_{np}\|_0 = 1, \|\mathbf{c}_{np}\|_1 = 1, \mathbf{c}_{np} \geq \mathbf{0}, \\ & \forall n \in [N], p \in [P], \end{aligned} \quad (2)$$

where $\|\mathbf{c}_{np}\|_0 = 1$ is the cardinality constraint and $\mathbf{c}_{np} \in \mathbb{R}^T$ is the sparse index vector that is used to characterize the optical flow field. The constraint in (2) suggests that \mathbf{c}_{np} is a binary index vector and only one element is nonzero, which means that it can only select one subatom within the searching window.

The optimization problem in (2) is a mixed-integer problem and NP-hard [13]. Therefore, we propose a heuristic algorithm to find an informative α and the sparse index vectors $\{\mathbf{c}_{np}\}_{n,p=1}^{N,P}$ for all vector pixels. As shown in Fig. (1), the optical flow field for each vector pixel is found by searching for the best match between neighboring subatoms and the corresponding input vector pixel. In practice, we found that searching for the best match without involving the sparse code is the key to render plausible performance in both classification accuracy and computational efficiency. Formally, we propose to find a local optimum of problem (2) by solving the following optimization problem:

$$\begin{aligned} \hat{\alpha} &= \arg \min_{\alpha} \frac{1}{2} \sum_{p=1}^P \left\| \sum_{n=1}^N \alpha_n \mathbf{B}_{np} \hat{\mathbf{c}}_{np} - \mathbf{x}_p \right\|_2^2 + \lambda \|\alpha\|_1 \\ \text{s.t. } \hat{\mathbf{c}}_{np} &= \arg \min_{\mathbf{c}_{np}} \frac{1}{2} \|\mathbf{B}_{np} \mathbf{c}_{np} - \mathbf{x}_p\|_2^2, \\ \|\mathbf{c}_{np}\|_0 &= 1, \|\mathbf{c}_{np}\|_1 = 1, \mathbf{c}_{np} \geq \mathbf{0}, \\ \forall n \in [N], p \in [P]. \end{aligned} \quad (3)$$

In our approach, the sparse coding part of (3) is solved by using the alternating direction method of multipliers (ADMM) [12]. One important advantage of the above model is that it is highly computational efficient because it only takes $\mathcal{O}(T)$ operations to search for the best match for each vector pixel.

3. SUPERVISED DICTIONARY LEARNING

In order to improve the efficiency of sparse coding and discriminability of the dictionary, we employ the supervised dictionary learning framework [6, 14, 15] to optimize the dictionary and the classifier parameters simultaneously. Formulated as a bilevel optimization problem, the dictionary is updated using back propagation to minimize the classification error. Formally, the supervised dictionary learning problem can be formulated as follows:

$$\min_{\mathbf{W}, \mathbf{D}} \mathbb{E}_{\mathbf{y}, \mathbf{X}} [\ell(\mathbf{y}, \mathbf{W} \hat{\alpha}(\mathbf{X}, \{\hat{\mathbf{c}}_{np}(\mathbf{D})\}, \mathbf{D}))] + \frac{\mu}{2} \|\mathbf{W}\|_F^2, \quad (4)$$

where $\ell(\cdot)$ is some smooth and convex function that is used to define the classification error and $\mu > 0$ is the regularization parameter used to alleviate the overfitting of the classifier. Due to the triviality of updating classifier parameters, here we only state the update for the dictionary:

$$\mathbf{D} \leftarrow \Pi(\mathbf{D} - \rho^t \cdot \partial \ell / \partial \mathbf{D}), \quad (5)$$

where $\rho > 0$ is the learning rate, t is the iteration counter and Π is the projection that regulate the Frobenius norm of

every tensor atom to be one. Similar to [6, 14, 15], (4) suggests that the update of both the dictionary and the classifier are driven by reducing classification error. The local optima can be solved by using descent method [16] based on error backpropagation. The sparse code α is an implicit function of \mathbf{X} , $\{\mathbf{c}_{np}\}$ and \mathbf{D} . In addition, each optical flow field \mathbf{c}_{np} is an implicit function of \mathbf{D} and \mathbf{x}_{np} . Therefore, given an input image \mathbf{X} and an optimal sparse code $\hat{\alpha}$, apply the chain rule of differentiation, the direction along which the upper-level cost decreases can be formulated as:

