

AN AUGMENTED LAGRANGIAN METHOD FOR ℓ_1 -REGULARIZED OPTIMIZATION PROBLEMS WITH ORTHOGONALITY CONSTRAINTS

WEIQIANG CHEN*, HUI JI*, AND YANFEI YOU[†]*

Abstract. A class of ℓ_1 -regularized optimization problems with orthogonality constraints has been used to model various applications arising from physics and information sciences, e.g., compressed modes for variational problems. Such optimization problems are difficult to solve due to the non-smooth objective function and non-convex constraints. Existing methods either are not applicable to such problems, or lack convergence analysis, e.g., the splitting of orthogonality constraints (SOC) method. In this paper, we propose a proximal alternating minimized augmented Lagrangian (PAMAL) method that hybridizes the augmented Lagrangian method and the proximal alternating minimization scheme. It is shown that the proposed method has the so-called sub-sequence convergence property, i.e., there exists at least one convergent sub-sequence and any convergent sub-sequence converges to a Karush-Kuhn Tucker (KKT) point of an equivalent minimization problem. Experiments on the problem of compressed modes illustrate that the proposed method is noticeably faster than the SOC method.

Key words. augmented Lagrangian; ℓ_1 -regularization; orthogonality constraints

AMS subject classifications. 65K10, 90C26, 49R05, 65L15

1. Introduction. In the last few decades, the concept of sparsity has been extensively exploited in a wide range of applications in imaging and information science. Most of these methods focus on the sparsity of the coefficients used for representing the corresponding vector with a set of atoms. The majority of sparsity-driven applications use the ℓ_1 -norm as the convex relaxation of the sparsity-prompting function in their variational formulations. Such applications include compressed sensing [10, 12, 14], model selection and learning [11, 27, 38], and image recovery [13, 35, 36]. Most of the optimization problems arising from these applications are convex problems. In the last ten years, there has been a huge growth in literature on efficient numerical solvers for these problems; see e.g., [9, 17, 21, 29].

Nevertheless, there are also many applications in which the data must satisfy non-convex constraints. One commonly seen non-convex constraint is the orthogonality constraint, i.e., the data for estimation can be expressed as an orthogonal matrix. Examples of such applications include sparse principal component analysis [24], eigenvalue problems in sub-space tracking [43] and mathematical physics [31], and orthogonal Procrustes problems in shape analysis [15]. Because orthogonality constraints are non-convex, such problems can be difficult, except in a few simple cases. Recently, the idea of sparsity is also exploited for problems with orthogonality constraints [31, 32, 37, 41, 42, 44], and ℓ_1 -regularization is introduced in the resulting variational model to regularize the sparsity of the data. We briefly describe two representative applications that involve ℓ_1 -regularized optimization problems with orthogonality constraints.

(a) *Compressed modes (waves) in physics* [31, 32]. Compressed modes are spatially localized solutions to the eigenvalue problem of the Schrödinger's equation. By considering the independent-particle Schrödinger's equation for a finite system of electrons, the corresponding eigenvalue problem can be reformulated as follows:

$$\min_{X \in \mathbb{R}^{n \times m}} \frac{1}{\mu} \|X\|_1 + \text{Tr}(X^\top H X) \quad \text{s.t.} \quad X^\top X = I_m, \quad (1.1)$$

where $\|X\|_1 := \sum_{i=1}^n \sum_{j=1}^m |X_{i,j}|$, μ is a pre-defined positive parameter that balances the

*Department of Mathematics, National University of Singapore, Singapore, 117542

[†]Corresponding author. School of Mathematics and Physics, Changzhou University, Jiangsu, China, 213164. This author was supported by the NSFC Grant 11471156.

sparsity and the accuracy of the solution, H denotes the (discrete) Hamiltonian, and the columns of X denote the eigenvectors with local support, the so-called compressed modes.

- (b) *Feature selection* [37, 44]. Feature selection seeks to choose a smaller subset of features with most information from high dimensional feature sets. It is used in computer vision [37] and social media data [44], etc. The models for feature selection in [37, 44] adhere to the following ℓ_1 -regularized optimization problem with (weighted) orthogonality constraints:

$$\min_{X \in \mathbb{R}^{n \times m}} \frac{1}{\mu} \|X\|_{2,1} + \text{Tr}(X^\top H X), \quad \text{s.t.} \quad X^\top M X = I_m, \quad (1.2)$$

where $\|X\|_{2,1} := \sum_{i=1}^n (\sum_{j=1}^m X_{i,j}^2)^{1/2}$, H is a symmetric matrix, and M is a symmetric positive definite matrix.

This paper aims at developing a numerical method to solve (1.1), as well as (1.2) with minor modifications. The proposed PAMAL method can be viewed as a method that hybridizes the augmented Lagrangian method [2] and the proximal alternating minimization (PAM) techniques proposed in [4]. The convergence analysis established in this paper shows that under very mild assumptions on the associated penalty parameters, the sequence generated by the proposed method has the *sub-sequence convergence property*, i.e., there exists at least one convergent sub-sequence and any convergent sub-sequence converges to a Karush-Kuhn Tucker (KKT) point of (2.1) (see (3.4) for details).

1.1. Related work on optimization problems with orthogonality constraints. We now give a brief review on the existing methods that can be applied to solve the problems with orthogonality constraints:

$$\min_{X \in \mathbb{R}^{n \times m}} J(X) \quad \text{s.t.} \quad X^\top X = I_m, \quad (1.3)$$

where J might be non-convex and non-differentiable. Existing numerical methods that are applicable to (1.3) can be classified under two categories: feasible and infeasible approaches.

The feasible approaches satisfy the constraints during each iteration, i.e., each point of the sequence generated by the approach satisfies the orthogonality constraints in (1.3). In fact, various optimization methods such as Newton's method, the conjugate gradient method, and the method of steepest descent have been used to solve (1.3) as feasible approaches. Most of the existing methods are based on the study of the manifold structures of the orthogonality constraints (see e.g., [1, 16, 20, 25, 40]). These methods require the objective function J to be differentiable, which is not applicable to the problem (1.1) studied in this paper. Furthermore, it is not trivial to satisfy the orthogonality constraints in (1.1) during each iteration, as suggested in [23]. Therefore, the feasible approach might not be ideal to solve (1.1) as its objective function is non-differentiable.

The PAMAL method proposed in this paper is an infeasible approach. Infeasible approaches simplify the constrained problem (1.3) by relaxing the constraints and iteratively diminish the degree of infeasibility. As a result, intermediate points of the generated sequence may not satisfy the orthogonality constraints. The penalty method (e.g., [7, 28]) approximates the problem (1.3) by penalizing the deviations from the constraints:

$$\min_{X \in \mathbb{R}^{n \times m}} J(X) + \frac{1}{2\kappa} \|X^\top X - I_m\|_F^2,$$

where κ denotes some penalty parameter decreasing to zero. If $J(X) = \text{Tr}(X^\top H X)$, then the quadratic penalty model can be viewed as an exact penalty method with a finite penalty pa-

parameter κ ; see e.g., [39]. While the penalty method is simple, it suffers from ill-conditioning issues, especially when the penalty parameter κ decreases to zero. Thus, the standard augmented Lagrangian method [18, 19] is often preferred as it does not require the parameter κ to decrease to zero. When applied to solve (1.3), the standard augmented Lagrangian method yields the following scheme:

$$\begin{cases} X^{k+1} & \in \operatorname{argmin}_X J(X) + \frac{\rho^k}{2} \|X^\top X - I_m\|_F^2 + \operatorname{Tr}((\Lambda^k)^\top (X^\top X - I_m)), \\ \Lambda^{k+1} & = \Lambda^k + \rho^k ((X^{k+1})^\top X^{k+1} - I_m). \end{cases}$$

The sub-problem of the above augmented Lagrangian scheme is rather complex and generally has no analytic solution. Indeed, it is not trivial to design an efficient solver for this sub-problem.

Aiming at a more computationally efficient method for solving (1.3), the SOC method [23] introduces auxiliary variables to split the orthogonality constraints, which leads to another formulation of (1.3):

$$\min_{X, P \in \mathbb{R}^{n \times m}} J(X) \quad \text{s.t.} \quad X = P, \quad P^\top P = I_m. \quad (1.4)$$

Using the ideas of alternating direction method of multipliers (ADMM) and the split Bregman method, the SOC method solves (1.4) by alternately updating the variables $\{X, P, B\}$:

$$\begin{cases} X^{k+1} = \operatorname{argmin}_X J(X) + \frac{\rho}{2} \|X - P^k + B^k\|_F^2; \\ P^{k+1} = \operatorname{argmin}_P \frac{\rho}{2} \|P - (X^{k+1} + B^k)\|_F^2, \quad \text{s.t.} \quad P^\top P = I; \\ B^{k+1} = B^k + X^{k+1} - P^{k+1}. \end{cases} \quad (1.5)$$

In contrast with the standard augmented Lagrangian method, each sub-problem in the iterations of the SOC method has an analytic solution. However, the trade-off is its challenging convergence analysis. To the best of our knowledge, it remains an open question whether the SOC method (1.5) has the sub-sequence convergence property.

1.2. Notations and preliminaries on non-smooth analysis. In this section, we introduce some notations and preliminaries on non-smooth analysis. Given a matrix $Y \in \mathbb{R}^{n \times m}$, its maximum (element-wise) norm is denoted by

$$\|Y\|_\infty := \max_{i,j} |Y_{i,j}|, \quad i = 1, \dots, n, \quad j = 1, \dots, m,$$

and we denote by $\operatorname{Vec}(Y)$ the $mn \times 1$ vector formed by stacking $\{Y_j\}_{j=1}^m$, the column vectors of Y , on top of one another, i.e., $\operatorname{Vec}(Y) := [Y_1^\top, Y_2^\top, \dots, Y_m^\top]^\top$. For any $v \in \mathbb{R}^n$, let $[v]_i$ denote its i -th component and let $\operatorname{diag}(v) \in \mathbb{R}^{n \times n}$ denote the diagonal matrix with diagonal entries $\{[v]_i\}_{i=1}^n$. For an index sequence $\mathcal{K} = \{k_0, k_1, k_2, \dots\}$ that satisfies $k_{j+1} > k_j$ for any $j \geq 0$, we denote $\lim_{k \in \mathcal{K}} x_k := \lim_{j \rightarrow \infty} x_{k_j}$. For any set S , its indicator function is defined by

$$\delta_S(X) = \begin{cases} 0, & \text{if } X \in S, \\ +\infty, & \text{otherwise.} \end{cases} \quad (1.6)$$

For a given matrix X and a constant $\alpha > 0$, the soft-thresholding operator is defined by

$$T_\alpha^1(X) := [T_\alpha^1(X_{i,j})]_{i,j}, \quad \text{where } T_\alpha^1(x) := \operatorname{sign}(x) \max(0, |x| - \alpha), \quad x \in \mathbb{R}. \quad (1.7)$$

DEFINITION 1.1 ([34]). *Let $S \subseteq \mathbb{R}^n$ and $\bar{x} \in S$. A vector v is normal to S at \bar{x} in the regular sense, expressed as $v \in \widehat{N}_S(\bar{x})$, if*

$$\langle v, x - \bar{x} \rangle \leq o(\|x - \bar{x}\|) \quad \text{for } x \in S,$$

where $o(\|y\|)$ is defined by $\lim_{\|y\| \rightarrow 0} \frac{o(\|y\|)}{\|y\|} = 0$. A vector is normal to S at \bar{x} in the general sense, expressed as $v \in N_S(\bar{x})$, if there exist sequences $\{x^k\}_k \subset S$, $\{v^k\}_k$ such that $x^k \rightarrow \bar{x}$ and $v^k \rightarrow v$ with $v^k \in \widehat{N}_S(x^k)$. The cone $N_S(\bar{x})$ is called the normal cone to S at \bar{x} .

