

Block-Simultaneous Direction Method of Multipliers

A proximal primal-dual splitting algorithm for nonconvex problems with multiple constraints

Fred Moolekamp · Peter Melchior

the date of receipt and acceptance should be inserted later

Abstract We introduce a generalization of the linearized Alternating Direction Method of Multipliers to optimize a real-valued function f of multiple arguments with potentially multiple constraints g_\circ on each of them. The function f may be nonconvex as long as it is convex in every argument, while the constraints g_\circ need to be convex but not smooth. If f is smooth, the proposed Block-Simultaneous Direction Method of Multipliers (bSDMM) can be interpreted as a proximal analog to inexact coordinate descent methods under constraints. Unlike alternative approaches for joint solvers of multiple-constraint problems, we do not require linear operators L of a constraint function $g(L \cdot)$ to be invertible or linked between each other. bSDMM is well-suited for a range of optimization problems, in particular for data analysis, where f is the likelihood function of a model and L could be a transformation matrix describing e.g. finite differences or basis transforms. We apply bSDMM to the Non-negative Matrix Factorization task of a hyperspectral unmixing problem and demonstrate convergence and effectiveness of multiple constraints on both matrix factors. The algorithms are implemented in python and released as an open-source package.

Keywords Optimization · Proximal Algorithms · Nonconvex Optimization · Block Coordinate Descent · Non-negative Matrix Factorization

1 Introduction

In this paper we seek to numerically

$$\text{minimize } f(\mathbf{x}_1, \dots, \mathbf{x}_N) + \sum_{j=1}^N \sum_{i=1}^{M_j} g_{ij}(L_{ij}\mathbf{x}_j), \quad (1)$$

Fred Moolekamp
Department of Astrophysical Sciences, Princeton University, Princeton, NJ 08544, USA
E-mail: fredem@princeton.edu

Peter Melchior
Department of Astrophysical Sciences, Princeton University, Princeton, NJ 08544, USA
E-mail: peter.melchior@princeton.edu

where $f: \mathbb{R}^{d_1} \times \dots \times \mathbb{R}^{d_N} \rightarrow \mathbb{R}$ is a potentially nonconvex function that is a closed proper convex function in each of N independent arguments, and $g_{ij}: \mathbb{R}^{d_{ij}} \rightarrow \mathbb{R}$ are convex functions that encode M_j constraints for each of those variables after they are mapped by linear operators L_{ij} .¹

If f and g_{\circ} were all smooth, we could directly apply conventional gradient methods to solve the problem, but in many cases this is not possible. Examples are projections onto the positive orthant of \mathbb{R}^{d_j} or regularization with the ℓ_1 norm. Instead we will access the functions through their proximal operators, defined as

$$\text{prox}_{\lambda f}(\mathbf{v}) \equiv \underset{\mathbf{x}}{\text{argmin}} \left\{ f(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{v}\|_2^2 \right\} \quad (2)$$

with a scaling parameter λ . Their primary purpose is to turn any convex function f into a strongly convex function even if f may assume infinite values. The minimization in [Equation 2](#) may seem complicated, but in many cases exact and simple solvers exist. For instance, if $f(\mathbf{x})$ is the indicator function of the closed convex set \mathcal{C} , prox_f is simply the Euclidean projection operator onto \mathcal{C} . Whether [Equation 1](#) can efficiently be solved thus depends on the cost of evaluating the proximal operators involved. For more details and various interpretations of proximal operators we refer to [Parikh et al \(2014\)](#) and references therein.

Several proximal algorithms exist to minimize a function f of a single argument. If f is smooth, the proximal gradient method provides minimization of $f(\mathbf{x}) + g(\mathbf{x})$ with the *forward-backward* scheme, where at iteration k a step in the direction of ∇f is followed by the application of prox_g :

$$\mathbf{x}^{k+1} := \text{prox}_{\lambda^k g} \left(\mathbf{x}^k - \lambda^k \nabla f(\mathbf{x}^k) \right). \quad (3)$$

If the step size is $\lambda \in (0, 1/L]$ with L being the Lipschitz constant of ∇f , the convergence rate is $O(1/k)$, which can further be accelerated ([Nesterov 2013](#)). We are particularly interested in introducing linear operators, which encode typical image processing operations, e.g. finite differences, smoothing, or basis transforms. A well-known approach to such a situation is the Alternating Direction Method of Multipliers (ADMM), which we review in [Section 2.1](#); several such constraints can be imposed with the Simultaneous Direction Method of Multipliers (SDMM, [Section 2.2](#)). We will introduce previous works in the relevant sections and state here only our main contributions: We 1) generalize SDMM to cases where none of the L_{ij} need to have full rank or be linked between different constraints, and 2) extend SDMM to nonconvex functions f of several arguments \mathbf{x}_j , resulting in a proximal variant of inexact block optimization methods ([Section 2.3](#)). As an example of the proposed algorithm, in [Section 3](#) we implement a solver for the Non-negative Matrix Factorization problem, where we allow for an arbitrary number of constraints on either of the matrix factors. We demonstrate convergence of the primal and dual variables as well as the effectiveness of multiple constraints in the example case. We conclude in [Section 4](#).

¹ Throughout this work, indices denote different variables or constraints, not elements of vectors or tensors.

2 Generalizing ADMM

2.1 Linearized Alternating Direction Method of Multipliers

We start by introducing the well-known Alternating Direction Method of Multipliers (ADMM, [Gabay and Mercier 1976](#); [Glowinski and Marroco 1975](#); [Eckstein and Bertsekas 1992](#)) in the notation that will be used throughout the paper. It is applicable when $N = 1$, i.e. to

$$\underset{\mathbf{x}_1}{\text{minimize}} f(\mathbf{x}_1) + g_1(\mathbf{L}_{11}\mathbf{x}_1). \quad (4)$$

We first re-write [Equation 4](#) in consensus form as

$$\begin{aligned} &\underset{\mathbf{z}_{11}}{\text{minimize}} f(\mathbf{x}_1) + g_1(\mathbf{z}_{11}) \\ &\text{subject to } \mathbf{L}_{11}\mathbf{x}_1 - \mathbf{z}_{11} = 0. \end{aligned} \quad (5)$$

