

A Dimension Reduction Technique for Large-scale Structured Sparse Optimization Problems with Application to Convex Clustering

Yancheng Yuan[†] Tsung-Hui Chang[‡] Defeng Sun[§] Kim-Chuan Toh[¶]

August 18, 2021

Abstract

In this paper, we propose a novel adaptive sieving (AS) technique and an enhanced AS (EAS) technique, which are solver independent and could accelerate optimization algorithms for solving large scale convex optimization problems with intrinsic structured sparsity. We establish the finite convergence property of the AS technique and the EAS technique with inexact solutions of the reduced subproblems. As an important application, we apply the AS technique and the EAS technique on the convex clustering model, which could accelerate the state-of-the-art algorithm SSNAL by more than 7 times and the algorithm ADMM by more than 14 times.

Keywords: Adaptive sieving, structured sparsity, dimension reduction, convex optimization, convex clustering.

AMS subject classification: 90C06, 90C25, 90C90

1 Introduction and Related Work

Clustering is one of the most important and fundamental problems in data science, which plays important roles in numerous applications. Significant advances have been achieved in clustering during the last few decades, including K-means [11, 26], spectral clustering [15, 20], subspace clustering [21, 27] and so on. Despite these developments, some known drawbacks of these centroid based models, such as sensitivity to the initialization, limited effectiveness in high dimensional problems, and more importantly, the requirement on prior knowledge of the number of clusters, are

[†]Department of Applied Mathematics, The Hong Kong Polytechnic University, Hung Hom, Hong Kong (yancheng.yuan@polyu.edu.hk).

[‡]School of Science and Engineering, The Chinese University of Hong Kong (Shenzhen) and Shenzhen Research Institute of Big Data, China changtsunghui@cuhk.edu.cn. The research of this author is in part supported by the Shenzhen Research Institute of Big Data, under Grant 2019ORF01002.

[§]Department of Applied Mathematics, The Hong Kong Polytechnic University, Hung Hom, Hong Kong (defeng.sun@polyu.edu.hk). The research of this author is supported in part by Hong Kong Research Grant Council under grant number 15303720 and the Shenzhen Institute of Big Data, China under Grant 2019ORF01002.

[¶]Department of Mathematics and Institute of Operations Research and Analytics, National University of Singapore, 10 Lower Kent Ridge Road, Singapore (mattohkc@nus.edu.sg). The research of this author is supported by the Ministry of Education, Singapore, under its Academic Research Fund Tier 3 grant call (MOE-2019-T3-1-010).

still challenging to overcome. Here, we want to emphasize that the requirement on prior knowledge of the number of clusters is impractical for most real applications and to estimate the number of clusters itself is as important as clustering. One could argue that we can run classical clustering algorithms, such as K-means, with a few guesses on the number of clusters, but these clustering results are usually independent. Thus, users still need to determine the final clustering results subjectively based on their own preference.

Recently, the convex clustering approach has been proposed [4, 10, 17] and becomes more and more popular due to its good empirical performance and nice theoretical guarantees [3, 6, 16, 22, 23, 32]. More recently, some nonconvex variants of the convex clustering model have also been proposed [9, 18]. Specifically, for a given collection of n data points which are put as columns of a matrix $A \in \mathbb{R}^{d \times N}$, the convex clustering model is to solve the following optimization problem

$$\min_{X \in \mathbb{R}^{d \times N}} \frac{1}{2} \sum_{i=1}^N \|X_{:i} - A_{:i}\|^2 + \lambda \sum_{1 \leq i < j \leq N} w_{ij} \|X_{:i} - X_{:j}\|_p, \quad (1)$$

where $X_{:i}$ (or $A_{:i}$) is the i -th column of X (or A), $w_{ij} = w_{ji} \geq 0$ are given weights and $\lambda \geq 0$ is the hyper-parameter to control the strength of the diffusion penalty. Here $\|\cdot\|_p$ is the vector p -norm and we require $p \geq 1$ to guarantee the convexity of the model. After solving the model (1) and obtaining the solution X^* , we assign the data points $A_{:i}$ and $A_{:j}$ to the same cluster if $X_{:i}^* = X_{:j}^*$. Readers who are interested in more details about cluster identification based on the convex clustering model with an inexact solution could refer to [2, 5, 22]. It has been proved in [2] that the convex clustering model (1) can generate a continuous clustering path with respect to the hyper-parameter λ . Thus, prior knowledge on the number of clusters is **not** required, which is a highly desirable property.