For a proper and lower semi-continuous function, denoted by $\sigma : \mathbb{R}^n \rightarrow (-\infty, +\infty]$, the domain of σ is defined as $\text{dom } \sigma := \{x \in \mathbb{R}^n : \sigma(x) < +\infty\}$.

DEFINITION 1.2 ([34]). Consider a proper and lower semi-continuous function $\sigma : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ and a point \bar{x} with finite $\sigma(\bar{x})$. Let $v \in \mathbb{R}^n$.

1. The vector v is said to be a regular sub-gradient of σ at \bar{x} , expressed as $v \in \widehat{\partial}\sigma(\bar{x})$, if $\sigma(x) \geq \sigma(\bar{x}) + \langle v, x - \bar{x} \rangle + o(\|x - \bar{x}\|)$.
2. The vector v is said to be a (limiting) sub-gradient of σ at \bar{x} , expressed as $v \in \partial\sigma(\bar{x})$, if there exist sequences $\{x^k\}_k$, $\{v^k\}_k$ such that $x^k \rightarrow \bar{x}$, $\sigma(x^k) \rightarrow \sigma(\bar{x})$ and $v^k \in \widehat{\partial}\sigma(x^k)$ with $v^k \rightarrow v$.
3. For each $x \in \text{dom } \sigma$, x is called a (limiting)-critical point of σ if $0 \in \partial\sigma(x)$.

We end this section with a result used for subsequent discussion.

REMARK 1.3 ([34, Example 6.7]). Let S be a closed non-empty subset of \mathbb{R}^n , then

$$\partial\delta_S(\bar{x}) = N_S(\bar{x}), \quad \bar{x} \in S.$$

Furthermore, for a smooth mapping $G : \mathbb{R}^n \rightarrow \mathbb{R}^m$, i.e., $G(x) := (g_1(x), \dots, g_m(x))^\top$, define $S = G^{-1}(0) \subset \mathbb{R}^n$. Set $\nabla G(x) := [\frac{\partial g_j}{\partial x_i}(x)]_{i,j=1}^{n,m} \in \mathbb{R}^{n \times m}$. If $\nabla G(\bar{x})$ has full rank m at a point $\bar{x} \in S$, with $G(\bar{x}) = 0$, then its normal cone to S can be explicitly written as

$$N_S(\bar{x}) = \{\nabla G(\bar{x})y \mid y \in \mathbb{R}^m\}. \quad (1.8)$$

2. PAMAL method. To solve (1.1), we present the PAMAL method, which hybridizes an augmented Lagrangian method with proximal alternating minimization. Let

$$\mathcal{S} := \{P \in \mathbb{R}^{n \times m} : P^\top P = I_m\}$$

denote the Stiefel manifold, and let $\delta_{\mathcal{S}}$ denote its indicator function (1.6). Then, one may re-write (1.1) as

$$\begin{aligned} \min_{X, Q, P \in \mathbb{R}^{n \times m}} \quad & \frac{1}{\mu} \|Q\|_1 + \text{Tr}(X^\top HX) + \delta_{\mathcal{S}}(P) \\ \text{s.t.} \quad & Q - X = 0, \quad \text{and} \quad P - X = 0. \end{aligned} \quad (2.1)$$

Let $\Lambda := (\Lambda_1, \Lambda_2) \subset \mathbb{R}^{n \times 2m}$. The augmented Lagrangian of (2.1) can be expressed as

$$\begin{aligned} L(X, Q, P, \Lambda; \rho) = \quad & \frac{1}{\mu} \|Q\|_1 + \text{Tr}(X^\top HX) + \langle \Lambda_1, Q - X \rangle + \frac{\rho}{2} \|Q - X\|_F^2 \\ & + \langle \Lambda_2, P - X \rangle + \frac{\rho}{2} \|P - X\|_F^2 + \delta_{\mathcal{S}}(P), \end{aligned} \quad (2.2)$$

where ρ is a positive penalty parameter. In the proposed method, (2.2) is solved via the augmented Lagrangian scheme [2], which alternately updates the estimate of (X, Q, P) , the multiplier Λ and the penalty parameter ρ . The main step is about how to update (X, Q, P) . In the proposed method, it is solved by the PAM method [4], which deals with a class of non-smooth and non-convex optimization problems. We now describe Algorithm 1.

Algorithm 1 : Method for solving (2.2)

Given pre-defined parameters $\{\epsilon^k\}_{k \in \mathbb{N}}$, $\bar{\Lambda}^1 := (\bar{\Lambda}_1^1, \bar{\Lambda}_2^1)$, ρ^1 , $\bar{\Lambda}_{p,\min}$, $\bar{\Lambda}_{p,\max}$, τ , γ that satisfy the conditions in Remark 2.1, for $k = 1, 2, \dots$,

1. Compute (X^k, Q^k, P^k) such that there exists $\Theta^k \in \partial L(X^k, Q^k, P^k, \bar{\Lambda}^k; \rho^k)$ satisfying

$$\|\Theta^k\|_\infty \leq \epsilon^k, \quad (P^k)^\top P^k = I_m, \quad (2.3)$$

where $\{\epsilon^k\}_{k \in \mathbb{N}}$ is a sequence of positive tolerance parameters.

2. Update the multiplier estimates:

$$\Lambda_1^{k+1} = \bar{\Lambda}_1^k + \rho^k(Q^k - X^k), \quad \Lambda_2^{k+1} = \bar{\Lambda}_2^k + \rho^k(P^k - X^k),$$

where $\bar{\Lambda}_p^{k+1}$ is the projection of Λ_p^{k+1} on $\{\Lambda_p : \bar{\Lambda}_{p,\min} \leq \Lambda_p \leq \bar{\Lambda}_{p,\max}\}$, $p = 1, 2$.

3. Update the penalty parameter:

$$\rho^{k+1} := \begin{cases} \rho^k, & \text{if } \|R_i^k\|_\infty \leq \tau \|R_i^{k-1}\|_\infty, i = 1, 2; \\ \gamma \rho^k, & \text{otherwise,} \end{cases} \quad (2.4)$$

where $R_1^k := Q^k - X^k$, $R_2^k := P^k - X^k$.

Step 1 of Algorithm 1 seeks the updates of primal variables such that there is an associated sub-gradient element of L , which satisfies a specified level of tolerance. Step 2 of Algorithm 1 updates the multiplier estimates by first computing the first-order approximations of the multipliers, which are then projected on a suitable box to ensure compactness. Step 3 of Algorithm 1 updates the penalty parameter ρ^k according to the degree of infeasibility. Note that the pre-defined parameters in Algorithm 1 will impact its convergence property. Remark 2.1 discusses the setting of these parameters.

REMARK 2.1 (Parameter setting). *The parameters in Algorithm 1 are set as follows. The sequence of positive tolerance parameters $\{\epsilon^k\}_{k \in \mathbb{N}}$ in (2.3) is chosen such that $\lim_{k \rightarrow \infty} \epsilon^k = 0$. The parameters $\bar{\Lambda}_1^1$, $\bar{\Lambda}_2^1$, $\bar{\Lambda}_{p,\min}$, $\bar{\Lambda}_{p,\max}$ are finite-valued matrices satisfying*

$$-\infty < [\bar{\Lambda}_{p,\min}]_{i,j} < [\bar{\Lambda}_{p,\max}]_{i,j} < \infty, \forall i, j, \quad p = 1, 2.$$

As we shall see in Section 2.2, for Step 1 of Algorithm 1 to be well defined, it suffices to have $\tau \in [0, 1)$, $\gamma > 1$ and ρ^1 to be sufficiently large so that

$$\rho^1 I_n + 2H \succ 0,$$

where H is the (discrete) Hamiltonian in (1.1).

REMARK 2.2 (Relation with the SOC method). *Both the PAMAL and the SOC methods [31] use the same splitting technique employed in (2.1). The main difference between the PAMAL method and the SOC method lies in how (X^k, Q^k, P^k) is updated. In the PAMAL method, the update (Step 1) is done by calling Algorithm 2, which runs several inner iterations to obtain an approximate solution to a critical point (X^k, Q^k, P^k) for L_k with a pre-defined tolerance ϵ_k , i.e.,*

$$\Theta^k \in \partial L(X^k, Q^k, P^k, \bar{\Lambda}^k; \rho^k), \text{ s.t. } \|\Theta^k\| \leq \epsilon^k. \quad (2.5)$$

The tolerance parameter sequence $\{\epsilon_k\}_{k \in \mathbb{N}}$ can be set to decrease to zero. In contrast, the SOC method only uses a single inner iteration in every outer iteration to solve the problem in

Step 1. Thus, there is no guarantee that its corresponding tolerance sequence will converge to zero, which makes the convergence analysis of the SOC method a very challenging task. Despite the fact that multiple inner iterations might be used in Algorithm 1, the flexibility on the accuracy of the solution in Step 1 makes it more computationally efficient in practice. For example, when applied to solve the compressed modes problem, the PAMAL method uses much fewer outer iterations to attain the stopping criterion, and the number of inner iterations in most outer iterations is only 1 or 2. As a result, the total number of inner iterations of the PAMAL method is less than that of the SOC method (See Tables 4.1 and 4.2).

Step 1 is the most crucial and difficult step of Algorithm 1, where the constraints (2.3) can also be viewed as relaxed KKT conditions for minimizing the augmented Lagrangian L (2.2). Thus, there are two questions to answer when executing Step 1 of Algorithm 1:

1. For each $k \in \mathbb{N}$, is Step 1 of Algorithm 1 well defined? In other words, is the existence of the points (X^k, Q^k, P^k) satisfying (2.3) guaranteed?
2. How can we efficiently compute such points with arbitrarily given accuracy, i.e., can the perturbation ϵ^k be arbitrarily small?

In the next subsection, we first describe a method for solving (2.3), which answers Question 2. Then, we establish Proposition 2.5, which shows that the answer to Question 1 is positive for the proposed method.

2.1. Algorithm for Step 1 of Algorithm 1. It can be seen that the constraint (2.3) is actually an ϵ^k -perturbation of the so-called critical point property

$$0 \in \partial L(X^k, Q^k, P^k, \bar{\Lambda}^k; \rho^k). \quad (2.6)$$

Thus, we need a method that can evaluate the corresponding critical points (X^k, Q^k, P^k) of the function $L(X, Q, P, \bar{\Lambda}^k; \rho^k)$ with arbitrary accuracy. Based on the PAM method [4], we propose a coordinate-descent method with proximal regularization. The PAM method [4] is proposed for solving a class of non-smooth and non-convex optimization problems. Under certain conditions on the objective function, it is shown [4, Theorem 6.2] that the PAM method has global convergence, i.e., the sequence generated by the method converges to some critical point. In Section 2.2, we will show that the function $L(X, Q, P, \bar{\Lambda}^k, \rho^k)$ indeed satisfies the sufficient conditions for the global convergence of the PAM method, provided that the penalty parameters $\{\rho^k\}_{k \in \mathbb{N}}$ satisfy a mild condition. In other words, Step 1 of Algorithm 1 is well defined provided that the parameters in Algorithm 1 are appropriately chosen when the PAM method is employed. Algorithm 2 gives the outline of the method for solving (2.3).