The central idea is to split the optimization in two separate tasks: one that minimizes f and another that satisfies g_1 by introducing the auxiliary variable \mathbf{z}_{11} ([Douglas and Rachford 1956](#)). This can be done by introducing Lagrange multipliers for each constraint, plus a quadratic term that is of critical importance for the convergence of the algorithm when g_1 is not strongly convex. The resulting Augmented Lagrangian for [Equation 4](#) is

$$\mathcal{L}(\mathbf{x}_1, \mathbf{z}_{11}, \boldsymbol{\lambda}_{11}) = f(\mathbf{x}_1) + g_1(\mathbf{z}_{11}) + \boldsymbol{\lambda}_{11}^\top (\mathbf{L}_{11}\mathbf{x}_1 - \mathbf{z}_{11}) + \frac{1}{2\rho_{11}} \|\mathbf{L}_{11}\mathbf{x}_1 - \mathbf{z}_{11}\|_2^2, \quad (6)$$

where $\rho_{11} \in \mathbb{R} > 0$ and $\boldsymbol{\lambda}_{11} \in \mathbb{R}^{d_1}$. We then need to find unique minimizers of \mathcal{L} with respect to its variables. In an iterative sequence one would formally need to update both \mathbf{x}_1 and \mathbf{z}_{11} simultaneously, which is often numerically difficult. [Chen and Teboulle \(1994\)](#) demonstrated that it is sufficient to update \mathbf{x}_1 first, then use the updated value of \mathbf{x}_1 for the \mathbf{z}_{11} -update, and then both values for the last update of $\boldsymbol{\lambda}_{11}$:

$$\begin{aligned} \mathbf{x}_1^{k+1} &:= \underset{\mathbf{x}_1}{\text{argmin}} \left\{ f(\mathbf{x}_1) + \boldsymbol{\lambda}_{11}^{k\top} \mathbf{L}_{11}\mathbf{x}_1 + \frac{1}{2\rho_{11}} \|\mathbf{L}_{11}\mathbf{x}_1 - \mathbf{z}_{11}^k\|_2^2 \right\} \\ \mathbf{z}_{11}^{k+1} &:= \underset{\mathbf{z}_{11}}{\text{argmin}} \left\{ g_1(\mathbf{z}_{11}) - \boldsymbol{\lambda}_{11}^{k\top} \mathbf{z}_{11} + \frac{1}{2\rho_{11}} \|\mathbf{L}_{11}\mathbf{x}_1^{k+1} - \mathbf{z}_{11}\|_2^2 \right\} \\ \boldsymbol{\lambda}_{11}^{k+1} &:= \boldsymbol{\lambda}_{11}^k + \frac{1}{\rho_{11}} (\mathbf{L}_{11}\mathbf{x}_1^{k+1} - \mathbf{z}_{11}^{k+1}). \end{aligned} \quad (7)$$

We can simplify the \mathbf{z}_{11} -update because we can add or subtract any terms independent of \mathbf{z}_{11} :

$$\begin{aligned} \mathbf{z}_{11}^{k+1} &:= \underset{\mathbf{z}_{11}}{\text{argmin}} \left\{ g_1(\mathbf{z}_{11}) - \frac{1}{2\rho_{11}} \|\mathbf{L}_{11}\mathbf{x}_1^{k+1} + \mathbf{u}_{11}^k - \mathbf{z}_{11}\|_2^2 \right\} = \text{prox}_{\rho_{11}g_1}(\mathbf{L}_{11}\mathbf{x}_1^{k+1} + \mathbf{u}_{11}^k) \\ \mathbf{u}_{11}^{k+1} &:= \mathbf{u}_{11}^k + \mathbf{L}_{11}\mathbf{x}_1^{k+1} - \mathbf{z}_{11}^{k+1} \end{aligned} \quad (8)$$

where we have introduced a scaled variable $\mathbf{u}_{11}^k \equiv \rho_{11}\boldsymbol{\lambda}_{11}^k$ and applied [Equation 2](#).

One cannot solve for \mathbf{x}_1^{k+1} in the same way because of the presence of the linear operator \mathbf{L}_{11} . [Stephanopoulos and Westerberg \(1975\)](#) showed that linearizing $\frac{1}{2\rho_{11}} \|\mathbf{L}_{11}\mathbf{x}_1 - \mathbf{z}_{11}^k\|_2^2$ at the current-iteration \mathbf{x}_1^k is a practical solution that preserves the general convergence of the algorithm while providing a separable update sequence:

$$\begin{aligned} \mathbf{x}_1^{k+1} &:= \underset{\mathbf{x}_1}{\text{argmin}} \left\{ f(\mathbf{x}_1) + \boldsymbol{\lambda}_{11}^{k\top} \mathbf{L}_{11}\mathbf{x}_1 + \frac{1}{\rho_{11}} \mathbf{L}_{11}^\top (\mathbf{L}_{11}\mathbf{x}_1^k - \mathbf{z}_{11}^k) \mathbf{x}_1 + \frac{1}{2\mu_1} \|\mathbf{x}_1 - \mathbf{x}_1^k\|_2^2 \right\} \\ &= \text{prox}_{\mu_1 f} \left(\mathbf{x}_1^k - \frac{\mu_1}{\rho_{11}} \mathbf{L}_{11}^\top (\mathbf{L}_{11}\mathbf{x}_1^k - \mathbf{z}_{11}^k + \mathbf{u}_{11}^k) \right), \end{aligned} \quad (9)$$

Algorithm 1 Linearized ADMM

```

1: procedure ADMM( $\mathbf{x}_1, \mu_1, \rho_{11}, \mathbf{L}_{11}$ )
2:    $\mathbf{x}_1^1 \leftarrow \mathbf{x}_1; \mathbf{z}_{11}^1 \leftarrow \mathbf{L}_{11}\mathbf{x}_1; \mathbf{u}_{11}^1 \leftarrow \mathbf{0}$ 
3:   for  $k=1, 2, \dots$  do
4:      $\mathbf{x}_1^{k+1} \leftarrow \text{prox}_{\mu_1 f}(\mathbf{x}_1^k - \mu_1/\rho_{11}\mathbf{L}_{11}^\top(\mathbf{L}_{11}\mathbf{x}_1^k - \mathbf{z}_{11}^k + \mathbf{u}_{11}^k))$ 
5:      $\mathbf{z}_{11}^{k+1} \leftarrow \text{prox}_{\rho_{11}g_1}(\mathbf{L}_{11}\mathbf{x}_1^{k+1} + \mathbf{u}_{11}^k)$ 
6:      $\mathbf{u}_{11}^{k+1} \leftarrow \mathbf{u}_{11}^k + \mathbf{L}_{11}\mathbf{x}_1^{k+1} - \mathbf{z}_{11}^{k+1}$ 
7:     if  $\|\mathbf{r}_{11}^{k+1}\|_2 \leq \epsilon^{\text{pri}} \wedge \|\mathbf{s}_{11}^{k+1}\|_2 \leq \epsilon^{\text{dual}}$  then break

```

where we have introduced the parameter μ_1 with $0 < \mu_1 \leq \rho_{11}/\|\mathbf{L}_{11}\|_s^2$.² The algorithm for a single variable \mathbf{x}_1 and a single constraint $g_1(\mathbf{L}_{11}\mathbf{x}_1)$ is also known as *split inexact Uzawa method* (e.g. [Esser et al 2010](#); [Parikh et al 2014](#)) and listed as [Algorithm 1](#).

