

# Supplementary Material for *High-resolution Hyperspectral Imaging via Matrix Factorization*

November 18, 2010

In this supplementary note, we describe in more detail our approach to constructing a factorization  $\mathbf{Y} = \mathbf{A}\mathbf{X} \in \mathbb{R}^{m \times p}$  of the given observation into a sparsifying basis  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and a matrix  $\mathbf{X} \in \mathbb{R}^{n \times p}$  of sparse coefficients. We formulate this problem as a constrained optimization which seeks a basis  $\mathbf{A}$  which best sparsifies the observations, subject to the data constraint  $\mathbf{Y} = \mathbf{A}\mathbf{X}$ . We also impose the additional constraint that each of the columns of  $\mathbf{A}$  have unit  $\ell^2$ -norm, to eliminate scale ambiguities between the two unknowns  $\mathbf{A}$ ,  $\mathbf{X}$ . This leaves us with a manifold of possible solution pairs

$$\mathcal{M} = \{(\mathbf{A}, \mathbf{X}) \mid \mathbf{Y} = \mathbf{A}\mathbf{X}, \|\mathbf{A}\mathbf{e}_i\|_2 = 1 \forall i\} \subset \mathbb{R}^{m \times n} \times \mathbb{R}^{n \times p}. \quad (0.1)$$

Motivated by the success of  $\ell^1$ -minimization in recovering signals which are sparse in a known basis [BDE08], as well as emerging theoretical results on the good properties of  $\ell^1$ -minimization for sparse matrix factorization [GS10], relax the search for the best sparsifying basis  $\mathbf{A}$  to the following optimization problem:

$$\text{minimize } \|\mathbf{X}\|_1 \quad \text{subject to } (\mathbf{A}, \mathbf{X}) \in \mathcal{M}. \quad (0.2)$$

This problem asks that we minimize a convex function, subject to a nonlinear (nonconvex) constraint. In the worst case, this is a very difficult problem. However, surprisingly encouraging empirical results suggest that there are large classes of problems that actually can be solved globally via this approach. These problems correspond to instances where the input  $\mathbf{Y}$  is indeed generated by a sufficiently sparse  $\mathbf{X}$ , and the dimension is sufficiently large (say,  $m \geq 20$  for a 3-sparse  $\mathbf{X}$ ). Figure 1 (right) shows one example of this behavior: for well-structured simulated problems, locally minimizing the  $\ell^1$ -norm correctly recovers the globally optimal solution  $(\mathbf{A}, \mathbf{X})$  with high probability. We will report on this phenomenon more extensively in a separate, forthcoming work.

However, this still leaves us with a very challenging optimization problem: the number of variables is very large – potentially in the hundreds of thousands or even millions – and the objective function is nonsmooth. Our approach to this problem is a very simple and natural one, which has been studied in the optimization literature at least since the 1970’s, as a nonsmooth Gauss-Newton method [Cro78, JO80]. This approach repeatedly linearizes the constraint  $\mathbf{x} \doteq (\mathbf{A}, \mathbf{X}) \in \mathcal{M}$ , replacing  $\mathcal{M}$  with its tangent space at the current iterate. This yields a simple iteration:<sup>1</sup>

$$\boldsymbol{\delta}_k = \arg \min \{ \|\mathbf{X} + \boldsymbol{\Delta}_X\|_1 \mid (\boldsymbol{\Delta}_A, \boldsymbol{\Delta}_X) \in T_{\mathbf{x}_k} \mathcal{M}, \|\boldsymbol{\Delta}_A\|_F^2 + \|\boldsymbol{\Delta}_X\|_F^2 \leq \eta_k \} \quad (0.3)$$

$$\mathbf{x}_{k+1} = \mathcal{P}_{\mathcal{M}}[\mathbf{x}_k + \boldsymbol{\delta}_k], \quad (0.4)$$

---

<sup>1</sup>We have identified  $T_{\mathbf{x}_k} \mathcal{M}$  with a subspace of  $\mathbb{R}^{m \times n} \times \mathbb{R}^{n \times p}$  in the natural way.