$$\frac{\partial \ell(\mathbf{y}, \mathbf{W} \alpha)}{\partial \mathbf{D}} = \frac{\partial \ell}{\partial \alpha} \frac{\partial \alpha}{\partial \mathbf{D}} + \sum_{p=1}^P \frac{\partial \ell}{\partial \mathbf{C}_p} \frac{\partial \mathbf{C}_p}{\partial \mathbf{D}}, \quad (6)$$

where $\mathbf{C}_p = \bigoplus_{n=1}^N \bar{\mathbf{c}}_{np} \in \mathbb{R}^{NP \times N}$ and \bigoplus denotes the direct sum. Also, $\bar{\mathbf{c}}_{np} \in \mathbb{R}^{NP}$ is obtained by zero-padding with \mathbf{c}_{np} , where $(N-1)P+1$ to NP elements of $\bar{\mathbf{c}}_{np}$ are from those of \mathbf{c}_{np} . Due to the binary constraints on $\{\mathbf{c}_{np}\}$, every element of the gradient $\partial \mathbf{C}_p / \partial \mathbf{D}$ equals to zero. On the other hand, the first part of the derivative can be solved by applying fixed point differentiation [17]. Due to the page limitation of the paper and the triviality for deriving the term $\partial \ell / \partial \alpha$, we only show the final derivation of $\partial \alpha / \partial \mathbf{D}$ as follows:

$$\frac{\partial \alpha_{\Lambda}}{\partial d_{mnp}} = \Theta_{\Lambda, \Lambda}^{-1} \left(\frac{\partial (\mathbf{D} \mathbf{C}_p)_{\Lambda}^{\top}}{\partial d_{mnp}} \mathbf{x}_p - \frac{\partial \Theta_{\Lambda, \Lambda}}{\partial d_{mnp}} \alpha_{\Lambda} \right), \quad (7)$$

where Λ is the index set of active atoms of the sparse code α . $(\mathbf{D} \mathbf{C}_p)_{\Lambda}$ is the matrix obtained by collecting the active columns of $\mathbf{D} \mathbf{C}_p$, $\Theta = \sum_{p=1}^P \mathbf{C}_p^{\top} \mathbf{D}^{\top} \mathbf{D} \mathbf{C}_p$ and $\Theta_{\Lambda, \Lambda}$ is the submatrix obtained by selecting the active columns and rows of Θ . The matrix $\Theta_{\Lambda, \Lambda}$ is always nonsingular since the total number of measurement MP is always significantly larger than the number of active atoms. Combining (6) with (7) for each dictionary element, the gradient for updating the dictionary can be achieved. For a large dataset, the dictionary and the classifier parameters are updated in an online manner.

4. EXPERIMENTS

In this section, we evaluate the proposed algorithm on handwritten digits datasets including the MNIST and USPS. The sparse coding is performed with a single dictionary and linear SVM is used for classification. For a fair comparison, we only compare with the results that are produced with the same SRC strategy. The dictionary size in our paper is set to be no larger than those used in other methods. Similar to [6], parameters in our experiments are chosen heuristically. The batch size for updating the dictionary is 512. Initial learning rate ρ is set to 0.001 and $\lambda = 0.01$.

4.1. Evaluation on the MNIST Database

MNIST [18] consists of a total number of 70,000 images of digits, of which 60,000 are training set and the rest 10,000