Following [Boyd et al \(2011\)](#), we implement stopping criteria based on primal residual $\mathbf{r}_{11}^{k+1} = \mathbf{L}_{11}\mathbf{x}_1^{k+1} - \mathbf{z}_{11}^{k+1}$ and the dual residual $\mathbf{s}_{11}^{k+1} = \frac{1}{\rho_{11}}\mathbf{L}_{11}^\top(\mathbf{z}_{11}^{k+1} - \mathbf{z}_{11}^k)$. To assess primal and dual feasibility, we require

$$\begin{aligned} \|\mathbf{r}_{11}^{k+1}\|_2 \leq \epsilon^{\text{pri}} &\equiv \sqrt{p}\epsilon^{\text{abs}} + \epsilon^{\text{rel}} \max\{\|\mathbf{L}_{11}\mathbf{x}_1^{k+1}\|_2, \|\mathbf{z}_{11}^{k+1}\|_2\} \quad \text{and} \\ \|\mathbf{s}_{11}^{k+1}\|_2 \leq \epsilon^{\text{dual}} &\equiv \sqrt{n}\epsilon^{\text{abs}} + \epsilon^{\text{rel}}/\rho_{11}\|\mathbf{L}_{11}^\top\mathbf{u}_{11}^{k+1}\|_2, \end{aligned} \quad (10)$$

where p and n are the number of elements in \mathbf{z}_{11} and \mathbf{x}_1 , respectively. The error thresholds ϵ^{abs} and ϵ^{rel} can be set at suitable values, depending on the precision and runtime constraints of the application.

2.2 Several constraint functions

It is often necessary to impose several constraints simultaneously. [Condat \(2013\)](#) proposed primal-dual split algorithms that can solve a restricted version of [Equation 1](#), in which $f(\mathbf{x})$ is convex and differentiable with a Lipschitz-continuous gradient, and only two additional, potentially non-smooth functions are present, $g(\mathbf{x})$ and $h(\mathbf{L}\mathbf{x})$. More generally applicable, [Combettes and Pesquet \(2011\)](#) introduced the Simultaneous Direction Method of Multipliers (SDMM) to

$$\underset{\mathbf{x}_1}{\text{minimize}} \sum_{i=1}^{M_1} g_i(\mathbf{L}_{i1}\mathbf{x}_1), \quad (11)$$

for which g_i only have to be convex. This corresponds to $N = 1$ and $f(\mathbf{x}_1) = g_1(\mathbf{x}_1)$ with $\mathbf{L}_{11} = \mathbf{I}$ in our notation, and enables the adoption of an arbitrary number of constraint functions, which is expressly what we seek. However, their algorithm requires that $\mathbf{Q} = \sum_{i=1}^{M_1} \mathbf{L}_{i1}^\top \mathbf{L}_{i1}$ be invertible, a limitation that is too restrictive for many problems of interest. We can dispense with this requirement by adopting the same linearization strategy as before with the ADMM. Such a strategy is sensible if at least one function g_l in [Equation 11](#) can act as f in [Equation 1](#).³

² We use $\|\cdot\|_s$ to denote the spectral norm, $\|\cdot\|_2$ for the element-wise ℓ_2 norm of vectors and tensors.

³ While it is always possible to reformulate the problem thusly because we can set $f(\mathbf{x}_1) = g_l(\mathbf{L}_{j1}\mathbf{x}_1)$ for any l , it may render inefficient the minimization of f by means of an proximal operator. This is the limitation of the algorithm we derive in this section.

Algorithm 2 Linearized SDMM

```

1: procedure SDMM( $\mathbf{x}_1, \mu_1, [\rho_{11} \dots, \rho_{M_1 1}], [L_{11} \dots, L_{M_1 1}]$ )
2:    $\mathbf{x}_1^1 \leftarrow \mathbf{x}_1; \mathbf{z}_{i1}^1 \leftarrow L_{i1} \mathbf{x}_1 \forall i \in \{1, \dots, M_1\}; \mathbf{u}_{i1}^1 \leftarrow \mathbf{0} \forall i \in \{1, \dots, M_1\}$ 
3:   for  $k = 1, 2, \dots$  do
4:      $\mathbf{x}_1^{k+1} \leftarrow \text{prox}_{\mu_1 f} \left( \mathbf{x}_1^k - \sum_{i=1}^{M_1} \mu_1 / \rho_{i1} L_{i1}^\top (L_{i1} \mathbf{x}_1^k - \mathbf{z}_{i1}^k + \mathbf{u}_{i1}^k) \right)$ 
5:     for  $i = 1, \dots, M_1$  do
6:        $\mathbf{z}_{i1}^{k+1} \leftarrow \text{prox}_{\rho_{i1} g_{i1}} (L_{i1} \mathbf{x}_1^{k+1} + \mathbf{u}_{i1}^k)$ 
7:        $\mathbf{u}_{i1}^{k+1} \leftarrow \mathbf{u}_{i1}^k + L_{i1} \mathbf{x}_1^{k+1} - \mathbf{z}_{i1}^{k+1}$ 
8:     if  $\bigwedge_i \{ \|\mathbf{r}_{i1}^{k+1}\|_2 \leq \varepsilon^{\text{pri}} \wedge \|\mathbf{s}_{i1}^{k+1}\|_2 \leq \varepsilon^{\text{dual}} \}$  then break

```

Without loss of generality, we take $l = M_1$ and redefine $M_1 \rightarrow M_1 - 1$ in Equation 11 to maintain our notation. We introduce M_1 primal variables \mathbf{z}_{i1} and solve the problem in consensus form:

$$\begin{aligned} & \underset{\mathbf{x}_1}{\text{minimize}} \quad f(\mathbf{x}_1) + \sum_{i=1}^{M_1} g_i(\mathbf{z}_{i1}) \\ & \text{subject to} \quad L_{i1} \mathbf{x}_1 - \mathbf{z}_{i1} = \mathbf{0} \quad \forall i \in \{1, \dots, M_1\}. \end{aligned} \quad (12)$$