The PAM method can be applied to solve (2.6) as follows. Indeed, at the k -th outer iteration, the problem (2.6) can be solved with arbitrary accuracy using the following set of inner iterations, which can be viewed as a proximal regularization of a three block Gauss-Seidel method:

$$\begin{cases} X^{k,j} \in \operatorname{argmin}_X L(X, Q^{k,j-1}, P^{k,j-1}, \bar{\Lambda}^k; \rho^k) + \frac{c_1^{k,j-1}}{2} \|X - X^{k,j-1}\|_F^2; \\ Q^{k,j} \in \operatorname{argmin}_Q L(X^{k,j}, Q, P^{k,j-1}, \bar{\Lambda}^k; \rho^k) + \frac{c_2^{k,j-1}}{2} \|Q - Q^{k,j-1}\|_F^2; \\ P^{k,j} \in \operatorname{argmin}_P L(X^{k,j}, Q^{k,j}, P, \bar{\Lambda}^k; \rho^k) + \frac{c_3^{k,j-1}}{2} \|P - P^{k,j-1}\|_F^2, \end{cases} \quad (2.7)$$

where the proximal parameters $\{c_i^{k,j}\}_{k,j}$, can be arbitrarily chosen as long as they satisfy

$$0 < \underline{c} \leq c_i^{k,j} \leq \bar{c} < \infty, \quad k, j \in \mathbb{N}, \quad i = 1, 2, 3,$$

for some pre-determined positive constants \underline{c} and \bar{c} .

It turns out that all sub-problems in (2.7) have analytic solutions. The solution to the first sub-problem is the least squares solution, the solution to the second sub-problem is obtained by soft-thresholding, and the solution to the last sub-problem can be obtained by the singular value decomposition (SVD). The iteration (2.7) is terminated when there exists $\Theta^{k,j} \in \partial L(X^{k,j}, Q^{k,j}, P^{k,j}, \bar{\Lambda}^k; \rho^k)$ satisfying

$$\|\Theta^{k,j}\|_\infty \leq \epsilon^k, \quad (P^{k,j})^\top P^{k,j} = I_m.$$

An explicit expression of the term $\Theta^{k,j} := (\Theta_1^{k,j}, \Theta_2^{k,j}, \Theta_3^{k,j})$ is given by

$$\begin{cases} \Theta_1^{k,j} := \rho^k(Q^{k,j-1} - Q^{k,j} + P^{k,j-1} - P^{k,j}) + c_1^{k,j-1}(X^{k,j-1} - X^{k,j}); \\ \Theta_2^{k,j} := c_2^{k,j-1}(Q^{k,j-1} - Q^{k,j}); \\ \Theta_3^{k,j} := c_3^{k,j-1}(P^{k,j-1} - P^{k,j}). \end{cases} \quad (2.8)$$

Algorithm 2 describes the method for solving (2.7), which completes Algorithm 1.

Algorithm 2 : Proposed method for solving (2.3)

1. Let $(X^{1,0}, Q^{1,0}, P^{1,0})$ be any initialization. For $k \geq 2$, set $(X^{k,0}, Q^{k,0}, P^{k,0}) := (X^{k-1}, Q^{k-1}, P^{k-1})$.
2. Re-iterate on j until $\|\Theta^{k,j}\|_\infty \leq \epsilon^k$, where $\Theta^{k,j}$ is defined by (2.8).
 1. $X^{k,j} = Z^{-1}(\bar{\Lambda}_1^k + \bar{\Lambda}_2^k + \rho^k Q^{k,j-1} + \rho^k P^{k,j-1} + c_1^{k,j-1} X^{k,j-1})$,
where $Z := Z^{k,j-1} = 2H + (\rho^k + \rho^k + c_1^{k,j-1})I_n$.
 2. $Q^{k,j} = T_\eta^1(\frac{\rho^k X^{k,j} - \bar{\Lambda}_1^k + c_2^{k,j-1} Q^{k,j-1}}{\rho^k + c_2^{k,j-1}})$, $\eta := \eta^{k,j} := \mu \cdot (\rho^k + c_2^{k,j-1})^{-1}$, where T_η^1 is the soft-thresholding operator defined by (1.7).
 3. $P^{k,j} = UI_{n \times m}V^\top$, where the matrices U, V are obtained from the SVD of

$$\frac{\rho^k X^{k,j} + c_3^{k,j-1} P^{k,j-1} - \bar{\Lambda}_2^k}{\rho^k + c_3^{k,j-1}} =: USV^\top.$$

3. Set $(X^k, Q^k, P^k) := (X^{k,j}, Q^{k,j}, P^{k,j})$ and

$$\Theta^k := \Theta^{k,j}. \quad (2.9)$$

2.2. Well definedness of Step 1 of Algorithm 1. In this subsection, we will show that Step 1 of Algorithm 1 is well defined with the use of Algorithm 2, provided that the initial positive penalty parameter ρ^1 satisfies $\rho^1 I_n + 2H \succ 0$ and $\gamma > 1$. In other words, we will show that the solutions for the constraint (2.3) are non-empty and Algorithm 2 can always find one solution, provided ρ^1 is appropriately chosen. For Step 1 to be well defined, it needs an important property of Algorithm 2, i.e., for each $k \in \mathbb{N}$,

$$(X^{k,j}, Q^{k,j}, P^{k,j}) \rightarrow (\bar{X}^k, \bar{Q}^k, \bar{P}^k), \quad \text{as } j \rightarrow \infty, \quad (2.10)$$

where $(\bar{X}^k, \bar{Q}^k, \bar{P}^k)$ is a critical point of $L(X, Q, P, \bar{\Lambda}^k; \rho^k)$. The proof of the limiting property (2.10) is based on the result [4, Theorem 6.2], which considers the minimization of a function $f : \mathbb{R}^{n_1} \times \dots \times \mathbb{R}^{n_p} \rightarrow \mathbb{R} \cup \{+\infty\}$ of the form

$$f(x) = g(x_1, \dots, x_p) + \sum_{i=1}^p f_i(x_i), \quad (2.11)$$

where the functions f_1, \dots, f_p and g satisfy the following assumptions:

- (i) $f_i : \mathbb{R}^{n_i} \rightarrow \mathbb{R} \cup \{+\infty\}$ is a proper lower semi-continuous function, $i = 1, 2, \dots, p$;
- (ii) g is a C^1 -function with locally Lipschitz continuous gradient;
- (iii) f is a K-Ł (Kurdyka-Łojasiewicz) function (see Remark 2.4 for more details), and $\inf_{\mathbb{R}^{n_1} \times \dots \times \mathbb{R}^{n_p}} f > -\infty$.

For a function f that satisfies these assumptions, it is shown [4] that the PAM scheme generates a critical point of f .

THEOREM 2.3. [4, Theorem 6.2] *Suppose that f is a K-Ł function of the form (2.11). Let $\{x^k\}_{k \in \mathbb{N}}$ be a sequence generated by [4, Algorithm 4]. If the sequence $\{x^k\}_{k \in \mathbb{N}}$ is bounded, then the following assertions hold:*

- (i) *The sequence $\{x^k\}_{k \in \mathbb{N}}$ has finite length, i.e., $\sum_{k=1}^{\infty} \|x^{k+1} - x^k\|_F < \infty$.*
- (ii) *The sequence $\{x^k\}_{k \in \mathbb{N}}$ converges to a critical point \bar{x} of f .*

Indeed, Algorithm 2 is a specific case of the PAM method [4].

REMARK 2.4 (K-Ł Property). *The global convergence of the PAM method established in [4, Theorem 6.2] requires the objective function f to satisfy the Kurdyka-Łojasiewicz (K-Ł) property on its effective domain; see [8, Definition 3] for more details on the K-Ł property. It is shown [8, Definition 5] that the so-called semi-algebraic functions satisfy the K-Ł property. Indeed, all functions in (2.1) are semi-algebraic functions, which includes the ℓ_1 -norm $\|x\|_1$, the quadratic functions $x^\top Hx$ and $\|Ax - b\|_2^2$, and δ_S , the indicator function of the Stiefel manifold S . Since a finite sum of semi-algebraic functions is also semi-algebraic, the objective function in (2.1) also satisfies the K-Ł property.*

Define $W := (X, P, Q)$. For the k -th iteration, then the augmented Lagrangian (2.2) can be expressed as

$$L_k(W) = L(X, Q, P, \bar{\Lambda}^k; \rho^k) = f_1(X) + f_2(Q) + f_3(P) + g_k(X, Q, P), \quad (2.12)$$

where

$$\begin{cases} f_1(X) := \text{Tr}(X^\top HX), & f_2(Q) := \frac{1}{\mu} \|Q\|_1, & f_3(P) := \delta_S(P), \\ g_k(X, Q, P) := \langle \bar{\Lambda}_1^k, Q - X \rangle + \frac{\rho^k}{2} \|Q - X\|_F^2 \\ \quad + \langle \bar{\Lambda}_2^k, P - X \rangle + \frac{\rho^k}{2} \|P - X\|_F^2. \end{cases}$$

Then, the following result shows that Step 1 of Algorithm 1 is well defined.

PROPOSITION 2.5. *For each $k \in \mathbb{N}$, denote the function given by (2.12) by L_k , and denote the sequence generated by Algorithm 2 by $\{(X^{k,j}, Q^{k,j}, P^{k,j})\}_{j \in \mathbb{N}}$. Then, $\Theta^{k,j}$ defined by (2.8) satisfies*

$$\Theta^{k,j} \in \partial L(X^{k,j}, Q^{k,j}, P^{k,j}, \bar{\Lambda}^k; \rho^k), \quad \forall j \in \mathbb{N}.$$

If the parameters γ, ρ^1 in Algorithm 1 are chosen such that

$$\gamma > 1, \quad \rho^1 > 0, \quad \rho^1 I_n + 2H \succ 0, \quad (2.13)$$

then for each $k \in \mathbb{N}$,

$$\|\Theta^{k,j}\|_\infty \rightarrow 0 \quad \text{as } j \rightarrow \infty.$$

Proof. To establish the first part of this proposition, recall the functions g_k, f_1, f_2, f_3 , as defined by L_k in (2.12). Then, a direct calculation shows that $\Theta^{k,j} = (\Theta_1^{k,j}, \Theta_2^{k,j}, \Theta_3^{k,j})$

defined by (2.8) can be expressed in terms of partial derivatives of $g := g_k$ as

$$\begin{aligned}\Theta_1^{k,j} &= -\nabla_X g(X^{k,j}, Q^{k,j-1}, P^{k,j-1}) - c_1^{k,j-1}(X^{k,j} - X^{k,j-1}) + \nabla_X g(X^{k,j}, Q^{k,j}, P^{k,j}); \\ \Theta_2^{k,j} &= -\nabla_Q g(X^{k,j}, Q^{k,j}, P^{k,j-1}) - c_2^{k,j-1}(Q^{k,j} - Q^{k,j-1}) + \nabla_Q g(X^{k,j}, Q^{k,j}, P^{k,j}); \\ \Theta_3^{k,j} &= -\nabla_P g(X^{k,j}, Q^{k,j}, P^{k,j}) - c_3^{k,j-1}(P^{k,j} - P^{k,j-1}) + \nabla_P g(X^{k,j}, Q^{k,j}, P^{k,j}).\end{aligned}\quad (2.14)$$