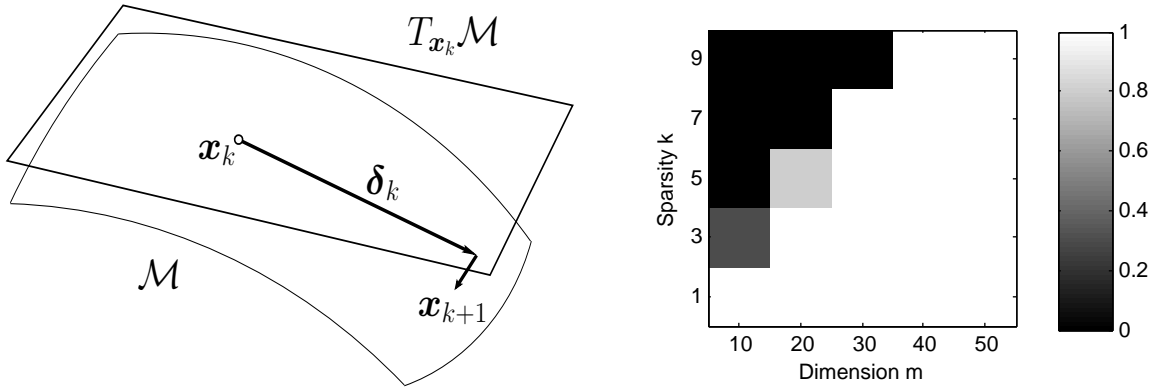


Figure 1: **Left: Illustration of our method.** At each iteration  $k$ , our algorithm linearizes the constraint manifold  $\mathcal{M}$ , and finds the perturbations  $\delta_k = (\Delta_A, \Delta_X)$  within this linearized space that make the objective  $\|\mathbf{X} + \Delta_X\|_1$  as small as possible. This minimizing point,  $\mathbf{x}_k + \delta_k$ , is then projected back onto  $\mathcal{M}$  to yield the next iterate,  $\mathbf{x}_{k+1}$ . **Right: Illustration of correct recovery.** We generate synthetic examples with randomly chosen  $\mathbf{A} \in \mathbb{R}^{m \times m}$ , for varying  $m = 10, 20, \dots, 50$ , and  $\mathbf{X} \in \mathbb{R}^{m \times p}$  of varying sparsity  $k = 1, 3, \dots, 9$ . Recovery is considered correct if the relative errors in  $\mathbf{A}$  and  $\mathbf{X}$  are both smaller than  $10^{-6}$ . The figure plots the fraction of correct recoveries over 10 independent trials - white corresponds to perfect recovery in all trials. We observe that when the dimension is large enough, the algorithm indeed correctly recovers  $(\mathbf{A}, \mathbf{X})$  with high probability.

where  $\mathcal{P}_{\mathcal{M}} : \mathbb{R}^{m \times n} \times \mathbb{R}^{n \times p}$  is a projection operator onto  $\mathcal{M}$ . The basic idea of this approach is illustrated in Figure 1. This algorithm is known to converge quadratically in the neighborhood of any strict local minimum [JO80]. We describe the two key steps of the algorithm in more details below.

**Linearized Subproblem:** Solving the linearized subproblem (0.3) amounts to solving a large, equality-constrained  $\ell^1$ -norm minimization problem. It is not difficult to show that the tangent space  $T_{\mathbf{x}}\mathcal{M}$  at a given point  $\mathbf{x} = (\mathbf{A}, \mathbf{X})$  can be expressed as the set of pairs  $(\Delta_A, \Delta_X)$  satisfying

$$\mathbf{A}\Delta_X + \Delta_A\mathbf{X} = 0, \quad \langle \mathbf{A}\mathbf{e}_i\Delta_A, \mathbf{e}_i \rangle = 0, \quad \forall i. \quad (0.5)$$

Hence, in more concrete form, the subproblem that we need to solve at each step is

$$\begin{aligned} \text{minimize } \|\mathbf{X} + \Delta_X\|_1 \quad \text{subject to} \quad & \mathbf{A}\Delta_X + \Delta_A\mathbf{X} = 0, \\ & \langle \mathbf{A}\mathbf{e}_i\Delta_A, \mathbf{e}_i \rangle = 0 \quad \forall i, \\ & \|\Delta_A\|_F^2 + \|\Delta_X\|_F^2 \leq \eta^2. \end{aligned} \quad (0.6)$$

For this purpose, we employ an Augmented Lagrange Multiplier (ALM) algorithm [Ber82], which was introduced into the literature on  $\ell^1$ -minimization by Yin, Osher and collaborators [YOGD08], and studied by those authors as a Bregman iterative algorithm. To describe this algorithm more fully, it is helpful to introduce a linear operator  $\Psi : \mathbb{R}^{m \times n} \times \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^{m \times p} \times \mathbb{R}^n$  corresponding to the above tangent space constraint:

$$\Psi(\Delta_A, \Delta_X) = (\mathbf{A}\Delta_X + \Delta_A\mathbf{X}, \text{diag}[\mathbf{A}^*\Delta_A]). \quad (0.7)$$

Introduce the convex function  $J : \mathbb{R}^{m \times n} \times \mathbb{R}^{n \times p}$  via

$$J(\Delta_A, \Delta_X) \doteq \begin{cases} \|\mathbf{X} + \Delta_X\|_1, & \|\Delta_X\|_F^2 + \|\Delta_A\|_F^2 \leq \eta^2 \\ \infty, & \text{else.} \end{cases} \quad (0.8)$$