Updating the primal and dual variables \mathbf{u}_{i1} is exactly the same as in Equation 8,

$$\begin{aligned} \mathbf{z}_{i1}^{k+1} &:= \text{prox}_{\rho_{i1} g_{i1}} (L_{i1} \mathbf{x}_1^{k+1} + \mathbf{u}_{i1}^k) \\ \mathbf{u}_{i1}^{k+1} &:= \mathbf{u}_{i1}^k + L_{i1} \mathbf{x}_1^{k+1} - \mathbf{z}_{i1}^{k+1}, \end{aligned} \quad (13)$$

where each update can be performed in parallel, justifying the ‘‘S’’ in SDMM. Because we isolated f as a function of \mathbf{x}_1 alone, we can minimize $\mathcal{L}(\mathbf{x}_1, \mathbf{z}_{11}, \dots, \mathbf{z}_{M_1 1}, \boldsymbol{\lambda}_{11}, \dots, \boldsymbol{\lambda}_{M_1 1})$ with respect to \mathbf{x}_1 with the update

$$\begin{aligned} \mathbf{x}_1^{k+1} &:= \underset{\mathbf{x}_1}{\text{argmin}} \left\{ f(\mathbf{x}_1) + \sum_{i=1}^{M_1} \left(\boldsymbol{\lambda}_{i1}^{k\top} L_{i1} \mathbf{x}_1 + \frac{1}{2\rho_{i1}} \|L_{i1} \mathbf{x}_1 - \mathbf{z}_{i1}^k\|_2^2 \right) \right\} \\ &= \underset{\mathbf{x}_1}{\text{argmin}} \left\{ f(\mathbf{x}_1) + \sum_{i=1}^{M_1} \left(\boldsymbol{\lambda}_{i1}^{k\top} L_{i1} \mathbf{x}_1 + \frac{1}{\rho_{i1}} L_{i1}^\top (L_{i1} \mathbf{x}_1^k - \mathbf{z}_{i1}^k) \mathbf{x}_1 + \frac{1}{2\mu_1} \|\mathbf{x}_1 - \mathbf{x}_1^k\|_2^2 \right) \right\} \\ &= \text{prox}_{\mu_1 f} \left(\mathbf{x}_1^k - \sum_{i=1}^{M_1} \mu_1 / \rho_{i1} L_{i1}^\top (L_{i1} \mathbf{x}_1^k - \mathbf{z}_{i1}^k + \mathbf{u}_{i1}^k) \right), \end{aligned} \quad (14)$$

where we linearized in the second step and added a quadratic penalty to introduce $\text{prox}_{\mu_1 f}$ in the third step. The parameters μ_1 and ρ_{i1} are bound by

$$\rho_{i1} / \|L_{i1}\|_s^2 \geq \beta_S \mu_1 \quad \text{with } 1 \leq \beta_S \leq M_1. \quad (15)$$

The parameter β_S is necessary to account for potentially correlated contributions of different g_i in Equation 14. If the g_i are partially degenerate and one would adopt the naive threshold $\beta_S = 1$, $\text{prox}_{\mu_1 f}$ will lose its contracting property. The most conservative option $\beta_S = M_1$ will always lead to a convergent minimizer, even if $g_i = g \forall i$, albeit generally at the expense of reduced convergence speeds.

The linearized form of SDMM is listed in Algorithm 2. We note that the resulting algorithm is not identical and thus not suited to the same set of problems as the original SDMM by Combettes and Pesquet (2011) because of the isolation of $f(\mathbf{x}_1)$ in Equation 12 that is not present in Equation 11. However, a strong similarity persists with the exception of Line 4 in Algorithm 2, we thus consider it appropriate to call this algorithm Linearized SDMM.

2.3 Nonconvex, multi-argument functions

We now seek to solve [Equation 1](#), which includes the treatment of the function f having several arguments. Even with our requirement that f be a closed proper convex function in each of its arguments, f itself is generally not convex (we refer the reader to e.g. [Zhang et al 2016](#) for cases when a constraint function g is not convex).

The nonconvexity that arises in multi-argument functions has been addressed with an ADMM variant first by [Hong et al \(2016\)](#), who provide a provably convergent solution for

$$\begin{aligned} & \underset{\mathbf{x}_1, \dots, \mathbf{x}_N}{\text{minimize}} && f(\mathbf{x}_1, \dots, \mathbf{x}_N) + \sum_{j=1}^N g_j(\mathbf{x}_j) \\ & \text{subject to} && \sum_j^N \mathbf{L}_j \mathbf{x}_j = \mathbf{b}, \end{aligned} \quad (16)$$

and later by [Wang et al \(2015\)](#), who solve

$$\begin{aligned} & \underset{\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{y}}{\text{minimize}} && f(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{y}) \\ & \text{subject to} && \sum_j^N \mathbf{L}_j \mathbf{x}_j + \mathbf{B} \mathbf{y} = \mathbf{b}. \end{aligned} \quad (17)$$

While close to our problem in [Equation 1](#), especially the form of [Equation 16](#), these approaches need the linking $\sum_j^N \mathbf{L}_j \mathbf{x}_j = \mathbf{b}$ across the variables \mathbf{x}_j (and \mathbf{y} for [Wang et al 2015](#)), which conflicts with our desire to have independent constraint functions. In addition, inequality constraints like $\mathbf{L} \mathbf{x} \leq 0$ cannot be expressed in either of the forms above.

By now our strategy for solving [Equation 1](#) with N variables \mathbf{x}_j should be apparent. We first bring the problem into consensus form, i.e. we seek to

$$\begin{aligned} & \underset{\mathbf{x}_1, \dots, \mathbf{x}_N}{\text{minimize}} && f(\mathbf{x}_1, \dots, \mathbf{x}_N) + \sum_{j=1}^N \sum_{i=1}^{M_j} g_{ij}(\mathbf{z}_{ij}) \\ & \text{subject to} && \mathbf{L}_{ij} \mathbf{x}_j - \mathbf{z}_{ij} = 0 \quad \forall j \in \{1, \dots, N\}, i \in \{1, \dots, M_j\}. \end{aligned} \quad (18)$$

This results in $\sum_{j=1}^N M_j$ primal and dual variables that are updated using

$$\begin{aligned} \mathbf{z}_{ij}^{k+1} &:= \text{prox}_{\rho_{ij} g_{ij}} \left(\mathbf{L}_{ij} \mathbf{x}_j^{k+1} + \mathbf{u}_{ij}^k \right) \\ \mathbf{u}_{ij}^{k+1} &:= \mathbf{u}_{ij}^k + \mathbf{L}_{ij} \mathbf{x}_j^{k+1} - \mathbf{z}_{ij}^{k+1}. \end{aligned} \quad (19)$$