On the other hand, given $(X^{k,j-1}, Q^{k,j-1}, P^{k,j-1})$, using [34, 8.8(c)], the PAM scheme (2.7) yields the following necessary first order optimality condition:

$$\begin{cases} \nabla f_1(X^{k,j}) + \nabla_X g(X^{k,j}, Q^{k,j-1}, P^{k,j-1}) + c_1^{k,j-1}(X^{k,j} - X^{k,j-1}) = 0; \\ \nu^{k,j} + \nabla_Q g(X^{k,j}, Q^{k,j}, P^{k,j-1}) + c_2^{k,j-1}(Q^{k,j} - Q^{k,j-1}) = 0; \\ \omega^{k,j} + \nabla_P g(X^{k,j}, Q^{k,j}, P^{k,j}) + c_3^{k,j-1}(P^{k,j} - P^{k,j-1}) = 0, \end{cases}\quad (2.15)$$

where $\nu^{k,j} \in \partial f_2(Q^{k,j})$ and $\omega^{k,j} \in \partial f_3(P^{k,j})$. Replacing the corresponding terms in (2.14) by (2.15) gives

$$\begin{cases} \Theta_1^{k,j} = \nabla f_1(X^{k,j}) + \nabla_X g(X^{k,j}, Q^{k,j}, P^{k,j}) \in \partial_X L_k(W^{k,j}); \\ \Theta_2^{k,j} = \nu^{k,j} + \nabla_Q g(X^{k,j}, Q^{k,j}, P^{k,j}) \in \partial_Q L_k(W^{k,j}); \\ \Theta_3^{k,j} = \omega^{k,j} + \nabla_P g(X^{k,j}, Q^{k,j}, P^{k,j}) \in \partial_P L_k(W^{k,j}). \end{cases}$$

By [3, Proposition 3] and [34], we have

$$\partial L_k(W) = \partial_X L_k(W) \times \partial_Q L_k(W) \times \partial_P L_k(W)$$

which implies that for each $k \in \mathbb{N}$,

$$\Theta^{k,j} \in \partial L(X^{k,j}, Q^{k,j}, P^{k,j}, \bar{\Lambda}^k; \rho^k), \quad \forall j \in \mathbb{N}.$$

For the second part of the proposition, since $\Theta^{k,j}$ is explicitly given by (2.8), to prove for each $k \in \mathbb{N}$, $\|\Theta^{k,j}\|_\infty \rightarrow 0$, as $j \rightarrow \infty$, it suffices to show that for each $k \in \mathbb{N}$, the sequence $\{(X^{k,j}, Q^{k,j}, P^{k,j})\}_{j \in \mathbb{N}}$ is convergent. Then it remains to verify that the functions $\{L_k(W)\}_{k \in \mathbb{N}}$ satisfy the conditions and assumptions made in Theorem 2.3. From its definition (2.12), it can be seen that the function L_k satisfies the assumptions (i) and (ii) of the function given by (2.11), and L_k is also a K-Ł function according to Remark 2.4. Thus, we only need to verify that each $k \in \mathbb{N}$, L_k is bounded below and the sequence $\{W^{k,j}\}_{j \in \mathbb{N}}$ is bounded.

For each $k \in \mathbb{N}$, the lower bound of L_k is proved by showing that L_k is a coercive function (i.e., $L_k(W) \rightarrow +\infty$, when $\|W\|_\infty \rightarrow \infty$) provided that the parameters γ, ρ^1 satisfy (2.13). Clearly, the two terms f_2 and f_3 of L_k in (2.12) are coercive. We may rewrite the remaining terms

$$f_1(X) + g_k(X, Q, P) := g_{1,k}(X, P) + g_{2,k}(X, Q), \quad (2.16)$$

where

$$\begin{cases} g_{1,k}(X, P) := \frac{1}{2} \text{Tr}(X^\top (2H + \rho^k I_n) X) - \langle \rho^k P + \bar{\Lambda}_2^k, X \rangle + \langle \bar{\Lambda}_2^k, P \rangle + \frac{\rho^k}{2} \|P\|_F^2, \\ g_{2,k}(X, Q) := \frac{\rho^k}{2} \|Q - X + \bar{\Lambda}_1^k / \rho^k\|_F^2 - \|\bar{\Lambda}_1^k / \rho^k\|_F^2. \end{cases}$$

It can be seen that $g_{2,k}$ is bounded below. Since $P \in \mathcal{S}$ (i.e., $P^\top P = I_m$), so $\|P\|_\infty = 1$ and $\|P\|_F = \sqrt{m}$. Thus we have

$$g_{1,k}(X, P) \geq \frac{1}{2} \text{Tr}(X^\top (2H + \rho^k I_n) X) - \|X\|_1 - \langle \bar{\Lambda}_2^k, X \rangle - \|\bar{\Lambda}_2^k\|_1 + \frac{\rho^k m}{2}.$$

Therefore, $g_{1,k}$ is coercive as long as $2H + \rho^k I_n$ is positive definite and $P \in \mathcal{S}$. Notice that the sequence $\{\rho^k\}_{k \in \mathbb{N}}$ is set in Step 3 of Algorithm 1 such that it is non-decreasing when $\gamma > 1$ which implies $\rho^k \geq \rho^1$ for any $k > 1$. If the initial parameter ρ^1 is set sufficiently large such that $2H + \rho^1 I_n \succ 0$, we have the positive definiteness of $2H + \rho^k I_n$ for any $k \geq 1$ and thus the term $f_1 + g_k$ is also coercive. In short, the functions $\{L_k\}_{k \in \mathbb{N}}$ defined by (2.12) are all coercive.

The boundedness of the sequence $\{W^{k_0,j}\}_{j \in \mathbb{N}}$ is proved by contradiction. Suppose on the contrary that the sequence $\{W^{k_0,j}\}_{j \in \mathbb{N}}$ is not bounded, and so $\lim_{j \rightarrow \infty} \|W^{k_0,j}\| = \infty$. As $L_{k_0}(W)$ is a coercive function, we have then $\lim_{j \rightarrow \infty} L_{k_0}(W^{k_0,j}) = +\infty$. However, according to [3, Lemma 3], we have that

$$L_{k_0}(W^{k_0,j+1}) + c\|W^{k_0,j+1} - W^{k_0,j}\|_F^2 \leq L_{k_0}(W^{k_0,j}), \quad j \in \mathbb{N},$$

which implies that $\{L_{k_0}(W^{k_0,j})\}_{j \in \mathbb{N}}$ is a non-increasing sequence, leading to a contradiction. This completes the proof. \square

REMARK 2.6. *The condition $2H + \rho^1 I_n \succ 0$ is a mild condition in the instance of the compressed modes problem [31]. In the case of the free-electron (FE) model, the discrete Hamiltonian $H \succeq 0$ and thus ρ^1 can be taken to be any positive number. In the case of the Kronig-Penney (KP) model, the magnitudes of the negative eigenvalues of the corresponding matrix H are generally less than 1. Thus, we may set $\rho^1 > 2$.*

3. Convergence analysis. For the convenience of notation and discussion, we rewrite the problem (2.1) using the notation of vectors. Let $x \in \mathbb{R}^{3mn}$ denote the column vector formed by concatenating the columns of X, Q, P , i.e.,

$$x := \text{Vec}([X|Q|P]). \quad (3.1)$$

Then, the problem (2.1) can be rewritten as the following:

$$\min_{x \in \mathbb{R}^{3mn}} f(x), \quad \text{subject to} \quad h_1(x) = 0; \text{ and } h_2(x) = 0; \quad (3.2)$$

where $h_1(x) \in \mathbb{R}^{2mn}$, denotes $\text{Vec}([Q - X|P - X])$, $h_2(x)$ denotes the $\frac{m(m+1)}{2} \times 1$ vector obtained by vectorizing only the lower triangular part of the symmetric matrix $P^\top P - I$ and $f(x) := \sum_{j=1}^m (\mu^{-1} \|Q_j\|_1 + X_j^\top H X_j)$. Let λ denote the concatenation of the two Lagrange multiplier vectors of Λ_1 and Λ_2 given by $\lambda := \text{Vec}([\Lambda_1|\Lambda_2])$. Then, the corresponding augmented Lagrangian of (3.2) can be expressed as

$$L(x, \lambda; \rho) := f(x) + \sum_{i=1}^{m_1} [\lambda]_i [h_1(x)]_i + \frac{\rho}{2} \sum_{i=1}^{m_2} [h_1(x)]_i^2, \quad \text{subject to} \quad x \in \Gamma,$$

where $m_1 := 2mn, m_2 := m(m+1)/2$, and

$$\Gamma := \{x : h_2(x) = 0\}. \quad (3.3)$$

Thus, $(X^*, Q^* P^*)$ is a KKT point for (2.1) if and only if the vector x^* defined by (3.1) is a KKT point for (3.2), i.e., there exists $w^* \in \partial f(x^*), \lambda^* \in \mathbb{R}^{m_1}, v^* \in \mathbb{R}^{m_2}$ such that

$$w^* + \sum_{i=1}^{m_1} [\lambda^*]_i \nabla [h_1(x^*)]_i + \sum_{i=1}^{m_2} [v^*]_i \nabla [h_2(x^*)]_i = 0; h_1(x^*) = 0; \text{ and } h_2(x^*) = 0. \quad (3.4)$$

In this section, we establish the sub-sequence convergence property of the PAMAL method, i.e., there exists at least one convergent sub-sequence of the sequence generated by Algorithm 1 and it converges to a KKT point of (3.2).

THEOREM 3.1. *Suppose that the positive parameters γ, ρ^1 in Algorithm 1 are chosen so that $\gamma > 1, 2H + \rho^1 I_n \succ 0$. Let $\{(X^k, Q^k, P^k)\}_{k \in \mathbb{N}}$ be the sequence generated by Algorithm 1. Then, the limit point set of $\{(X^k, Q^k, P^k)\}_{k \in \mathbb{N}}$ is non-empty, and every limit point is a KKT point of the original problem (2.1).*

Proof. [Sketch of the proof] The proof of the sub-sequence convergence property of the PAMAL method is organized as follows. Firstly, in Section 3.1, we establish a crucial ingredient needed by the convergence analysis, namely the linear independence of the gradient vectors $\{\nabla[h_1(x)]_i\}_{i=1}^{m_1} \cup \{\nabla[h_2(x)]_i\}_{i=1}^{m_2}$ when $x \in \Gamma$. Consequently, any locally optimal solution to (3.2) is necessarily a KKT point of (3.2). Secondly, in Section 3.2, we show that any limit point of a sequence generated by Algorithm 1 is also a KKT point of (2.1). Lastly, in Section 3.3, we show that for (3.2), the sequence $\{(X^k, Q^k, P^k)\}_{k \in \mathbb{N}}$ generated by Algorithm 1 must be bounded. These results together establish the sub-sequence convergence property of the PAMAL method. \square

3.1. Linear Independence and KKT first-order necessary conditions. It is noted that the objective function f in (3.2) is merely Lipschitz continuous on bounded sets, which is equivalent to the notion of strict continuity (see [34, Definition 9.1]). In order to establish that a locally optimal solution satisfies the KKT first-order necessary conditions in the non-smooth case, we need to invoke the following result.