Although the convex clustering model (1) is strongly convex, it is still quite challenging to solve since the number of terms in the diffusion penalty grows with n and could be extremely large (up to $O(n^2)$). A number of numerical optimization algorithms has been proposed for solving the convex clustering model. Among them, the alternating direction method of multipliers (ADMM) and the alternating minimization algorithm (AMA), which are proposed in [2], are very popular. Recently, a second-order based algorithm called SSNAL has been proposed [22, 29], which can efficiently solve (1) to achieve high accuracy for large n but moderate d , by taking advantage of the so-called second-order sparsity. However, the scalability could still be limited for those algorithms due to their need to handle all data points simultaneously. On the other hand, one may try to use some stochastic algorithms to solve (1) as in [16], however, the empirical performance is not so attractive since we need a rather accurate solution in order to determine the cluster memberships correctly based on the obtained solution X^* . Naturally one would ask whether we could design a deterministic algorithm which can scale as well as those stochastic algorithms. While this goal seems unattainable at the first glance, here we will give a positive answer. Now, we briefly explain the key idea on why this is possible. As demonstrated later in this paper, the same idea works not only for convex clustering model, but also for other optimization problems with special structures.

The key idea behind could actually be explained in a single sentence, that is, although the number of data points N could be extremely large, the number of clusters must be small, which is the purpose of clustering. In other words, for well-chosen values of λ , most columns of the optimal solution X^* for model (1) should be identical. In this case, most of the terms in the diffusion penalty should be zero. If we can remove those zero terms in advance and reduce the dimension

of X simultaneously, we only need to solve a small scale optimization problem, even for extremely large N . In this paper, we will propose an adaptive sieving (AS) technique and an enhanced adaptive sieving (EAS) technique, which are rigorous implementations of the aforementioned idea with theoretical guarantees. The details could be found later in this paper. We would like to emphasize that, the dimension reduction techniques proposed in this paper are solver independent, thus they could be applied to various algorithms which can solve the reduced optimization problems.

Motivated by the convex clustering problem, in this paper, we consider the optimization problems of the following form:

$$\min_{x \in \mathbb{R}^n} F_\lambda(x) := f(x) + \lambda p(Bx), \quad (2)$$

where $\lambda > 0$ is a hyper-parameter, $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a twice continuously differentiable convex function, $p : \mathbb{R}^m \rightarrow (-\infty, +\infty]$ is a closed proper convex function and $B : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a linear map. In some real applications, p is usually a regularizer which can enforce sparsity on Bx and B is a linear map which encodes desirable structures of x . This indicates the meaning of structured sparsity. With special designed matrices B , the optimization problem (2) includes many important models, such as the convex clustering model (1)¹, fused lasso model [25], clustered lasso model [19], and so on.

In this paper, we will propose solver independent techniques which can solve the optimization problem (2) via solving a sequence of subproblems with much smaller problem size. The main idea of this paper is inspired by a recent preprint [8]. Authors in [8] introduce an adaptive sieving technique to reduce the dimension of the optimization problem with sparse solutions (taking the linear map B to be identity I in (2)). However, the same idea cannot be directly applied to (2) for the cases where $B \neq I$. First, the optimal solutions of (2) may not be sparse at all, only with some special structures, such as being block-wise constant. Thus, the adaptive sieving technique in [8] could not apply on x . Second, one may try to apply the AS directly on Bx by introducing a new variable $y = Bx$. Although this idea may work for reducing the dimension of y (or Bx), this direct application cannot reduce the dimension of x simultaneously. If we cannot reduce the dimension of x simultaneously, we still need to solve large scale subproblems. Third, one of the keys for applying the AS technique is to check the optimality condition of (2) for a given $\bar{x} \in \mathbb{R}^n$. However, as one may see later, this is highly non-trivial if the inverse of B is not available (which is the case for most of problems with structured sparsity). In this paper, we will propose a new adaptive sieving technique and an enhanced adaptive sieving technique to address all of these issues.