Then we linearize the quadratic term $\frac{1}{2\rho_{ij}} \|\mathbf{L}_{ij} \mathbf{x}_j - \mathbf{z}_{ij}^k\|_2^2$, and utilize that f is convex in every argument to solve for each \mathbf{x}_j with a proximal operator as minimizer of f wrt \mathbf{x}_j :

$$\begin{aligned} \mathbf{x}_j^{k+1} &:= \underset{\mathbf{x}_j}{\text{argmin}} \left\{ f(\mathbf{x}_1, \dots, \mathbf{x}_N) + \sum_{i=1}^{M_j} \left(\boldsymbol{\lambda}_{ij}^{k\top} \mathbf{L}_{ij} \mathbf{x}_j + \frac{1}{2\rho_{ij}} \|\mathbf{L}_{ij} \mathbf{x}_j - \mathbf{z}_{ij}^k\|_2^2 \right) \right\} \\ &= \underset{\mathbf{x}_j}{\text{argmin}} \left\{ f(\mathbf{x}_1, \dots, \mathbf{x}_N) + \sum_{i=1}^{M_j} \left(\boldsymbol{\lambda}_{ij}^{k\top} \mathbf{L}_{ij} \mathbf{x}_j + \frac{1}{\rho_{ij}} \mathbf{L}_{ij}^\top (\mathbf{L}_{ij} \mathbf{x}_j - \mathbf{z}_{ij}^k) \mathbf{x}_j + \frac{1}{2\rho_{ij}} \|\mathbf{x}_j - \mathbf{x}_j^k\|_2^2 \right) \right\} \\ &= \text{prox}_{\mu_j f, j} \left(\mathbf{x}_j^k - \sum_{i=1}^{M_j} \frac{\mu_j}{\rho_{ij}} \mathbf{L}_{ij}^\top (\mathbf{L}_{ij} \mathbf{x}_j^k - \mathbf{z}_{ij}^k + \mathbf{u}_{ij}^k) \right), \end{aligned} \quad (20)$$

The entire algorithm, which we call bSDMM, is listed as [Algorithm 3](#). If f is separable, bSDMM amounts to N independent, and thus parallelizable, solutions for each \mathbf{x}_j using [Algorithm 2](#). If not, to the knowledge of the authors, general convergence guarantees do not exist since the solutions for each \mathbf{x}_j may not be unique. However, in the case of quadratic convex functions and of $N = 2$ convergence can be guaranteed ([Grippio and Sciandrone 2000](#); [Lin 2007](#)). If f is smooth, block-wise optimization by means of successive proximal forward-backward steps is convergent even with Nesterov-type acceleration ([Razaviyayn et al 2013](#); [Xu and Yin 2013](#)). In that case, these optimization steps constitute a proximal variant of a

Algorithm 3 Block-SDMM

```

1: procedure BSDMM( $[\mathbf{x}_1, \dots, \mathbf{x}_N]$ ,  $h$ ,  $\beta_G$ ,  $[\mathbf{L}_{11}, \dots, \mathbf{L}_{M_N N}]$ )
2:   for  $j = 1, \dots, N$  do
3:      $\mathbf{x}_j^1 \leftarrow \mathbf{x}_j$ 
4:      $\mathbf{z}_{ij}^1 \leftarrow \mathbf{L}_{ij} \mathbf{x}_j \forall i \in \{1, \dots, M_j\}$ 
5:      $\mathbf{u}_{ij}^1 \leftarrow \mathbf{0} \forall i \in \{1, \dots, M_j\}$ 
6:   for  $k = 1, 2, \dots$  do
7:     for  $j = 1, \dots, N$  do
8:        $\mu_j^{k+1} \leftarrow h(j; \mathbf{x}_1^k, \dots, \mathbf{x}_N^k)$ 
9:        $\rho_{ij}^{k+1} \leftarrow \beta_G \mu_j^{k+1} \|\mathbf{L}_{ij}\|_2^2$ 
10:       $\mathbf{x}_j^{k+1} \leftarrow \text{prox}_{\mu_j f, j} \left( \mathbf{x}_j^k - \sum_{i=1}^{M_j} \mu_j / \rho_{ij} \mathbf{L}_{ij}^\top (\mathbf{L}_{ij} \mathbf{x}_j^k - \mathbf{z}_{ij}^k + \mathbf{u}_{ij}^k) \right)$ 
11:      for  $i = 1, \dots, M_j$  do
12:         $\mathbf{z}_{ij}^{k+1} \leftarrow \text{prox}_{\rho_{ij} g_{ij}} (\mathbf{L}_{ij} \mathbf{x}_j^{k+1} + \mathbf{u}_{ij}^k)$ 
13:         $\mathbf{u}_{ij}^{k+1} \leftarrow \mathbf{u}_{ij}^k + \mathbf{L}_{ij} \mathbf{x}_j^{k+1} - \mathbf{z}_{ij}^{k+1}$ 
14:      if  $\bigwedge_{ij} \left\{ \|\mathbf{r}_{ij}^{k+1}\|_2 \leq \epsilon^{\text{pri}} \wedge \|\mathbf{s}_{ij}^{k+1}\|_2 \leq \epsilon^{\text{dual}} \right\}$  then break

```

block coordinate descent algorithm. Our approach is different because it uses a primal-dual split instead of a direct constraint projection, which can only deal with one constraint per optimization variable, and the nonconvex ADMM variants of [Hong et al \(2016\)](#) and [Wang et al \(2015\)](#) because of the independence of the constraints.

Furthermore, because the updates do not necessarily yield the minimum of f or the constraints in each step, e.g. for a single step of $\nabla_j f$, it is very efficient, but the limit point of the sequence is not necessarily a local minimum, just a stationary point. As long as the approximations are sufficiently precise, ADMM is still convergent ([Eckstein and Bertsekas 1992](#); [Eckstein and Yao 2017](#)). [Berry et al \(2007\)](#) studied approximate solvers in several different applications and found acceptable results at a fraction of the computational cost. We will inspect the convergence properties of the bSDMM algorithm with an suitable example in the next section.

As with the SDMM, the presence of several constraints g_{ij} for a single function f necessitates restraint when choosing ρ_{ij} so as not to overwhelm $\text{prox}_{\mu_j f, j}$:

$$\rho_{ij} / \|\mathbf{L}_{ij}\|_2^2 \geq \beta_{\text{BS}} \mu_j \text{ with } 1 \leq \beta_{\text{BS}} \leq NM_j. \quad (21)$$

However, because f is now a function with several arguments, the parameter μ_j may change with every iteration. For instance, if f is smooth, μ_j is bound by the Lipschitz constant of $\nabla_j f(\mathbf{x}_1, \dots, \mathbf{x}_N)$. The bSDMM algorithm therefore requires the function $h(j; \mathbf{x}_1, \dots, \mathbf{x}_N)$ to compute μ_j , from which it will then determine ρ_{ij} to satisfy [Equation 21](#) (lines 8 and 9 of [Algorithm 3](#)).

3 Non-negative Matrix Factorization

An important application of bSDMM is Non-negative Matrix Factorization (NMF, [Paatero and Tapper 1994](#)), which seeks to describe a data set $D \in \mathbb{R}^{B \times L}$ of L -dimensional features that are observed B times as a product of two non-negative matrices $A \in \mathbb{R}^{B \times K}$ and $S \in \mathbb{R}^{K \times L}$. The idea is to reduce the dimensionality of the problem to K prototypes, encoded in S , whose sum generates the data in each observations with relative amplitudes encoded in A .

