

Supplementary Materials for Block and Group Regularized Sparse Modeling for Dictionary Learning

Yu-Tseh Chi*

ychi@cise.ufl.edu

Mohsen Ali*

moali@cise.ufl.edu

Ajit Rajwade[†]

ajit.rajwade@daiict.ac.in

Jeffrey Ho*

jho@cise.ufl.edu

Abstract

This document is the supplementary materials for the CVPR 2013 paper "Group and Sparse Regularized Sparse Modeling for Dictionary Learning". The first section contains derivations of the theoretical proof and of the two optimization algorithms in the Sec. 2 of the main text. The second section contains extra results from the experiment in Sec 3.2 of the main text.

1. Supplementary Materials to the Method Section

1.1. Simple Proof of the Theoretical Guarantee

Without loss of generality, we will prove the condition of one group data \mathbf{X} . We first concatenate columns of $\mathbf{X} \in \mathbb{R}^{n \times s}$ into a vector $\mathbf{x}' \in \mathbb{R}^{(n \cdot s)}$. The dictionary $\mathbf{D} = [\mathbf{d}_1 \cdots \mathbf{d}_m]$ of m columns has to be converted accordingly to

$$\begin{aligned} \mathbf{D}' &= \left[\begin{array}{ccc|ccc|ccc} \mathbf{d}_1 & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{d}_2 & \mathbf{0} & \cdots & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{d}_1 & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{d}_2 & \cdots & \mathbf{0} & \cdots & \mathbf{0} \\ & & \vdots & & & \vdots & & & & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{d}_1 & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \cdots & \mathbf{d}_m \end{array} \right] \\ &= [\mathbf{I}_s \otimes \mathbf{d}_1, \mathbf{I}_s \otimes \mathbf{d}_2 \cdots \mathbf{I}_s \otimes \mathbf{d}_m], \end{aligned}$$

where \otimes denotes Kronecker product, and \mathbf{I}_s is a $s \times s$ identity matrix. A block of k atoms, $\mathbf{d}_i \cdots \mathbf{d}_{i+k-1}$, in \mathbf{D} is now a block of size $k \times s$ in \mathbf{D}' . We also concatenate rows of \mathbf{C} into a vector \mathbf{c}' as follows:

$$\mathbf{c}' = [\mathbf{C}_1 \ \mathbf{C}_2 \ \cdots \ \mathbf{C}_m]^\top,$$

where \mathbf{C}_i is the i -th row in \mathbf{C} . Assuming the number of atoms of the i -th block of the original dictionary \mathbf{D} is n_{di} , then $\mathbf{C}_{[i]}$, the i -th block of coefficients \mathbf{C} , is now a segment

of size $s \cdot n_{di}$ in \mathbf{c} . Our program

$$\mathcal{P}_{\ell_{0,p}} : \min_{\mathbf{C}} \sum_i I(\|\mathbf{C}_{[i]}\|_p) \text{ s. t. } \mathbf{X} = \mathbf{D}\mathbf{C} \quad (1)$$

is now

$$\min_{\mathbf{c}'} \sum_i I(\|\mathbf{c}'_{[i]}\|_p) \text{ s. t. } \mathbf{x}' = \mathbf{D}'\mathbf{c}'. \quad (2)$$

Likewise,

$$\mathcal{P}_{\ell_{1,p}} : \min_{\mathbf{C}} \sum_i \|\mathbf{C}_{[i]}\|_p \text{ s. t. } \mathbf{X} = \mathbf{D}\mathbf{C}. \quad (3)$$

is now

$$\min_{\mathbf{c}'} \sum_i \|\mathbf{c}'_{[i]}\|_p \text{ s. t. } \mathbf{x}' = \mathbf{D}'\mathbf{c}'. \quad (4)$$

$\mathcal{P}'_{\ell_{0,p}}$ and $\mathcal{P}'_{\ell_{1,p}}$ can be converted to

$$\min_{\mathbf{c}'} \sum_i I(\|\mathbf{D}'\mathbf{c}'_{[i]}\|_p) \text{ s. t. } \mathbf{x}' = \mathbf{D}'\mathbf{c}', \quad (5)$$

and

$$\min_{\mathbf{c}'} \sum_i \|\mathbf{D}'\mathbf{c}'_{[i]}\|_p \text{ s. t. } \mathbf{x}' = \mathbf{D}'\mathbf{c}', \quad (6)$$

respectively.

Due to the structure of \mathbf{D}' and the matrix norm we use here is element-wise norm, all the restricted isometry constants and block related coherence values defined in [2] of \mathbf{D}' are identical to those of \mathbf{D} . Therefore, given a unique k -block sparse solution of $\mathcal{P}_{\ell_{0,p}}$ and $\mathcal{P}'_{\ell_{0,p}}$, the proof of the equivalence between $\mathcal{P}_{\ell_{1,p}}$ ($\mathcal{P}'_{\ell_{1,p}}$) and $\mathcal{P}_{\ell_{0,p}}$ ($\mathcal{P}'_{\ell_{0,p}}$) under appropriate conditions can be constructed in a similar fashion as in [2].