The Bregman iterative approach solves a sequence of convex programs in unknowns  $\mathbf{h} = (\mathbf{H}_A, \mathbf{H}_X)$ , with varying inputs  $\mathbf{r}_j = (\mathbf{R}_A, \mathbf{R}_X)$ , starting from  $\mathbf{r}_0 = (\mathbf{0}, \mathbf{0})$ :

$$\mathbf{h}_{j+1} = \arg \min_{\mathbf{h}} J(\mathbf{h}) + \lambda \|\Psi(\mathbf{h}) - \mathbf{r}_j\|, \quad (0.9)$$

and updating the residual via  $\mathbf{r}_{j+1} = \mathbf{r}_j + \Psi(\mathbf{h}_{j+1}) - \Psi(\mathbf{h}_j)$ . This sequence of iterates converges to the solution to (0.6). Each of these problems (0.9) is solved via a simple soft-thresholding algorithm, which repeatedly sets

$$\mathbf{z}_{i+1} = \arg \min_{\mathbf{z}} J(\mathbf{z}) + \frac{\lambda}{2\beta} \|\mathbf{z} - \beta\Psi(\mathbf{z}_i - \mathbf{r}_j)\|. \quad (0.10)$$

This problem is a variant of shrinkage or soft-thresholding that arises in  $\ell^1$ -minimization. It can be solved via an efficient projection algorithm similar to the one given in [MBPS10].

**Projection:** We project a given pair  $\mathbf{A}, \mathbf{X}$  onto  $\mathcal{M}$  as follows. We first scale the columns of  $\mathbf{A}$  to obtain a matrix  $\mathbf{A}'$  whose columns have unit  $\ell^2$  norm. We then select  $\mathbf{X}'$  to be the matrix  $\tilde{\mathbf{X}}$  satisfying  $\mathbf{Y} = \mathbf{A}'\tilde{\mathbf{X}}$  which is closest to  $\mathbf{X}$  in Frobenius norm. In notation:

$$\mathbf{X}' = \mathbf{X} + (\mathbf{A}')^\dagger (\mathbf{Y} - \mathbf{A}'\mathbf{X}). \quad (0.11)$$

We then set  $\mathcal{P}(\mathbf{A}, \mathbf{X}) = (\mathbf{A}', \mathbf{X}')$ .

It should be noted that this operator does not find a pair  $(\mathbf{A}', \mathbf{X}')$  that is closest to  $(\mathbf{A}, \mathbf{X})$  in terms of the  $\ell^2$  norm  $\|(\mathbf{A}, \mathbf{X})\| = (\|\mathbf{A}\|_F^2 + \|\mathbf{X}\|_F^2)^{1/2}$ . Finding a closest point on  $\mathcal{M}$  in terms of this (rather natural) norm would require solving a separate nonconvex optimization problem. Rather, our operator  $\mathcal{P}(\mathbf{A}, \mathbf{X})$  can be viewed as the limit (as  $t \rightarrow \infty$ ) of projection operators with respect to weighted  $\ell^2$  norms  $\|(\mathbf{A}, \mathbf{X})\| = (t\|\mathbf{A}\|_F^2 + \|\mathbf{X}\|_F^2)^{1/2}$ . This substitution leaves us with a simple, tractable operator  $\mathcal{P}$ , and does not affect the qualitative behavior of the algorithm.

**Initialization:** We find it surprisingly effective to simply start the algorithm from a random initialization. We choose  $\mathbf{A}_0$  to be a matrix whose columns are independent random samples from the uniform distribution on the sphere  $\mathbb{S}^{m-1}$ . We then simply set  $\mathbf{X}_0 = \mathbf{A}_0^\dagger \mathbf{Y}$ .

## References

- [BDE08] A.M. Bruckstein, D. L. Donoho, and M. Elad. From sparse solutions of systems of equations to sparse modeling of signals and images. *To appear in SIAM Review*, 2008.
- [Ber82] D. Bertsekas. *Constrained Optimization and Lagrange Multiplier Methods*. Academic Press, 1982.
- [Cro78] L. Cromme. Strong uniqueness: A far-reaching criterion for the convergence analysis of iterative procedures. *Numerische Mathematik*, 29:179–193, 1978.

- [GS10] R. Gribonval and K. Schnass. Dictionary identification - sparse matrix factorization via  $\ell_1$ -minimization. <http://arxiv.org/abs/0904.4774>, 2010.
- [JO80] K. Jittorntrum and M. Osborne. Strong uniqueness and second order convergence in nonlinear discrete approximation. *Numerische Mathematik*, 34:439–455, 1980.
- [MBPS10] J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online learning for matrix factorization and sparse coding. *Journal of Machine Learning Research*, 11:19–60, 2010.
- [YOGD08] W. Yin, S. Osher, D. Goldfarb, and J. Darbon. Bregman iterative algorithms for  $\ell^1$ -minimization with applications to compressed sensing. *SIAM Journal on Imaging Sciences (SJIS)*, 2008.