We will adopt the Euclidean cost function, which corresponds to the negative log-likelihood under the assumption of standard Gaussian errors on each element of D (see [Blanton and Roweis \(2007\)](#); [Zhu \(2016\)](#) for extensions to heteroscedastic errors). The objective function is thus

$$f(A, S) = \|A \cdot S - D\|_2^2, \quad (22)$$

and the non-negative constraints can be expressed as

$$g_+(A) + g_+(S) \text{ where } g_+(X) = \begin{cases} 0 & \text{if } X_{mn} \geq 0 \forall m, n \\ \infty & \text{else.} \end{cases} \quad (23)$$

In the notation of [Equation 1](#), this corresponds to the minimally non-trivial case of $N = 2$ and $M_j = 1, L_{ij} = 1$ for $j = 1, 2$. The most basic NMF solver uses a ‘‘multiplicative update’’ (MU) rule that can be derived from a gradient descent argument ([Lee and Seung 2001](#)). However, MU has long been criticized for its often inferior convergence properties, which can be traced back to the implicit treatment of the constraint. Several alternative approaches have been brought forward to address the shortcomings of MU solvers, including the ability to impose constraints other than non-negativity (e.g. [Berry et al 2007](#)). A comparison of these different approaches is not the focus of this work (see [Xu and Yin \(2013\)](#) for a recent overview). Instead, we demonstrate that bSDMM can successfully and efficiently impose several constraints on both matrix factors.

To do so, we need the proximal-operator forms of the desired constraints. In cases where the constraint is given by an indicator function of a convex set \mathcal{C} , the proximal operator is simply the projection operator onto \mathcal{C} under the Euclidean norm, and the step size λ is irrelevant. For example, the non-negativity constraint becomes the (element-wise) projection onto the non-negative orthant:

$$\text{prox}_{\lambda g_+}(\mathbf{x}) = \max(\mathbf{0}, \mathbf{x}). \quad (24)$$

Many other proximal operators can be evaluated analytically, e.g. for penalty functions involving ℓ_p norms, Total Variation, Maximum Entropy ([Combettes and Pesquet 2011](#); [Parikh et al 2014](#)). If K is large, the data are noisy, or $B \ll L$, the NMF factors are generally degenerate. Additional constraints then become necessary for reasonable results.

3.1 Example: Hyperspectral unmixing

Hyperspectral data are images taken of the same scene at several wavelengths ($B \sim 100$), often beyond the range visible to humans. This extended and more fine-grained spectral information can be used to robustly identify distinct components in the images. If the spatial resolution of the imager is high enough, those components can be observed in their pure form, i.e. each pixel received contributions only from one component, but in general the spectral information of pixels is mixed. Under the assumption that the components do not interfere, the so-called ‘‘linear mixing model’’, the NMF allows us to unmix the spectral contributions of each pixel, simultaneously inferring the pure spectrum of each of K components, called ‘‘endmembers’’, as well as their amplitude in each pixel ([Berry et al 2007](#); [Jia and Qian 2009](#); [Gillis 2014](#)).

[Figure 1](#) shows a false-color image of the National Mall in Washington D.C. made from hyperspectral HYDICE ([Mitchell 1995](#)) data comprising $B = 191$ wavelengths from 400 to

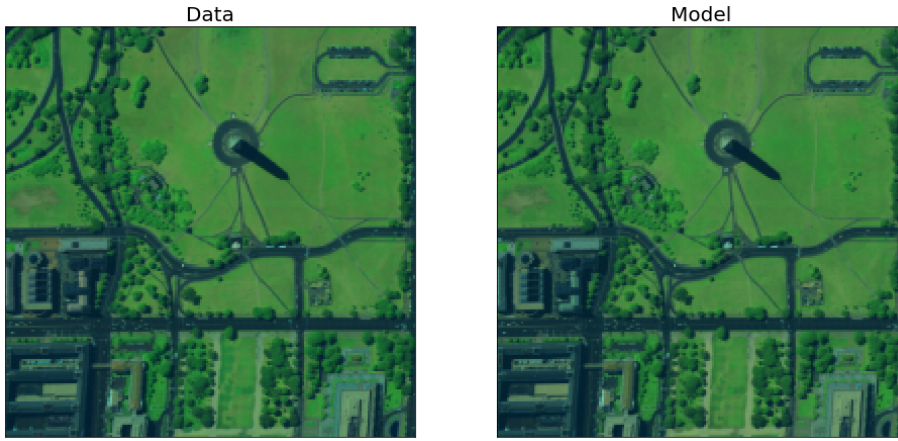


Fig. 1 False-colored images of hyperspectral data and a four-component model comprising “concrete”, “soil”, “vegetation”, and a spatially flat “background”. The false-color image maps the sum of the first 50 wavelengths to blue, the next 50 to green, and the remaining 91 wavelengths to red.

2475 nm.⁴ By our convention, the columns of matrix A describe the endmember spectra, and the rows of S the endmember amplitude per pixel. Despite the large number of wavelengths, the problem is still strongly underconstrained even with a small number of endmembers. We therefore impose several constraints:

- To prevent the degeneracy between A and S that stems from the transformation $(A, S) \rightarrow (AQ, Q^{-1}S)$ with an arbitrary invertible matrix Q , we normalize the endmember spectra, i.e. the columns of A . This normalization is different from the one usually adopted in hyperspectral unmixing applications, where the endmember amplitudes are normalized in each pixel, which results in endmembers being defined by both shape and amplitude of the spectrum. We prefer our approach because it maintains spectral similarity between regions of different brightness.
- The radiation recorded by the hyperspectral camera is a combination of light reflected off the ground and the atmosphere. Since we are interested in the former, and the latter is not expected to vary over the image, we add a “background” component that we constrain to be spatially flat.
- As the scene on the ground is mostly coherent over large areas, we add a two-dimensional anisotropic total variation (TV) penalty (Chambolle and Lions 1997; Chambolle 2004).

In summary, we minimize

$$\|A \cdot S - D\|_2^2 + g_+(A) + g_+(S) + g_{\text{norm}}(L_{\text{norm}}A) + g_{\text{bg}}(S) + \lambda (\|G_x S\|_1 + \|G_y S\|_1), \quad (25)$$

where g_+ is given in Equation 23, and g_{norm} and g_{bg} are additional indicator functions. In detail, we combine the NMF fidelity term and the positivity constraints into forward-backward operators of the form of Equation 3, one for A and one for S . This is the proximal-gradient technique for solving the NMF (e.g. Xu and Yin 2013), implemented in Line 10 of Algorithm 3.