1.2. Derivation of Reconstructed Block/Group Sparse Coding Algorithm

Continue from Section 2.3 in the main text. For clarity of presentation, we will again derive the reconstructed block/group sparse coding (R-BGSC) algorithm for only

*University of Florida, Gainesville, Florida, U. S. A.

[†]Dhirubhai Ambani Institute of Information and Communication Technology, Gandhinagar, India

one single group of data. $\mathcal{P}'_{\ell_1, p}$ can be cast as an optimization problem that minimizes

$$\frac{1}{2} \|\mathbf{X} - \sum_{i \neq r} \mathbf{D}_{[i]} \mathbf{C}_{[i]} - \mathbf{D}_{[r]} \mathbf{C}_{[r]}\|_F^2 + \lambda \sum_i \|\mathbf{D}_{[i]} \mathbf{C}_{[i]}\|_p. \quad (7)$$

We now derive the crucial steps in the optimization algorithm for $p = 2$ (as the $p = 1$ case is straightforward). Similar to the derivation in the previous section, we first assume the norm $\|\mathbf{D}_{[r]} \mathbf{C}_{[r]}\|_F$ is positive for the optimal solution for $\mathbf{C}_{[r]}$. Taking the gradient of the objective function with respect to $\mathbf{C}_{[r]}$ and equating it to zero, we have

$$-\mathbf{D}_{[r]}^\top \mathbf{X} + \mathbf{D}_{[r]}^\top \sum_{i \neq r} \mathbf{D}_{[i]} \mathbf{C}_{[i]} + \mathbf{D}_{[r]}^\top \mathbf{D}_{[r]} \mathbf{C}_{[r]} + \lambda \mathbf{D}_{[r]}^\top \frac{\mathbf{D}_{[r]} \mathbf{C}_{[r]}}{\|\mathbf{D}_{[r]} \mathbf{C}_{[r]}\|_F} = 0. \quad (8)$$

Now denoting the first two terms with $-\mathbf{N}$ and computing the singular value decomposition of $\mathbf{D}_{[r]} = \mathbf{U} \mathbf{S} \mathbf{V}^\top$, we have

$$\mathbf{V} \mathbf{S}^2 \mathbf{V}^\top \mathbf{C}_{[r]} + \lambda \mathbf{V} \mathbf{S} \frac{\mathbf{S} \mathbf{V}^\top \mathbf{C}_{[r]}}{\|\mathbf{S} \mathbf{V}^\top \mathbf{C}_{[r]}\|_F} = \mathbf{N}. \quad (9)$$

Multiplying both sides of the above equation with \mathbf{V}^\top , and letting $\mathbf{Y} = \frac{\mathbf{S} \mathbf{V}^\top \mathbf{C}_{[r]}}{\|\mathbf{S} \mathbf{V}^\top \mathbf{C}_{[r]}\|_F}$, $\kappa = \|\mathbf{S} \mathbf{V}^\top \mathbf{C}_{[r]}\|_F$, and $\hat{\mathbf{N}} = \mathbf{V}^\top \mathbf{N}$, we have

$$(\kappa \mathbf{S} + \lambda \mathbf{S}) \mathbf{Y} = \hat{\mathbf{N}} \implies \mathbf{Y} = (\kappa \mathbf{S} + \lambda \mathbf{S})^{-1} \hat{\mathbf{N}}, \quad \text{s. t. } \|\mathbf{Y}\|_F = 1. \quad (10)$$

Using the same method as in Section 2.2 in the main text, we can solve for κ first and then compute the iterate of $\mathbf{C}_{[r]}$.

1.3. Derivation of ICS-DL algorithm.

Continue from Eq. (17) in Section 2.4 in the main text, denoting the first two terms with $-\nu_i$, $\mathbf{c}_r \mathbf{c}_r^\top$ with t , and $\sum \mathbf{d}_j \mathbf{d}_j^\top$ with Φ_r , we have

$$t \mathbf{d}_r + \gamma \frac{\mathbf{d}_r}{\|\mathbf{d}_r\|_2} + \beta \Phi_r \mathbf{d}_r = \nu_i \implies t \mathbf{U}^\top \mathbf{d}_r + \gamma \frac{\mathbf{U}^\top \mathbf{d}_r}{\|\mathbf{U}^\top \mathbf{d}_r\|_2} + \beta \Sigma_\Phi \mathbf{U}^\top \mathbf{d}_r = \mathbf{U}^\top \nu_i, \quad (11)$$

where $\mathbf{U} \Sigma_\Phi \mathbf{U}^\top$ is the eigen decomposition of Φ_r and Σ_Φ is a diagonal matrix only containing the non-zero eigen-values of Φ_r , and \mathbf{U} are the corresponding eigen-vectors. Denoting $\frac{\mathbf{U}^\top \mathbf{d}_r}{\|\mathbf{U}^\top \mathbf{d}_r\|_2}$ by \mathbf{y} , $\|\mathbf{U}^\top \mathbf{d}_r\|_2$ by κ , and $\mathbf{U}^\top \nu_i$ by $\tilde{\nu}_i$, respectively, Eq. (11) then becomes

$$\kappa t \mathbf{y} + \gamma \mathbf{y} + \kappa \beta \Sigma_\Phi \mathbf{y} = \tilde{\nu}_i \implies \mathbf{y} = ((\kappa t + \gamma) \mathbf{I} + \kappa \beta \Sigma_\Phi)^{-1} \tilde{\nu}_i, \quad \text{s. t. } \|\mathbf{y}\|_2 = 1. \quad (12)$$

We can use the same methods as in previous sections to solve for the iterate of \mathbf{d}_r , and if the solution κ is ≤ 0 , we will set $\mathbf{d}_k = \mathbf{0}$.