THEOREM 3.2. [34, Exercise 10.52] (*non-smooth Lagrange multiplier rule*). *For a nonempty, closed set $\mathcal{X} \in \mathbb{R}^n$ and strictly continuous functions $f_0 : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ with $F = (f_1, \dots, f_m)$, consider the problem*

$$\min_{x \in \mathcal{X}} f_0(x) + \theta(F(x)),$$

where $\theta : \mathbb{R}^m \rightarrow \mathbb{R}$ is proper, lower semi-continuous and convex with effective domain D . Suppose \bar{x} is a locally optimal solution at which the following constraint qualification is satisfied:

$$0 \in \partial(yF)(\bar{x}) + N_{\mathcal{X}}(\bar{x}), \quad y \in N_D(F(\bar{x})) \implies y = 0,$$

where $y \in \mathbb{R}^m$ and $yF := \sum_{i=1}^m y_i f_i$. Then there exists a vector \bar{y} such that

$$0 \in \partial(f_0 + \bar{y}F)(\bar{x}) + N_{\mathcal{X}}(\bar{x}), \quad \bar{y} \in \partial\theta(F(\bar{x})).$$

Moreover, the set of such vectors \bar{y} is compact.

Before applying the above result, we first show that the gradient vectors of h_1, h_2 defined in (3.2) satisfy a linear independence constraint qualification whenever $x \in \Gamma$, i.e., it satisfies the orthogonality constraints. This leads to the KKT first order necessary conditions.

LEMMA 3.3. *Suppose that $x \in \Gamma$. Then, $\{\nabla[h_1(\bar{x})]_i\}_{i=1}^{m_1} \cup \{\nabla[h_2(\bar{x})]_i\}_{i=1}^{m_2}$ are linearly independent, where h_1 and h_2 are defined in (3.2). Consequently, if \bar{x} is a locally optimal solution of the problem (3.2), then \bar{x} is a KKT point for (3.2).*

Proof. From the definition (3.1) of x , by setting $m_3 := mn$, it can be seen that

$$\nabla h_1(x) = \begin{bmatrix} -I_{m_3 \times m_3} & -I_{m_3 \times m_3} \\ I_{m_3 \times m_3} & 0_{m_3 \times m_3} \\ 0_{m_3 \times m_3} & I_{m_3 \times m_3} \end{bmatrix}, \quad \text{and} \quad \nabla h_2(x) = \begin{bmatrix} 0_{m_3 \times m_2} \\ 0_{m_3 \times m_2} \\ M(x) \end{bmatrix}, \quad (3.5)$$

where $M(x) \in \mathbb{R}^{m_3 \times m_2}$ given by

$$M(x) = \begin{bmatrix} 2P_1 & P_2 & P_3 & \dots & P_m & 0 & 0 & \dots & 0 & 0 & 0 & 0 \\ 0 & P_1 & 0 & \dots & 0 & 2P_2 & P_3 & \dots & P_m & \vdots & \vdots & \vdots \\ 0 & 0 & P_1 & \ddots & \vdots & 0 & P_2 & \dots & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & 0 & \vdots & \ddots & \ddots & \vdots & 2P_{m-1} & P_m & 0 \\ 0 & 0 & \dots & 0 & P_1 & 0 & \dots & 0 & P_2 & 0 & P_{m-1} & 2P_m \end{bmatrix}.$$

Since $x \in \Gamma$, the column vectors $\{P_i\}_{i=1}^m$ are orthogonal to each other, and thus the columns of $M(x)$ are orthogonal to each other. Furthermore, the first $2m_3$ rows of $\nabla h_2(x)$ form a zero matrix. Thus, from the matrix structures of $\nabla h_1(x)$, $\nabla h_2(x)$ in (3.5), one can easily show that $\{\nabla[h_1(x)]_i\}_{i=1}^{m_1} \cup \{\nabla[h_2(x)]_i\}_{i=1}^{m_2}$ are linearly independent for any $x \in \Gamma$.

Secondly, if \bar{x} is a locally optimal solution of the problem (3.2), then $\bar{x} \in \Gamma$. It can be seen from the arguments above that $\nabla h_2(\bar{x})$ is of full column rank. Furthermore, applying (1.8) in Remark 1.3 on the smooth function h_2 leads to

$$N_\Gamma(\bar{x}) = \{\nabla h_2(\bar{x})z \mid z \in \mathbb{R}^{m_2}\} = \left\{ \sum_{i=1}^{m_2} [z]_i \nabla[h_2(\bar{x})]_i \mid z \in \mathbb{R}^{m_2} \right\}. \quad (3.6)$$

A direct calculation shows that the constraint qualification for (3.2) amounts to verifying

$$0 \in \nabla h_1(\bar{x})y + \nabla h_2(\bar{x})z, y \in \mathbb{R}^n \implies y = 0,$$

which holds true by the linear independence of $\{\nabla[h_1(\bar{x})]_i\}_{i=1}^{m_1} \cup \{\nabla[h_2(\bar{x})]_i\}_{i=1}^{m_2}$ as $\bar{x} \in \Gamma$. Notice that $\delta_{\{0\}}$ is proper, lower semi-continuous and convex with effective domain $D = \{0\}$. Then, by applying Theorem 3.2 with the setting: $f_0 := f, \theta := \delta_{\{0\}}, F := h_1$ and $\mathcal{X} := \Gamma$, we established (3.2). Together with (3.6), we have the existence of vectors $\bar{w} \in \partial f(\bar{x}), \bar{y} \in \mathbb{R}^{m_1}, \bar{z} \in \mathbb{R}^{m_2}$ such that

$$\bar{w} + \sum_{i=1}^{m_1} [\bar{y}]_i \nabla[h_1(\bar{x})]_i + \sum_{i=1}^{m_2} [\bar{z}]_i \nabla[h_2(\bar{x})]_i = 0.$$

In other words, the locally optimal point \bar{x} is also a KKT point of (3.2). \square

3.2. Limit points as KKT points. In this section, we show that any limit point generated by Algorithm 1 is also a KKT point of (2.1), i.e., any limit point x^* of the corresponding sequence $\{x^k\}_{k \in \mathbb{N}}$ with respect to (X^k, Q^k, P^k) is a KKT point for (3.2). Recall that the normal cone $\partial \delta_S(X, Q, P) = N_S(X, Q, P)$ in vector notation is given by (3.6). Thus, in vector notation, finding the solution satisfying the constraint (2.3) at Step 1 of Algorithm 1 is equivalent to calculating a solution x^k such that there exist vectors $\omega^k \in \partial f(x^k)$ and v^k which satisfy

$$\|\omega^k + \sum_{i=1}^{m_1} ([\bar{\lambda}^k]_i + \rho^k [h_1(x^k)]_i) \nabla[h_1(x^k)]_i + \sum_{i=1}^{m_2} [v^k]_i \nabla[h_2(x^k)]_i\|_2 \leq \epsilon^k \quad (3.7)$$

with $h_2(x^k) = 0, k \in \mathbb{N}$.

REMARK 3.4. *Algorithm 1 can recast as an equality constrained version of [2, Algorithm 3.1] in vector notation. However, we cannot directly apply the results in [2, Theorems 4.1–4.2] on our problem, as our objective function f defined by (3.2) is not in C^1 .*

In vector notation, the main result is stated as follows.

THEOREM 3.5. *Suppose $\{x^k\}_{k \in \mathbb{N}}$ is a sequence generated by Algorithm 1. Let x^* be a limit point of this sequence, i.e., there exists a sub-sequence $\mathcal{K} \subseteq \mathbb{N}$ such that $\lim_{k \in \mathcal{K}} x^k = x^*$. Then x^* is also a KKT point of (3.2).*

Proof. The proof consists of two main parts. The first part shows that x^* is a feasible point of (3.2), i.e., $h_1(x^*) = 0$ and $h_2(x^*) = 0$. The second part shows that x^* satisfies the remaining KKT property in (3.4).

We start with the proof of the feasibility of x^* for h_2 . After running Step 1 of Algorithm 1, we obtain $h_2(x^k) = 0$ for all $k \in \mathcal{K}$, therefore, $h_2(x^*) = 0$, i.e., $x^* \in \Gamma$. The next step is to show $h_1(x^*) = 0$, which is discussed in two cases.

We now prove the first case where the sequence $\{\rho^k\}_{k \in \mathbb{N}}$ is bounded. Recall that $\gamma > 1$ in Algorithm 1. Thus, the update rule (2.4) on ρ^k in Step 3 of Algorithm 1 suggests from some iteration k_0 onwards, the penalty parameter ρ^k will remain the same, which implies that $\|h_1(x^{k+1})\|_\infty \leq \tau \|h_1(x^k)\|_\infty$, $k \geq k_0$ for some constant $\tau \in [0, 1)$. The feasibility $h_1(x^*) = 0$ is then proved.

In the other case where the sequence $\{\rho^k\}_{k \in \mathbb{N}}$ is not bounded, for each $k \in \mathcal{K}$, there exist vectors $\{\delta^k\}_{k \in \mathbb{N}}$ with $\|\delta^k\|_\infty \leq \epsilon^k$ and $\epsilon^k \downarrow 0$ such that

$$w^k + \sum_{i=1}^{m_1} ([\bar{\lambda}^k]_i + \rho^k [h_1(x^k)]_i) \nabla [h_1(x^k)]_i + \sum_{i=1}^{m_2} [v^k]_i \nabla [h_2(x^k)]_i = \delta^k, \quad (3.8)$$

for some $w^k \in \partial f(x^k)$. Dividing both sides of (3.8) by ρ^k , we have

$$\sum_{i=1}^{m_1} ([\bar{\lambda}^k / \rho^k]_i + [h_1(x^k)]_i) \nabla [h_1(x^k)]_i + \sum_{i=1}^{m_2} [\hat{v}^k]_i \nabla [h_2(x^k)]_i = \frac{\delta^k - w^k}{\rho^k}, \quad (3.9)$$

where $\hat{v}^k = (\rho^k)^{-1} v^k$. Define

$$\Xi(x)^\top := [\nabla h_1(x) \quad \nabla h_2(x)]$$

and

$$\eta^k := ([\bar{\lambda}^k / \rho^k]_1 + [h_1(x^k)]_1, \dots, [\bar{\lambda}^k / \rho^k]_{m_1} + [h_1(x^k)]_{m_1}, [\hat{v}^k]_1, \dots, [\hat{v}^k]_{m_2})^\top.$$

Then the equality (3.9) can be re-written as

$$\Xi(x^k)^\top \eta^k = (\delta^k - w^k) / \rho^k. \quad (3.10)$$

By Lemma 3.3, $\{\nabla [h_1(x^*)]_i\}_{i=1}^{m_1} \cup \{\nabla [h_2(x^*)]_i\}_{i=1}^{m_2}$ are linearly independent as $x^* \in \Gamma$. Moreover, the gradient vectors $\nabla h_1, \nabla h_2$ are continuous and $h_2(x^k) = 0$ for all $k \in \mathcal{K}$. Note that by the continuity of the gradient vectors ∇h_1 and ∇h_2 , $\Xi(x^k) \rightarrow \Xi(x^*)$, which has full rank as $x^* \in \Gamma$. Thus, $\Xi(x^k) \Xi(x^k)^\top \rightarrow \Xi(x^*) \Xi(x^*)^\top \succ 0$.