⁴ Data set obtained from <https://engineering.purdue.edu/~biehl/MultiSpec/>

The normalization $L_{\text{norm}}A = 1_K$ is achieved by $L_{\text{norm}} = 1_B^\top$, for which the proximal operator is simply a projection onto 1_K .

The background component requires that $S_{\text{bg}} = \text{const}$ for all pixels. Since the proximal operator must yield the closest point on the submanifold in the Euclidean norm, it alters the background component row $S_{\text{bg}} \rightarrow \langle S_{\text{bg}} \rangle_L 1_L$, where the expectation value is carried out over all pixels, while leaving all other components unchanged.

For the TV penalty, we use the gradient operators in horizontal (G_x) and vertical (G_y) direction and make use of the analytic form of the proximal operator for the l_1 norm, the soft-thresholding operator (e.g. [Combettes and Wajs 2005](#)). While we could adjust the penalty parameter λ for the horizontal and vertical direction, as well as for every component, we have not found it necessary to explore that option.

The last four constraint and penalty functions are implemented in the SDMM fashion, giving rise to one auxiliary variable for A and three for S .

We run bSDMM with a TV penalty of $\lambda = 10$ until feasibility with $\epsilon^{\text{rel}} = 0.01$ for primal and dual residuals is reached (ϵ^{abs} was set to zero). We chose the number of endmembers to be three, which we label as “concrete”, “soil”, and “vegetation”, plus the flat background.⁵

While random initialization of the matrices works reasonably well, we found better results when we initialize the spectra by a three-step approach: First, we determine the background spectrum by the minimum value over all pixels at a given wavelength. Second, we select a reference pixel that, from its location in the scene, should be a pure representation of one of the three other components, and subtract from that pixel the background spectrum. Third, we normalize all spectra to sum up to one. The initialization of S is less important. We start with a zero matrix and then utilize bSDMM to make suitable updates given the initialized spectra.

The convergence of A , S , and the entire model $A \cdot S$ is shown in the top row of [Figure 2](#), primal and dual feasibility according to [Equation 10](#) in the middle and bottom row. The residual requirements are shown as dashed lines in [Figure 2](#), from which we can conclude that primal feasibility is generally achieved after about 30 iterations, while the dual feasibility of the background component requires almost 150 iterations. The resulting endmember spectra and amplitudes are shown in [Figure 3](#) and [Figure 4](#). The complete code to reproduce this hyperspectral unmixing example with the bSDMM-NMF approach is available at <https://github.com/fred3m/hyperspectral>.

4 Conclusion

In this work we have built upon the ADMM as a fast and flexible solver for constrained optimization problems. We have extended it in two directions. First, we allow for multiple constraint functions to be applied. Unlike the previously proposed SDMM approach by [Combettes and Pesquet \(2011\)](#), we do not require that linear operators of the constraint functions, which may be needed for an efficient proximal operator formulation, be invertible. Second, we address the case of a function that is convex in multiple arguments through an inexact block optimization method. The proposed method, bSDMM, is effective in a range of constrained optimization problems that cannot fully be solved with e.g. proximal gradient methods. As a result of its ADMM heritage, it is particularly suitable for applications where a fast and approximate solution is more important than an accurate one.

⁵ The choice of $K = 4$ is somewhat arbitrary, and we have not attempted to find the optimal number of components since that is not the relevant aspect of this work.

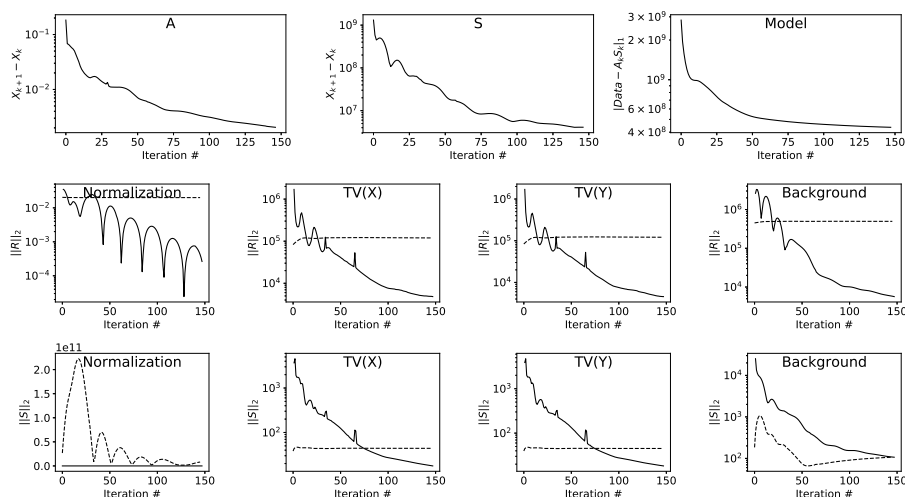


Fig. 2 Convergence to a stationary point of the problem in Equation 25: A, S, and the model $A \cdot S$ (top); the primal residual r_{ij}^k (middle) and the dual residual s_{ij}^k (bottom) for the four constraint variables. The limits ϵ^{pri} and ϵ^{dual} for primal and dual feasibility (cf. Equation 10) are shown as dashed lines with $\epsilon^{rel} = 0.01$.

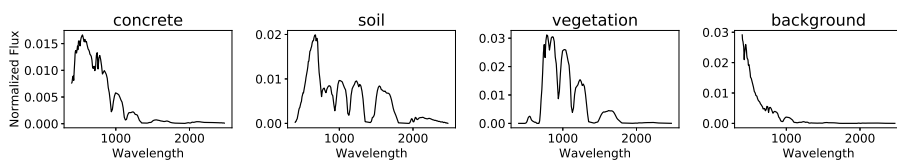


Fig. 3 Endmember spectrum for every component. Each spectrum is normalized to unity. The labels are approximate descriptions given the regions in the image the endmembers mostly represent.

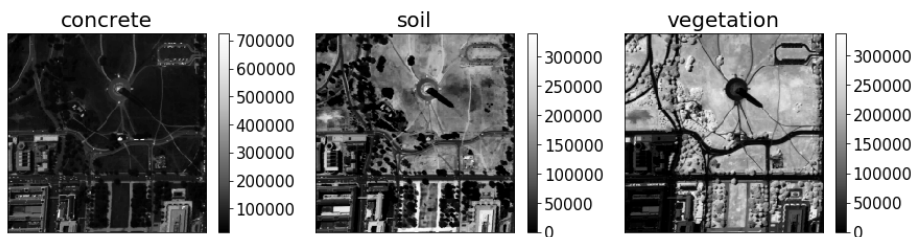


Fig. 4 Intensity of the three spatially variable components “concrete”, “soil”, and “vegetation”. Several features are prominent in the endmember intensities, such as trees, trails, and rooftops. These plots do not represent an endmember classification: because of our spectrum normalization, a region appears dark in these plots if it has a different spectrum than the endmember or if it reflects very little light (e.g. road surfaces).