Note that it is not uncommon to add a post-processing step to make atoms in \mathbf{D} to have unit norms or simply requiring $\|\mathbf{d}_r\|_2$ to be 1. This changes the iterate of \mathbf{d}_r to $\mathbf{d}_r = (t \mathbf{I} + \beta \Phi_r)^{-1} \nu_i$ as $\|\mathbf{d}_r\|_2 = 1$. Note that it is not uncommon to add a post-processing step to make atoms in \mathbf{D} to have unit norms or simply constraint $\|\mathbf{d}_r\|_2$ to be 1. This makes the whole algorithm much more efficient as computing eigen-decomposition of a typically large matrix Φ_r can now be avoided.

2. Supplementary Experimental Results

2.1. More results from face recognition

This section shows the supplementary results to Sec. 3.2 in the main text. Fig 1 shows the classification rates with dimensionality reduced to $m = \{100, \dots, 500\}$. The "λ" for each algorithm are similar to what are listed in the main text because we normalized the vectors after the PCA projection.

References

- [1] S. Bengio, F. Pereira, Y. Singer, and D. Strelow. Group sparse coding. *Advances in Neural Information Processing Systems*, 22:82–89, 2009.
- [2] E. Elhamifar and R. Vidal. Block-sparse recovery via convex optimization. *Signal Processing, IEEE Transactions on*, PP(99):1, 2012.

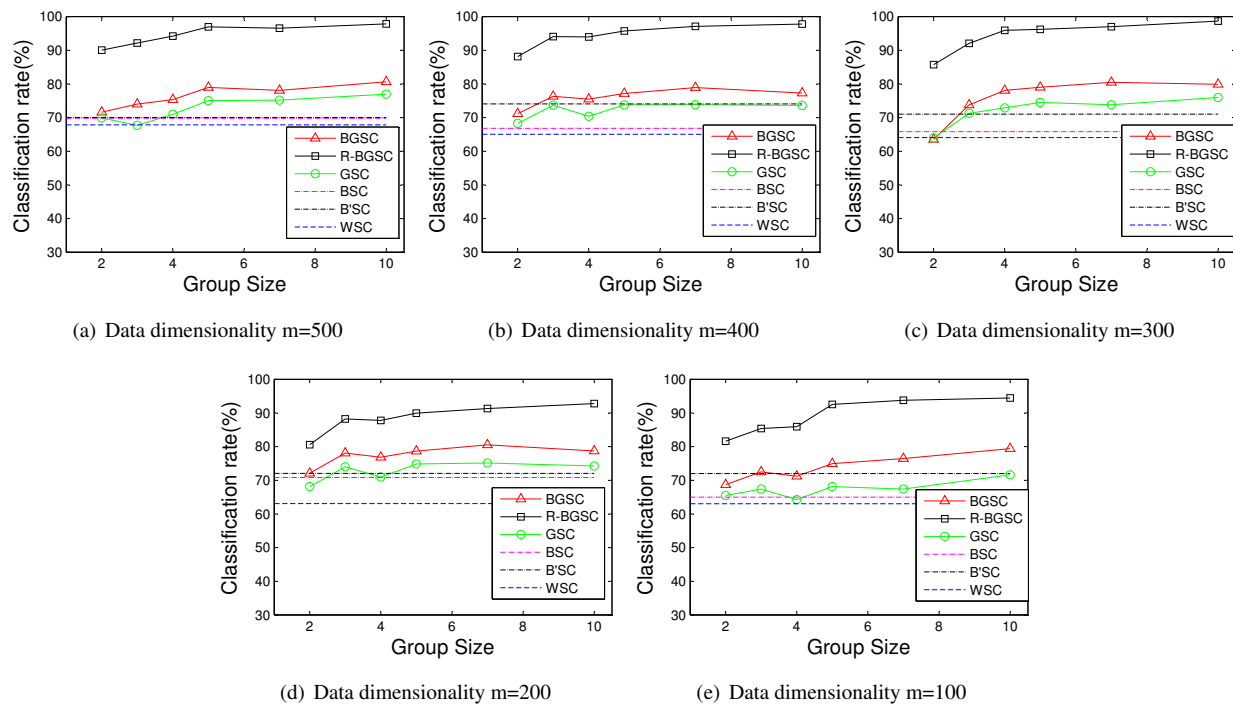


Figure 1: Classification rates by different algorithms under different data dimensionality reduction ($m = 500 \dots 100$).