It is known that the eigenvalues of a symmetric matrix vary continuously with its matrix values. Thus, we deduce that $\Xi(x^k) \Xi(x^k)^\top$ is non-singular for sufficiently large $k \in \mathcal{K}$, which leads to

$$\eta^k = [\Xi(x^k) \Xi(x^k)^\top]^{-1} \Xi(x^k) (\delta^k - w^k) / \rho^k.$$

Since f is the summation of a convex function and a continuously differentiable function, the set $\bigcup_{x \in \mathcal{M}} \partial f(x)$ is a bounded set whenever \mathcal{M} is bounded. This can be seen by invoking [6, Proposition B.24 (b)]. Here, we set $\mathcal{M} = \{x^k\}_{k \in \mathcal{K}}$ which is clearly a bounded set. Thus,

we have $\{w^k\}_{k \in \mathcal{K}}$ is bounded. Together with $\|\delta^k\| \leq \epsilon^k \downarrow 0$, taking limits as $k \in \mathcal{K}$ goes to infinity gives $\eta^k \rightarrow 0$. The boundedness of the safeguard Lagrange multipliers $\{\bar{\lambda}^k\}_k$ implies that $[h_1(x^*)]_i = 0 = [\hat{v}]_j$ for all i, j . Thus, $h_1(x^*) = 0$ and this ends the first part of the proof.

Next, we will show that x^* is a KKT point of the problem (3.2). Since $\{w^k\}_{k \in \mathcal{K}}$ is bounded, there exists a sub-sequence $\mathcal{K}_2 \subseteq \mathcal{K}$ such that $\lim_{k \in \mathcal{K}_2} w^k = w^*$. Recall that $\lim_{k \in \mathcal{K}_2} x^k = x^*$ and $w^k \in \partial f(x^k)$. Thus,

$$w^* \in \partial f(x^*),$$

by the closedness property of the limiting sub-differential. Together with $[\lambda^{k+1}]_i = [\bar{\lambda}^k]_i + \rho^k [h_1(x^k)]_i$, for $\forall i$, it can be seen from Algorithm 1 that for $k \in \mathcal{K}_2$,

$$w^k + \sum_{i=1}^{m_1} [\lambda^{k+1}]_i \nabla [h_1(x^k)]_i + \sum_{i=1}^{m_2} [v^k]_i \nabla [h_2(x^k)]_i = \delta^k, \quad (3.11)$$

for some vectors δ^k with $\|\delta^k\|_\infty \leq \epsilon^k \downarrow 0$ and $w^k \in \partial f(x^k)$. Define

$$\pi^k := ([\lambda^{k+1}]_1, \dots, [\lambda^{k+1}]_{m_1}, [v^k]_1, \dots, [v^k]_{m_2})^\top. \quad (3.12)$$

Then (3.11) can be re-written as

$$\Xi(x^k)^\top \pi^k = \delta^k - w^k.$$

By the same arguments in the first part, the matrix $\Xi(x^k) \Xi(x^k)^\top$ is non-singular for sufficiently large $k \in \mathcal{K}_2$ and

$$\pi^k = [\Xi(x^k) \Xi(x^k)^\top]^{-1} \Xi(x^k) (\delta^k - w^k). \quad (3.13)$$

Hence, by taking limits on (3.13) as $k \in \mathcal{K}_2$ goes to infinity, we have

$$\pi^k \rightarrow \pi^* = -[\Xi(x^*) \Xi(x^*)^\top]^{-1} \Xi(x^*) w^*.$$

By the definition (3.12) of π^k , taking limits as $k \in \mathcal{K}_2$ approaches infinity on both sides of (3.11) leads to

$$w^* + \sum_{i=1}^{m_1} [\lambda^*]_i \nabla [h_1(x^*)]_i + \sum_{i=1}^{m_2} [v^*]_i \nabla [h_2(x^*)]_i = 0,$$

where λ^*, v^* are obtained from π^* similar to (3.12). Thus x^* is a KKT point of (3.2) and this completes the second part of the proof. \square

3.3. Existence of Limit points. The results presented in the previous sections assume the existence of a limit point of the sequence $\{x_k\}_{k \in \mathbb{N}}$, i.e., the sequence generated by Algorithm 1 contains at least one convergent sub-sequence. We now prove this existence by showing that the generated sequence is bounded.

PROPOSITION 3.6. *Let $\{(X^k, Q^k, P^k)\}_{k \in \mathbb{N}}$ be the sequence generated by Algorithm 1. Suppose that the parameters γ, ρ^1 in Algorithm 1 are chosen so that $\gamma > 1$ and $2H + \rho^1 I_n \succ 0$. Then, $\{(X^k, Q^k, P^k)\}_{k \in \mathbb{N}}$ is bounded and thus contains at least one convergent sub-sequence.*

Proof. The boundedness of $\{P^k\}_{k \in \mathbb{N}}$ is easy to see from Step 1 of Algorithm 1. It remains to show that $\{(X^k, Q^k)\}_{k \in \mathbb{N}}$ is bounded. Using a direct extension of the result

[3, Proposition 3], the first two partial sub-differentials of L in (2.3) yield the following: there exist $\nu^k \in \frac{1}{\mu} \partial \|Q^k\|_1$ and $\zeta^k = (\zeta_1^k, \zeta_2^k) \in \mathbb{R}^{n \times m} \times \mathbb{R}^{n \times m}$ such that

$$\begin{cases} \zeta_1^k = 2HX^k - \bar{\Lambda}_1^k + \rho^k(X^k - Q^k) - \bar{\Lambda}_2^k + \rho^k(X^k - P^k); \\ \zeta_2^k = \nu^k + \bar{\Lambda}_1^k + \rho^k(-X^k + Q^k), \end{cases} \quad (3.14)$$

where $\|\zeta^k\|_\infty \leq \epsilon^k$. Summing the above two equations gives

$$(2H + \rho^k)X^k = \zeta_1^k + \zeta_2^k + \bar{\Lambda}_2^k + \rho^k P^k - \nu^k.$$

Together with $2H + \rho^k I_n \succ 0$, we have

$$X^k = (2H + \rho^k I_n)^{-1}(\zeta_1^k + \zeta_2^k + \bar{\Lambda}_2^k + \rho^k P^k - \nu^k).$$

Let $H := V \text{diag}(\lambda_1, \dots, \lambda_n) V^\top$ denote the SVD of the symmetric matrix H . Then,

$$\begin{aligned} X^k &= V \text{diag}(1/(2\lambda_1 + \rho^k), \dots, 1/(2\lambda_n + \rho^k)) V^\top (\zeta_1^k + \zeta_2^k + \bar{\Lambda}_2^k - \nu^k) \\ &\quad + V \text{diag}(\rho^k/(2\lambda_1 + \rho^k), \dots, \rho^k/(2\lambda_n + \rho^k)) V^\top P^k. \end{aligned} \quad (3.15)$$

Recall that $\{\rho^k\}_{k \in \mathbb{N}}$ is non-decreasing and $2H + \rho^1 I_n \succ 0$. We have then, for $k \in \mathbb{N}$, $2H + \rho^k I_n \succ 0$, which gives $2\lambda_i + \rho^k > 0$, $i = 1, \dots, n$. Thus, for all $k \in \mathbb{N}$,

$$\begin{aligned} 0 &< 1/(2\lambda_i + \rho^k) \leq 1/(2\lambda_i + \rho^1) < +\infty, \quad i = 1, \dots, n, \\ 0 &< \rho^k/(2\lambda_i + \rho^k) \leq \max(\rho^1/(2\lambda_i + \rho^1), 1), \quad i = 1, \dots, n. \end{aligned}$$

Together with the fact that $\{\zeta^k\}_{k \in \mathbb{N}}$, $\{\bar{\Lambda}^k\}_{k \in \mathbb{N}}$ and $\{\nu^k\}_{k \in \mathbb{N}}$ are bounded, combining the two inequalities (3.15) and (3.16) shows that the sequence $\{X^k\}_{k \in \mathbb{N}}$ is bounded. Then, the boundedness of the sequence $\{Q^k\}_{k \in \mathbb{N}}$ can be also derived from (3.14). \square

It is noted that Proposition 3.6 still holds if the ℓ_1 -term $\frac{1}{\mu} \|Q\|_1$ in (2.1) is replaced by any convex function with bounded sub-gradients on its domain, e.g., $\frac{1}{\mu} \|Q\|_{2,1}$.

4. The compressed modes for variational problems in physics. This section is organized as follows. In subsection 4.1, we present some background information on compressed modes. In subsection 4.2, we review some existing methods to obtain compressed modes, which include the SOC method [23] introduced by Lai and Osher. Finally, we compare the numerical performance of the PAMAL method against that of the SOC method in subsection 4.3 on the compressed modes problem.

4.1. Background on compressed modes. Motivated by the localized Wannier functions [26] used in solid state physics and quantum chemistry, a variational approach is developed in [31] to produce the so-called *compressed modes*, which are spatially localized solutions to the time-independent Schrödinger's equation:

$$\hat{H}\phi(x) = \lambda\phi(x), \quad x \in \Omega. \quad (4.1)$$

In (4.1), Ω is a bounded subset of \mathbb{R}^d and \hat{H} denotes the corresponding Hamiltonian

$$\hat{H} = -\frac{1}{2}\Delta + V,$$

where Δ denotes the Laplacian operator and V denotes the potential energy function, represented by a multiplication operator with a bounded measurable function. Spatially localized solutions to the eigenvalue problem (4.1) not only enable efficient computations related to the

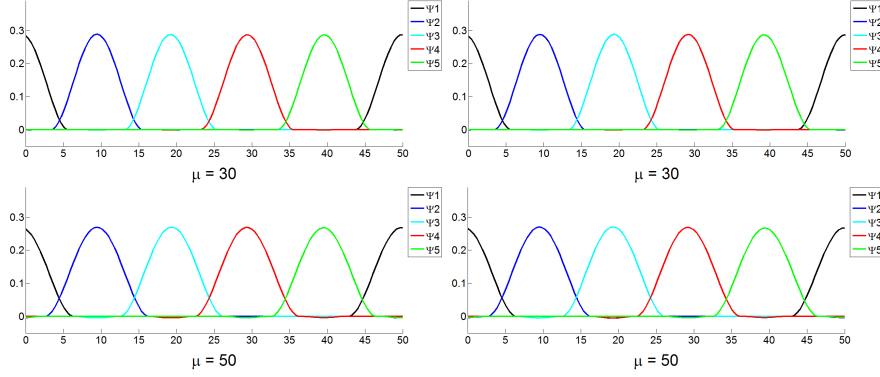


FIG. 4.1. The comparison of the first five modes obtained for the 1D FE model with different values of μ . The first column shows the results computed by the SOC method cited from [31]; the second column shows the results computed by the PAMAL method.

general Schrödinger's equation, but also fit certain observations in physics. For example, the screened correlations in condensed matter are typically short-ranged [33].