We showed its effectiveness to solve a hyperspectral unmixing problem, under the assumption of the linear mixing model, by performing a Non-Negative Matrix Factorization with multiple non-trivial constraints. In a future work (Melchior et al., in prep.) we will utilize the bSDMM-NMF approach to separate stars and galaxies in astronomical images, a similar additive mixing problem, which requires several restrictive constraints for adequate performance.

As we believe in the usefulness of the algorithm and want to endorse reproducible research, we release the python implementation of the algorithms presented here and the bSDMM-NMF as an open-source package at <https://github.com/pmelchior/proxmin>.

Acknowledgements

We would like to thank Robert Vanderbei and Jonathan Eckstein for useful discussions regarding the algorithm, and Jim Bosch and Robert Lupton for comments on its astrophysical applications.

References

- Berry MW, Browne M, Langville AN, Pauca VP, Plemmons RJ (2007) Algorithms and applications for approximate nonnegative matrix factorization. *Computational statistics & data analysis* 52(1):155–173
- Blanton MR, Roweis S (2007) K-Corrections and Filter Transformations in the Ultraviolet, Optical, and Near-Infrared. *Astronomical Journal* 133:734–754, DOI 10.1086/510127, [astro-ph/0606170](https://arxiv.org/abs/astro-ph/0606170)
- Boyd S, Parikh N, Chu E, Peleato B, Eckstein J (2011) Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning* 3(1):1–122
- Chambolle A (2004) An algorithm for total variation minimization and applications. *Journal of Mathematical Imaging and Vision* 20(1):89–97, DOI 10.1023/B:JMIV.0000011325.36760.1e, URL <https://doi.org/10.1023/B:JMIV.0000011325.36760.1e>
- Chambolle A, Lions PL (1997) Image recovery via total variation minimization and related problems. *Numerische Mathematik* 76(2):167–188, DOI 10.1007/s002110050258, URL <https://doi.org/10.1007/s002110050258>
- Chen G, Teboulle M (1994) A proximal-based decomposition method for convex minimization problems. *Mathematical Programming* 64(1-3):81–101
- Combettes PL, Pesquet JC (2011) Proximal splitting methods in signal processing. In: *Fixed-point algorithms for inverse problems in science and engineering*, Springer, pp 185–212
- Combettes PL, Wajs VR (2005) Signal recovery by proximal forward-backward splitting. *Multiscale Modeling & Simulation* 4(4):1168–1200
- Condat L (2013) A primal–dual splitting method for convex optimization involving lipschitzian, proximable and linear composite terms. *Journal of Optimization Theory and Applications* 158(2):460–479
- Douglas J, Rachford HH (1956) On the numerical solution of heat conduction problems in two and three space variables. *Transactions of the American mathematical Society* 82(2):421–439
- Eckstein J, Bertsekas DP (1992) On the douglasrachford splitting method and the proximal point algorithm for maximal monotone operators. *Mathematical Programming* 55(1):293–318
- Eckstein J, Yao W (2017) Approximate admm algorithms derived from lagrangian splitting. *Computational Optimization and Applications* DOI 10.1007/s10589-017-9911-z, URL <https://doi.org/10.1007/s10589-017-9911-z>
- Esser E, Zhang X, Chan TF (2010) A general framework for a class of first order primal-dual algorithms for convex optimization in imaging science. *SIAM Journal on Imaging Sciences* 3(4):1015–1046
- Gabay D, Mercier B (1976) A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Computers & Mathematics with Applications* 2(1):17–40
- Gillis N (2014) *The Why and How of Nonnegative Matrix Factorization*, Chapman and Hall/CRC, pp 257–291. DOI doi:10.1201/b17558-13, URL <https://doi.org/10.1201/b17558-13>
- Glowinski R, Marroco A (1975) Sur l’approximation, par éléments finis d’ordre un, et la résolution, par pénalisation-dualité d’une classe de problèmes de dirichlet non linéaires. *Revue française d’automatique, informatique, recherche opérationnelle Analyse numérique* 9(2):41–76

- Grippo L, Sciandrone M (2000) On the convergence of the block nonlinear gauss–seidel method under convex constraints. *Operations research letters* 26(3):127–136
- Hong M, Luo ZQ, Razaviyayn M (2016) Convergence analysis of alternating direction method of multipliers for a family of nonconvex problems. *SIAM Journal on Optimization* 26(1):337–364
- Jia S, Qian Y (2009) Constrained nonnegative matrix factorization for hyperspectral unmixing. *IEEE Transactions on Geoscience and Remote Sensing* 47(1):161–173, DOI 10.1109/TGRS.2008.2002882
- Lee DD, Seung HS (2001) Algorithms for non-negative matrix factorization. In: Leen TK, Dietterich TG, Tresp V (eds) *Advances in Neural Information Processing Systems* 13, MIT Press, pp 556–562, URL <http://papers.nips.cc/paper/1861-algorithms-for-non-negative-matrix-factorization.pdf>
- Lin CJ (2007) Projected gradient methods for nonnegative matrix factorization. *Neural computation* 19(10):2756–2779
- Mitchell PA (1995) Hyperspectral digital imagery collection experiment (hydice). In: *Proc.SPIE*, vol 2587, pp 2587 – 2587 – 26, DOI 10.1117/12.226807, URL <http://dx.doi.org/10.1117/12.226807>
- Nesterov Y (2013) Gradient methods for minimizing composite functions. *Mathematical Programming* 140(1):125–161
- Paatero P, Tapper U (1994) Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics* 5(2):111–126
- Parikh N, Boyd S, et al (2014) Proximal algorithms. *Foundations and Trends® in Optimization* 1(3):127–239
- Razaviyayn M, Hong M, Luo ZQ (2013) A unified convergence analysis of block successive minimization methods for nonsmooth optimization. *SIAM Journal on Optimization* 23(2):1126–1153
- Stephanopoulos G, Westerberg AW (1975) The use of hestenes’ method of multipliers to resolve dual gaps in engineering system optimization. *Journal of Optimization Theory and Applications* 15(3):285–309
- Wang Y, Yin W, Zeng J (2015) Global convergence of admm in nonconvex nonsmooth optimization. *arXiv preprint arXiv:151106324*
- Xu Y, Yin W (2013) A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion. *SIAM Journal on imaging sciences* 6(3):1758–1789
- Zhang S, Qian H, Gong X (2016) An alternating proximal splitting method with global convergence for nonconvex structured sparsity optimization. In: *AAAI*, pp 2330–2336
- Zhu G (2016) Nonnegative Matrix Factorization (NMF) with Heteroscedastic Uncertainties and Missing data. *ArXiv e-prints* [1612.06037](https://arxiv.org/abs/1612.06037)