To demonstrate the effectiveness of the proposed idea, we evaluate the empirical performance of our proposed AS and EAS technique with state-of-the-art algorithms: SSNAL [29], ADMM [2] and AMA [2], for solving the convex clustering model (1). As the readers will see later in the numerical experiments section, the numerical results on both simulated and real data sets demonstrate that the proposed AS and EAS could substantially reduce the dimension of the optimization problems. As a result, the AS/EAS techniques can accelerate the state-of-the-art algorithm SSNAL by more than 7 times and the algorithm ADMM by more than 14 times for solving the convex clustering model.

The main contributions of our paper can be summarized as follows:

- We propose a new solver independent adaptive sieving (AS) technique which can be applied to solve large scale optimization problems (2) with structured sparsity by solving a sequence of subproblems with much smaller size.

¹We can take $x = \text{vec}(X) \in \mathbb{R}^{dN}$, where $\text{vec}(X)$ is the vectorization of the matrix X by stacking its columns.

- We show the details of how to reduce the dimension of x and Bx simultaneously. Also, we show how to construct the corresponding reduced subproblem of (2) based on the structured sparsity of x (i.e., the sparsity of Bx).
- Our proposed AS technique allows the reduced subproblems to be solved inexactly and we prove the finite convergence property of the proposed AS technique for solving (2).
- As one will see later, although the AS technique will converge in finite iterations for solving (2), the sieving procedure of the AS technique may continue even if we obtain an optimal solution x^* of (2). To address this issue, we propose an enhanced adaptive sieving (EAS) technique, which can certify the optimality of an obtained solution with low additional computational cost. This can potentially reduce the sieving iterations of the AS technique and further accelerate the algorithms. The finite convergence property of the EAS technique is also proved.
- Both the AS technique and the EAS technique are extended to obtain a solution path of the structured sparse optimization problem (2) for a sequence of hyper-parameters $+\infty > \lambda_1 > \lambda_2 > \dots > \lambda_k > 0$.
- As an important application, extensive numerical experiments on the convex clustering model for both simulated and real data sets are provided. The superior numerical experiment results demonstrate the power of the AS technique and the EAS technique for accelerating numerical optimization algorithms to generate the solution path for the convex clustering model (1).

The rest of this paper is organized as follows: In Section 2, we introduce the adaptive sieving technique and the enhanced adaptive sieving technique for optimization problems with structured sparsity. The application of the AS technique and the EAS technique on the convex clustering model will be shown in Section 3 and numerical results are summarized in Section 4. We conclude the paper in Section 5.

Notation. We use blackboard bold capital letters to denote finite dimensional real Euclidean spaces, e.g. \mathbb{X}, \mathbb{Y} . In particular, we use $\mathbb{R}^{m \times n}$ (\mathbb{R}) to denote the set of all real $m \times n$ matrices (real numbers). We denote column vectors by lowercase letters, e.g. $v \in \mathbb{R}^n$, and matrices by capital letters, e.g. $A \in \mathbb{R}^{m \times n}$. We denote the transpose of the matrix A as A^T ; the i -th (i j -th) element of a vector v (matrix A) by v_i (A_{ij}). For a given integer $n \geq 1$, we denote the collection of integers from 1 to n by $[n]$. We denote the complement of an index set $I \subseteq [m]$ as I^c . For given index sets $I \subseteq [m]$ and $J \subseteq [n]$, we denote the submatrix consisting with rows (columns) indexed by I (J) as A_I ($A_{\cdot J}$). We denote the range space and null space of A by $\text{Range}(A)$ and $\text{Null}(A)$, respectively. For a vector $x \in \mathbb{R}^n$ and a scalar $p > 0$, we define the vector p -norm as: $\|x\|_p := (\sum_{i=1}^n |x_i|^p)^{1/p}$. We use $\|\cdot\|$ to denote the vector 2-norm. For a closed proper convex function $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$, the conjugate of f is $f^*(z) := \sup_{x \in \mathbb{R}^n} \{ \langle x, z \rangle - f(x) \}$. For a closed convex set $C \subseteq \mathbb{R}^n$ and a given vector $a \in \mathbb{R}^n$, the projection of a onto the set C is $\Pi_C(a) := \arg \min_{x \in C} \frac{1}{2} \|x - a\|^2$.