In [31], the authors considered the independent-particle Schrödinger's equation for a finite system of N electrons, with the electron spin neglected for simplicity. The ground state energy of these electrons, denoted by E_0 , can be formulated as a variational problem, which minimizes the total energy subject to orthonormality conditions for the stationary states:

$$E_0 = \min_{\Phi_N} \sum_{j=1}^N \langle \phi_j, \hat{H} \phi_j \rangle \quad \text{s.t.} \quad \langle \phi_j, \phi_k \rangle = \delta_{jk}, \quad (4.2)$$

where $\langle \phi_j, \phi_k \rangle := \int_{\Omega} \phi_j(x)^* \phi_k(x) dx$. The solutions $\Phi_N = \{\phi_i\}_{i=1}^N$ form a set of orthonormal eigenfunctions which are usually not spatially localized. Therefore, an ℓ_1 -regularized model is proposed in [31] to obtain the solutions of (4.2) with better spatial localization:

$$E = \min_{\Psi_N} \sum_{j=1}^N \frac{1}{\mu} |\psi_j|_1 + \langle \psi_j, \hat{H} \psi_j \rangle \quad \text{s.t.} \quad \langle \psi_j, \psi_k \rangle = \delta_{jk}, \quad (4.3)$$

where $|\psi_j|_1 := \int_{\Omega} |\psi_j(x)| dx$ and the constant μ is a pre-defined parameter that balances the sparsity and the accuracy of the solution. It is shown [5, 30] that with fixed N , the approximation error of the energy E calculated by (4.3) to the ground state energy E_0 is decreasing as $\mu \rightarrow \infty$. By considering $\Omega = [0, L]^d$ with periodic boundary conditions and equally spaced nodes in each direction, the discretized version of (4.2) is expressed as

$$\Psi_N = \operatorname{argmin}_{\Psi \in \mathbb{R}^{n \times N}} \frac{1}{\mu} \|\Psi\|_1 + \operatorname{Tr}(\Psi^T H \Psi) \quad \text{s.t.} \quad \Psi^T \Psi = I, \quad (4.4)$$

where $\|\Psi\|_1 := \sum_{i,j} |\Psi_{i,j}|$, and H is a symmetric matrix formed by the discretization of the Hamiltonian \hat{H} . The solution $\Psi_N = \{\psi_j\}_{j=1}^N$ forms the first N compressed modes (CMs) of the discretized version of the eigenvector problem (4.2).

4.2. Existing methods for problems related to compressed modes. A scheme based on Bregman iterations, namely the splitting of orthogonality constraints (SOC) method [23]

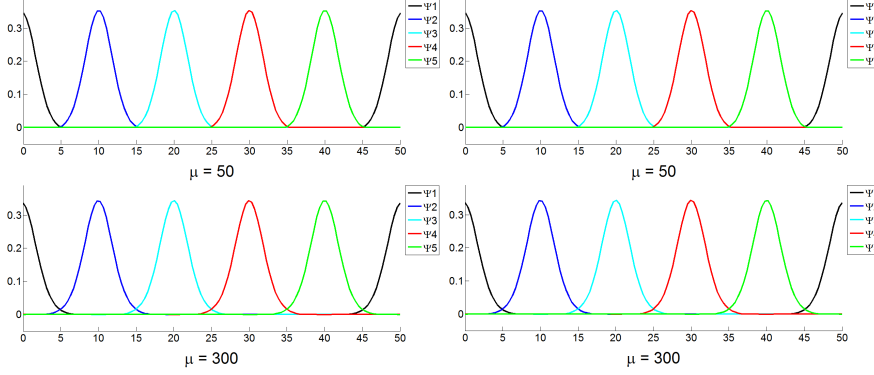


FIG. 4.2. The comparison of the first five modes obtained for the 1D KP model with two different values of μ . The first column shows the results computed by the SOC method [31]; the second column shows the results computed by the PAMAL method.

was used to solve the CMs problem (4.4) by considering the following optimization problem

$$\min_{\Psi, Q, P} \frac{1}{\mu} \|Q\|_1 + \text{Tr}(\Psi^\top H \Psi) \quad \text{s.t.} \quad Q - \Psi = 0, P - \Psi = 0, P^\top P = I.$$

It is demonstrated in the numerical experiments conducted in [31] that the SOC method can produce compressed modes of good quality. Nevertheless, to the best of our knowledge, there is no analysis of its convergence property provided in the literature.

More recently, a convex relaxation approach is proposed in [22], which re-models the CMs problem into a density matrix minimization problem with ℓ_1 -regularization:

$$\begin{aligned} \min_{P \in \mathbb{R}^{n \times n}} \quad & \text{Tr}(HP) + \frac{1}{\mu} \|P\|_1 \\ \text{s.t.} \quad & P = P^\top, \quad \text{Tr}(P) = N, \quad 0 \preceq P \preceq I. \end{aligned} \quad (4.5)$$

In [22], the convex model (4.5) is solved by the split Bregman method, with the convergence analysis provided.

4.3. Computations of CMs via the PAMAL method. In this section, we applied the PAMAL method to solve the CMs problem (4.4) under similar settings given in [31], which considered both the 1D free-electron (FE) and the 1D Kronig-Penney (KP) models. These two models adhere to the same Hamiltonian structure $\hat{H} = -\frac{1}{2}\partial_{x^2} + V(x)$, but differ in the potential energy function V . The FE model describes the behaviour of valence electrons in a crystal structure of a metallic solid, with $V \equiv 0$. The KP model is an idealized quantum-mechanical system that consists of a periodic array of rectangular potential wells, approximated by inverted Gaussians in [31] by setting

$$V = -V_0 \sum_{j=1}^{N_{el}} \exp\left[-\frac{(\cdot - 10j)^2}{2\delta^2}\right],$$

where $V_0 := 1$, $N_{el} := 5$, $\delta := 3$.

In our experiments, the domain $\Omega := [0, 50]$ is discretized with $n = 128$ equally spaced nodes. The parameters of the PAMAL method are set as follows: $\tau = 0.99$, $\gamma = 1.01$, $\rho^1 = 2|\lambda_{\min}(H)| + N/2$, $\bar{\Lambda}_{p,\min} = -100$, $\bar{\Lambda}_{p,\max} = 100$, $\Lambda_p^1 = 0_{n \times N}$, $p = 1, 2$, and

$$\epsilon^k = (0.999)^k, \quad k \in \mathbb{N}.$$

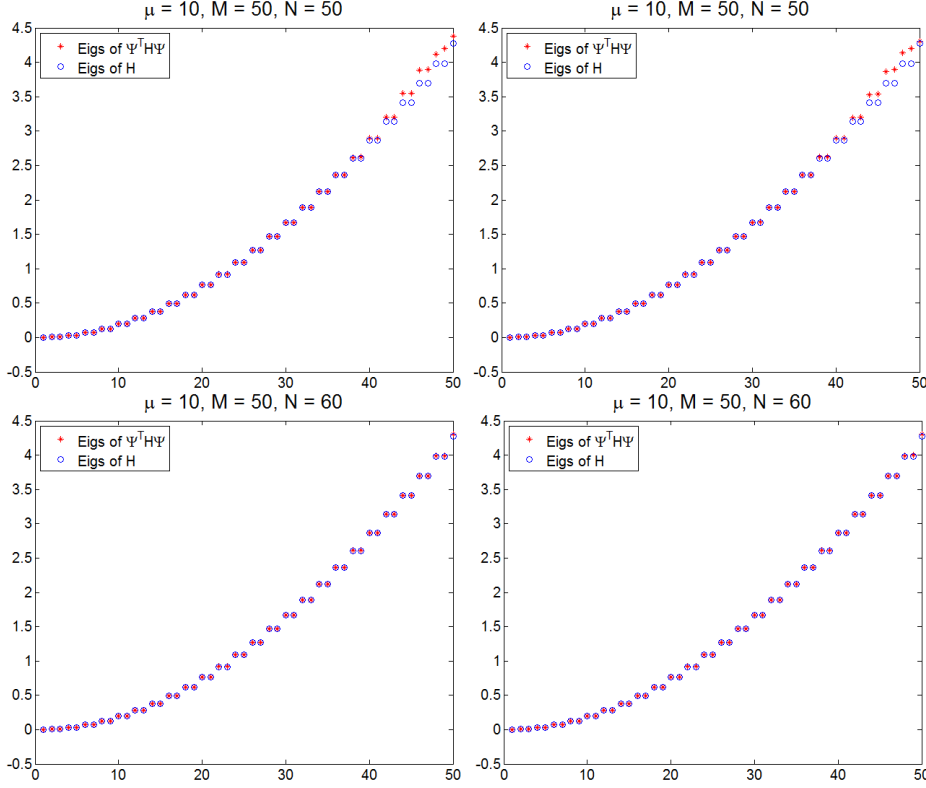


FIG. 4.3. The comparison of the first 50 eigenvalues obtained for the 1D FE model with different values of N . The first column shows the results computed by the SOC method [31]; the second column shows the results computed by the PAMAL method.

The parameters in Algorithm 2 are set as $\underline{c} = c_i^{k,j} = \bar{c} = 0.5$, for all k, j, i in both the FE and KP models. In the SOC method, we use the same penalty parameters ($\lambda = \mu N/20$, $r = \mu N/5$) as recommended by [31, equations 15-17]. In both the PAMAL method and the SOC method, the same random matrix initialization is used. In order to produce CMs of reasonable localization, we set the stopping criterion as $|J(P^k) - J(P^{k-1})| < 10^{-5}$, where J is the objective function given in (4.4), i.e., $J(\Psi) := \frac{1}{\mu} \|\Psi\|_1 + \text{Tr}(\Psi^\top H \Psi)$.

Both methods are implemented in MATLAB and the experiments are done on a PC with a 1.70GHz CPU and 4G of RAM. The number of outer iterations, total number of inner iterations and CPU time, of the PAMAL and SOC methods are averaged over 50 experimental trials. Table 4.1 and Table 4.2 display comparisons of the computational costs of the two methods. In general, with the same stopping criterion, the proposed PAMAL method is at least twice as fast as the SOC method. As discussed in Remark 2.2, the performance gain of the PAMAL methods comes from the flexibility on the accuracy of the solution for Step 1 in Algorithm 1. The first five CMs of the 1D FE and KP models computed by the SOC/PAMAL methods are shown in the first/second columns of Figure 4.1 and Figure 4.2 respectively. It can be seen that the CMs computed by the PAMAL method are compactly supported functions and their localization degree is largely similar to that of the CMs obtained via the SOC method, as shown in Figure 4.1 and Figure 4.2. We next examine the approximation behavior of the unitary transformations derived from the CMs to the eigenmodes of the Schrödinger operator. The approximation accuracy is demonstrated by comparing the first M eigenvalues

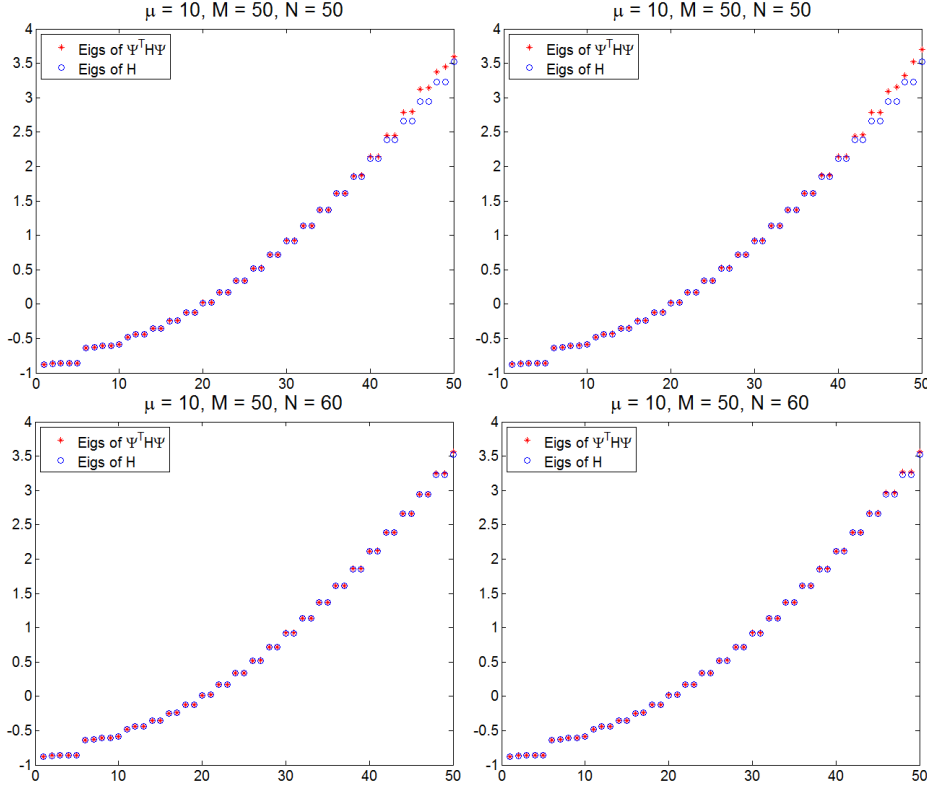


FIG. 4.4. The comparison of the first 50 eigenvalues obtained for the 1D KP model with different values of N . The first column shows the results computed by the SOC method [31]; the second column shows the results computed by the PAMAL method.