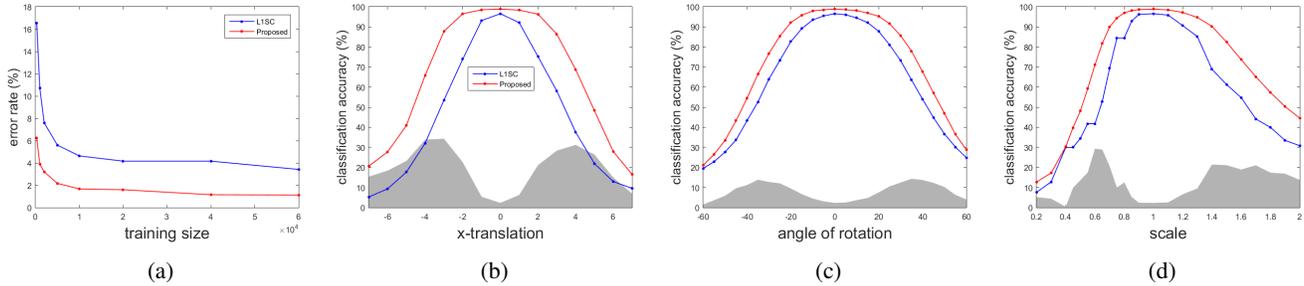


Fig. 2: The proposed method demonstrates plausible performance on MNIST digits recognition with a small number of training samples. It also demonstrates robustness towards various image deformations. Classification accuracy of different experimental settings are shown in the above sub-figures: (a) Error rate under various sizes of training samples. (b) Translation along x direction versus classification accuracy. (c) In-plane rotation only. (d) Scale variation only. In (b)-(c), red and blue lines are the results of the proposed method and L1SC, respectively. Gray shadow area at the bottom of each figure is the accuracy difference between the proposed method and L1SC.

are test set. Each digit is centered and normalized in a 28×28 field. The dictionary size N is set to be 150 for this database.

We first evaluate the performance of the proposed algorithm under various number of training samples. We follow the same experimental setting as in [19], examining the classification accuracy given the training size $\{300, 1K, 2K, 5K, 10K, 20K, 40K, 60K\}$. The performance is shown in Fig. 2 (a). The proposed method significantly outperforms the ℓ_1 sparse coding-based algorithm (L1SC) [15].

We then demonstrate the robustness of the proposed method towards various image deformations. Following a similar setting as in [4], we perform the translation along x direction, rotation and scaling separately only on the test samples. We report the classification accuracy with respect to various levels of deformation and compare the performance with L1SC. The experimental results are shown in Fig. 2(b)-(d). Performance of our method and L1SC are illustrated in red and blue lines, respectively. The shadow area at the bottom of each figure is the accuracy difference between the two methods. We can see for all three deformations, the proposed method consistently outperforms L1SC. In addition, the hump shape of the shadow area indicates that the proposed method is robust to numerous image deformations.

Finally, the error rate for the MNIST is shown in Table 1. Our method reaches the lowest error rate of 1.12%. On MNIST, differences of more than 0.1% are statistically significant [20]. Comparing with the second best algorithm, the proposed method reduces the error rate by 0.12%, exhibiting better generality and dictionary compactness.

4.2. Evaluation on the USPS Database

The USPS dataset has 7,291 training and 2,007 test images, where each of them is of size 16×16 . Being compared to MNIST, the USPS dataset has a much larger variance and a smaller training set, which challenges the dictionary generality. For a fair comparison, the dictionary size N is set

Method	MNIST	USPS
CBN	1.95 (3×10^4)	4.14 (7291)
ESC [21]	5.16 (150)	6.03 (80)
Ramirez <i>et al.</i> [22]	1.26 (800)	3.98 (80)
Deep Belief Network [8]	1.25 (-)	- (-)
MMDL [23]	1.24 (150)	-(-)
Proposed	1.12 (150)	3.43 (80)
Improvements	9.7%	13.8%

Table 1: Error Rate (%) on MNIST and USPS datasets. The dictionary size is shown in the parentheses. Improvements over the second best algorithm is shown in the last line.

to be 80. Local patch size is 5×5 ($M = 25$). Searching window size is 5×5 ($T = 25$). The performance of various approaches on USPS database are depicted in Table 1. Our algorithm achieves the lowest error rate 3.43% among other supervised learning-based methods. The experimental result validates the efficacy of our proposed algorithm on a dataset with a larger variance.

5. CONCLUSION

In this paper, we present a novel sparse coding algorithm that is able to dynamically select the dictionary subatoms to adapt to the misaligned image dataset. In the proposed method, both the dictionary atoms and the input test image are represented by tensors, and each vector pixel in the tensor image is a vectorized local patch. Each tensor atom is aligned with the input tensor image using large displacement optical flow, which is highly computationally efficient. Using the fixed point differentiation, a supervised dictionary learning algorithm is developed for the proposed sparse coding framework, which significantly reduces the required dictionary size.