2 An Adaptive Sieving Technique for Structured Sparsity

In this section, we will introduce a novel adaptive sieving technique for obtaining the solution path for the structured sparse convex programming problem (2). Equivalently, we can reformulate (2)

as follows:

$$\begin{aligned} \min_{x \in \mathbb{R}^n, y \in \mathbb{R}^m} \quad & f(x) + \lambda p(y) \\ \text{s.t.} \quad & Bx - y = 0. \end{aligned} \tag{P_\lambda}$$

The Lagrangian function corresponding to (P_λ) is defined as:

$$l(x, y; z) := f(x) + \lambda p(y) + \langle z, Bx - y \rangle, \tag{3}$$

where $z \in \mathbb{R}^m$ is the Lagrange multiplier. The corresponding dual problem is given by

$$\max_{z \in \mathbb{R}^m} \quad D_\lambda(z) := -f^*(-B^T z) - \lambda p^*(z/\lambda). \tag{D_\lambda}$$

Here, f^* and p^* are the conjugate of f and p , respectively. Denote the solution set to (P_λ) as Ω_λ . The Karush-Kuhn-Tucker (KKT) conditions imply that $(x^*, y^*) \in \Omega_\lambda$ if and only if there exists $z^* \in \mathbb{R}^m$ such that

$$\begin{cases} \nabla f(x^*) + B^T z^* = 0, \\ z^* \in \lambda \partial p(y^*), \\ Bx^* - y^* = 0. \end{cases} \tag{KKT}$$

For any given triplet $(x, y, z) \in \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^m$, we define the KKT residual function for problem (P_λ) as:

$$R_\lambda(x, y, z) := \begin{pmatrix} \nabla f(x) + B^T z \\ y - \text{Prox}_{\lambda p}(y + z) \\ Bx - y \end{pmatrix}. \tag{4}$$

We know that $(x^*, y^*) \in \Omega_\lambda$ if and only if there exists $z^* \in \mathbb{R}^m$ such that

$$R_\lambda(x^*, y^*, z^*) = 0.$$

In this paper, we make the following two mild assumptions.

Assumption 1. *For any given $\lambda > 0$, the optimal solution set Ω_λ to the optimization problem (P_λ) is non-empty and compact.*

Assumption 2. *For any given $\lambda > 0$ and $y \in \mathbb{R}^m$, we define*

$$I^c := \{i \in [m] \mid y_i \neq 0\},$$

if $I^c \neq \emptyset$, then $(\partial(\lambda p(y)))_{I^c}$ is a singleton.

Remark 1. *Here, we make some remarks on Assumption 2. For most of the common used regularizers, such as lasso [24], group lasso [28], exclusive lasso [31], the Assumption 2 is satisfied. Let us take the lasso regularizer as an example. If $p(y) = \|y\|_1$, then,*

$$(\partial(\lambda p(y)))_i = \begin{cases} \lambda \text{sign}(y_i) & \text{if } y_i \neq 0, \\ [-\lambda, \lambda] & \text{if } y_i = 0. \end{cases}$$

Thus, we know that for any $y \neq 0$,

$$(\partial(\lambda p(y)))_{I^c} = (\lambda \text{sign}(y))_{I^c}$$

is a singleton. Here $\text{sign}(\cdot)$ is the signum function.

When the matrix B is the identity mapping and p is a regularizer that can induce sparsity, it has been demonstrated in [8] that we can substantially reduce the dimension of the problem (P_λ) by applying the adaptive sieving technique. However, it is not clear whether a similar idea could benefit those models whose solutions are not sparse but have some special structures. In this paper, we will give a positive answer to this question in the following sections.

The theme of this paper is to design a technique which can reduce the dimension of a class of optimization problem (P_λ) with structured sparsity by exploring the intrinsic structure of the problem in an explicit way. Readers will see shortly that the key idea behind is quite simple but a rigorous realization of this simple idea is highly non-trivial.