Problems			No. of outer iterations		Total no. of inner iterations		CPU time (s)	
N	M	μ	PAMAL	SOC	PAMAL	SOC	PAMAL	SOC
5	5	30	77	237	82	237	0.07	0.15
5	5	50	87	499	92	499	0.07	0.27
50	50	10	512	3124	522	3124	1.35	7.25
60	50	10	484	4147	497	4147	1.54	11.02

TABLE 4.1

Computational costs of the PAMAL method and the SOC method for the FE model.

Problems			No. of outer iterations		Total no. of inner iterations		CPU time (s)	
N	M	μ	PAMAL	SOC	PAMAL	SOC	PAMAL	SOC
5	5	50	66	304	75	304	0.06	0.17
5	5	300	62	1826	71	1826	0.05	0.94
50	50	10	496	3179	507	3179	1.38	7.44
60	50	10	478	4118	491	4118	1.55	10.99

TABLE 4.2

Computational costs of the PAMAL method and the SOC method for the KP model.

$(\sigma_1, \dots, \sigma_M)$ of the matrix $\text{Tr}(\Psi_N^\top H \Psi_N)$ obtained by the M eigenvalues $(\lambda_1, \dots, \lambda_M)$ of the corresponding Schrödinger operators. Figure 4.3 and Figure 4.4 reveal that the approximation accuracies of the SOC and PAMAL methods are similar for the FE and KP models respectively, where it can be seen that $\{\sigma_i\}_{i=1}^m$ converges to $\{\lambda_i\}_{i=1}^m$ with increasing number N of CMs.

5. Conclusion. In this paper, we proposed the PAMAL method, a numerical method for solving a class of ℓ_1 -regularized optimization problems with orthogonality constraints. It is shown in this paper that the proposed method has the sub-sequence convergence property, which is not provided in the existing SOC method [31]. In addition, the experiments show that when applied to solve the compressed modes problem, the proposed PAMAL method is noticeably faster than the SOC method in producing modes of comparable quality.

References.

- [1] P.-A. Absil, R. Mahony, and R. Sepulchre, *Optimization Algorithms on Matrix Manifolds*, Princeton University Press, 2009.
- [2] R. Andreani, E. G. Birgin, J. M. Martínez, and M. L. Schuverdt, *On augmented Lagrangian methods with general lower-level constraints*, SIAM J. Optimiz. **18** (2007), no. 4, 1286–1309.
- [3] H. Attouch, J. Bolte, P. Redont, and A. Soubeyran, *Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the Kurdyka-Łojasiewicz inequality*, Math. Oper. Res. **35** (2010), no. 2, 438–457.
- [4] H. Attouch, J. Bolte, and B. F. Svaiter, *Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward-backward splitting, and regularized Gauss–Seidel methods*, Math. Program. **137** (2013), no. 1–2, 91–129.
- [5] F. Barekat, *On the consistency of compressed modes for variational problems*, arXiv preprint arXiv:1310.4552 (2013).
- [6] D. P. Bertsekas, *Nonlinear Programming*, Athena Scientific, 1999.
- [7] F. Bethuel, H. Brezis, and F. Hélein, *Asymptotics for the minimization of a Ginzburg-Landau functional*, Calc. Var. Partial Dif. **1** (1993), no. 2, 123–148.
- [8] J. Bolte, S. Sabach, and M. Teboulle, *Proximal alternating linearized minimization for nonconvex and nonsmooth problems*, Math. Program. (2013), 1–36.
- [9] J.-F. Cai, S. Osher, and Z. Shen, *Convergence of the linearized Bregman iteration for ℓ^1 -norm minimization*, Math. Comput. **78** (2009), no. 268, 2127–2136.
- [10] E. Candès, J. Romberg, and T. Tao, *Stable signal recovery from incomplete and inaccurate measurements*, Comm. Pure Appl. Math. **8** (2006), no. 59.
- [11] E. J. Candès, X. Li, Y. Ma, and J. Wright, *Robust principal component analysis?*, J. ACM **58** (2011), no. 3, 11.
- [12] E. J. Candès and B. Recht, *Exact matrix completion via convex optimization*, Found. Comput. Math. **9** (2009), no. 6, 717–772.
- [13] T. Chan, A. Marguina, and P. Mulet, *High-order total variation-based image restoration*, SIAM J. Sci. Comput. **22** (2000), no. 2, 503–516.
- [14] D. L. Donoho, *Compressed sensing*, IEEE T. Inform. Theory **52** (2006), no. 4, 1289–1306.
- [15] I. L. Dryden and K. V. Mardia, *Statistical Shape Analysis*, Vol. 4, John Wiley & Sons New York, 1998.
- [16] A. Edelman, T. A. Arias, and S. T. Smith, *The geometry of algorithms with orthogonality constraints*, SIAM J. MATRIX Anal. A. **20** (1998), no. 2, 303–353.
- [17] E. Esser, *Applications of Lagrangian-based alternating direction methods and connections to split Bregman*, CAM Reports **9** (2009), 31.
- [18] M. Fortin and R. Glowinski, *Augmented Lagrangian Methods: Applications to the Numerical Solution of Boundary-value Problems*, Elsevier, 2000.
- [19] R. Glowinski and P. Le Tallec, *Augmented Lagrangian and Operator-splitting Methods in Nonlinear Mechanics*, Vol. 9, SIAM, 1989.
- [20] D. Goldfarb, Z. Wen, and W. Yin, *A curvilinear search method for p -harmonic flows on spheres*, SIAM J. Imaging Sci. **2** (2009), no. 1, 84–109.
- [21] T. Goldstein and S. Osher, *The split Bregman method for L_1 -regularized problems*, SIAM J. Imaging Sci. **2** (2009), no. 2, 323–343.
- [22] R. Lai, J. Lu, and S. Osher, *Density matrix minimization with ℓ_1 regularization*, arXiv preprint arXiv:1403.1525 (2014).
- [23] R. Lai and S. Osher, *A splitting method for orthogonality constrained problems*, J. Sci. Comput. **58** (2014), no. 2, 431–449.
- [24] Z. Lu and Y. Zhang, *An augmented Lagrangian approach for sparse principal component analysis*, Math. Program. **135** (2012), 149–193.

- [25] J. H. Manton, *Optimization algorithms exploiting unitary constraints*, IEEE Trans. Signal Proces. **50** (2002), no. 3, 635–650.
- [26] N. Marzari and D. Vanderbilt, *Maximally localized generalized Wannier functions for composite energy bands*, Phys. Rev. B **56** (1997), no. 20, 12847.
- [27] M. Meinshausen and B. Yu, *Lasso-type recovery of sparse representations for high-dimensional data*, Ann. Stat. **1** (20081), no. 37, 246–270.
- [28] J. Nocedal and S. J. Wright, *Numerical Optimization*, 2nd ed., Springer Series in Operations Research and Financial Engineering, Springer New York, 2006.
- [29] S. Osher, Y. Mao, B. Dong, and W. Yin, *Fast linearized Bregman iteration for compressive sensing and sparse denoising*, Commun. Math. Sci. **8** (2010), no. 1, 93–111.
- [30] S. Osher and K. Yin, *On the completeness of the compressed modes in the eigenspace*, UCLA CAM Reports: 13-62 (2013).
- [31] V. Ozolinis, R. Lai, R. Caffisch, and S. Osher, *Compressed modes for variational problems in mathematics and physics*, P. Natl. Acad. Sci. USA **110** (2013), no. 46, 18368–18373.
- [32] ———, *Compressed plane waves - compactly supported multiresolution basis for the Laplace operator*, P. Natl. Acad. Sci. USA **111** (2014), no. 5, 1691–1696.
- [33] E. Prodan and W. Kohn, *Nearsightedness of electronic matter*, P. Natl. Acad. Sci. USA **102** (2005), no. 33, 11635–11638.
- [34] R. T. Rockafellar, R. J.-B. Wets, and M. Wets, *Variational Analysis*, Vol. 317, Springer, 1998.
- [35] L. I. Rudin, S. Osher, and E. Fatemi, *Nonlinear total variation based noise removal algorithms*, Physica D **60** (1992), no. 1, 259–268.
- [36] Z. Shen, *Wavelet Frames and Image Restorations*, Proc. ICM, 2010.
- [37] J. Tang and H. Liu, *Unsupervised feature selection for linked social media data*, Proc. 18th ACM SIGKDD International Conf. Knowl. Discov. Data Min., 2012, pp. 904–912.
- [38] R. Tibshirani, *Regression shrinkage and selection via the lasso*, J. Royal. Statist. Soc. B **1** (1996), no. 58, 267–288.
- [39] Z. Wen, X. Liu C. Yang, and Y. Zhang, *Trace-penalty minimization for large-scale eigenspace computation*, DTIC Document, 2013.
- [40] Z. Wen and W. Yin, *A feasible method for optimization with orthogonality constraints*, Math. Program. **142** (2013), no. 1-2, 397–434.
- [41] M. Yaghoobi and M. E. Davies, *Relaxed analysis operator learning*, Proceedings of the NIPS, Workshop on Analysis Operator Learning vs. Dictionary Learning: Fraternal Twins in Sparse Modeling, 2012.
- [42] M. Yaghoobi, S. Nam, R. Gribonval, and M. E. Davies, *Noise aware analysis operator learning for approximately cospase signals*, ICASSP - IEEE International Conference on Acoustics, Speech, and Signal Processing - 2012, Mar 2012, Kyoto, Japan. IEEE, 2012, 2012, pp. 5409–5412.
- [43] B. Yang, *Projection approximation subspace tracking*, IEEE T. Signal Proces. **43** (1995), no. 1, 95–107.
- [44] Y. Yang, H. T. Shen, Z. Ma, Z. Huang, and X. Zhou, *$\ell_{2,1}$ -norm regularized discriminative feature selection for unsupervised learning*, IJCAI, 2011, pp. 1589–1594.