References

- [1] J. Wright, A.Y. Yang, A. Ganesh, S.S. Sastry, and Y. Ma, "Robust face recognition via sparse representation," *IEEE Trans. on Pattern Anal. and Mach. Intell.*, vol. 31, no. 2, pp. 210–227, Feb. 2009.
- [2] J. Yang, K. Yu, Y. Gong, and T. Huang, "Linear spatial pyramid matching using sparse coding for image classification," in *CVPR*, pp. 1794–1801, Jun. 2009.
- [3] S. Agarwal and D. Roth, "Learning a sparse representation for object detection," in *ECCV*, vol. 4, pp. 113–130, May 2002.
- [4] A. Wagner, J. Wright, A. Ganesh, Z. Zhou, H. Mobahi, and Y. Ma, "Toward a practical face recognition system: Robust alignment and illumination by sparse representation," *IEEE Trans. on Pattern Anal. and Mach. Intell.*, vol. 34, no. 2, pp. 372–386, Feb. 2012.
- [5] H. Lee, R. B. Grosse, R. Ranganath, and A. Y. Ng, "Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations," in *ICML*, Jun. 2009.
- [6] J. Yang, K. Yu, and T. Huang, "Supervised translation-invariant sparse coding," in *CVPR*, pp. 3517–3524, Jun. 2010.
- [7] A. Coates and A. Y. Ng, "The importance of encoding versus training with sparse coding and vector quantization," in *ICML*, pp. 921–928, Jul. 2011.
- [8] G. E. Hinton and S. Osindero, "A fast learning algorithm for deep belief nets," *Neural Computation*, vol. 18, no. 7, pp. 527–1554, Jul. 2006.
- [9] J. Mairal, P. Koniusz, Z. Harchaoui, and C. Schmid, "Partial Face Recognition: A Sparse Representation-based Approach," in *ICASSP*, Mar. 2016.
- [10] K. Kavukcuoglu, P. Sermanet, Y. Boureau, K. Gregor, M. Mathieu, and Y. Lecun, "Learning convolutional feature hierarchies for visual recognition," in *NIPS*, pp. 1090–1098, Dec. 2010.
- [11] T. Brox and J. Malik, "Large displacement optical flow: Descriptor matching in variational motion estimation," *IEEE Trans. on Pattern Anal. and Mach. Intell.*, vol. 33, no. 3, pp. 500–513, Mar. 2011.
- [12] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Journal FTML*, vol. 3, no. 1, pp. 1–122, Jan. 2011.
- [13] D. Bertsimas and R. Weismantel, "Optimization over integers," *Athena Scientific*, 2005.
- [14] J. Yang, Z. Wang, Z. Lin, X. Shu, and T. Huang, "Bilevel sparse coding for coupled feature spaces," in *CVPR*, pp. 2360–2367, Jun. 2012.
- [15] J. Mairal, F. Bach, and J. Ponce, "Task-driven dictionary learning," *IEEE Trans. on Pattern Anal. and Mach. Intell.*, vol. 34, no. 4, pp. 791–804, Apr. 2012.
- [16] B. Colson, P. Marcotte, and G. Savard, "An overview of bilevel optimization," *Ann. of Operat. Res.*, vol. 153, no. 1, pp. 235–256, Apr. 2007.
- [17] D. M. Bradley and J. A. Bagnell, "Differentiable sparse coding," in *NIPS*, Dec. 2008.
- [18] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, Nov. 1998.
- [19] J. Mairal, P. Koniusz, Z. Harchaoui, and C. Schmid, "Convolutional Kernel Networks," in *NIPS*, 2014.
- [20] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle, "Greedy layer-wise training of deep networks," in *NIPS*, pp. 153–160, Dec. 2007.
- [21] H. Lee, A. Battle, R. Raina, and A. Y. Ng, "Efficient sparse coding algorithms," in *NIPS*, pp. 801–808, Dec. 2006.
- [22] I. Ramirez, P. Sprechmann, and G. Sapiro, "Classification and clustering via dictionary learning with structured incoherence and shared features," in *CVPR*, pp. 3501–3508, Jun. 2010.
- [23] Z. Wang, J. Yang, N. M. Nasrabadi, and T. Huang, "A max-margin perspective on sparse representation-based classification," in *ICCV*, pp. 1217–1224, Dec. 2013.