We first introduce our principal idea in a general way, then, we will propose a technique called adaptive sieving (AS) to rigorously implement the idea. We fix the parameter λ in (P_λ) for now. For a given index set $I \subseteq [m]$, if there is some prior knowledge for us to believe that $y_I = 0$, which we call the structured sparsity, then it is natural for us to consider the following constrained optimization problem generated by the index set I :

$$\begin{aligned} \min_{x \in \mathbb{R}^n, y \in \mathbb{R}^m} \quad & f(x) + \lambda p(y) \\ & Bx - y = 0, \\ & y_I = 0. \end{aligned} \tag{P_\lambda(I)}$$

We denote this problem as $(P_\lambda(I))$ to indicate its dependence on the index set I , we will denote the index set $[m] \setminus I$, which is the complement of I , as I^c .

Our principal idea is to obtain a solution to the original optimization problem (P_λ) by solving a sequence of subproblems with lower dimension, which are induced by $(P_\lambda(I))$. The key for a successful realization of this principal idea for solving (P_λ) depends on answering the following questions:

- Q1: For a given index set $I \in [m]$, how to effectively reduce the dimension of (P_λ) based on $(P_\lambda(I))$?
- Q2: If we can solve $(P_\lambda(I))$ to obtain a solution pair (\bar{x}, \bar{y}) , which is **not** yet a solution to (P_λ) , can we guarantee that we can update the index set I to construct a new reduced problem in the form of $(P_\lambda(I))$?
- Q3: If the solution pair (\bar{x}, \bar{y}) , which is obtained by solving $(P_\lambda(I))$, is indeed a solution to (P_λ) , can we certify its optimality and stop the whole procedure?
- Q4: Is the proposed technique robust to the inexactness of the obtained solution pair? In other words, if we can only obtain an inexact solution of $(RP_\lambda(I))$ (defined in Section 2.1) under a given tolerance $\epsilon > 0$, can we obtain a solution of (P_λ) under the tolerance $O(\epsilon)$?
- Q5: Is it possible for the proposed technique to be solver independent? In other words, the technique could be applied to any algorithms that can solve $(RP_\lambda(I))$ inexactly under a given tolerance.

Remark 2. *We make some remarks before we describe the proposed AS technique.*

1. *Although designing an efficient and convergent algorithm for solving $(P_\lambda(I))$ is also an important task, it is not the main purpose of this paper. Actually, as one may see later, our proposed AS technique is solver independent. There are also existing algorithms which can solve $(RP_\lambda(I))$ to a moderate accuracy [2, 22, 29].*

2. Although it seems unnecessary to raise the question Q3 at the first glance, it is actually essential for applying any dimension reduction technique to solve (P_λ) based on $(P_\lambda(I))$. In order to check the optimality of the solution pair (\bar{x}, \bar{y}) , we need to construct the corresponding dual solution and check the corresponding KKT condition (KKT). This is highly non-trivial since the dual solutions are not unique if structured sparsity exists. This is also one of the main difficulties for applying the AS technique to problem (P_λ) with structured sparsity as compared to [8].
3. The robustness mentioned in Q4 is also very important since the best we can expect in general is to obtain an inexact solution to $(RP_\lambda(I))$.

Now, we start to give the details of our realization of the principal idea.

2.1 A Dimension Reduction Technique for (P_λ) Based on $(P_\lambda(I))$

We first show how we can reduce the dimension of the variables x and y simultaneously for the problem (P_λ) based on the constrained optimization problem $(P_\lambda(I))$, which is one of the core ideas of this paper. These details also answer the question Q1.

Assume that the rank of B_I is $r > 0$. Then there exists three index sets α , β and γ with $|\gamma| = r$, that forms a partition of $[n]$, such that $B_{I\beta} = 0$ and $B_{I\gamma}$ has full column rank. Here, we also assume that the index set α is nonempty; otherwise, we must have

$$x_\gamma = 0.$$

Since $B_{I\gamma}$ has full column rank, there is a unique $|\gamma| \times |\alpha|$ matrix $M_{\gamma\alpha}$,² such that

$$B_{I\alpha} + B_{I\gamma}M_{\gamma\alpha} = 0. \quad (5)$$

Then, we can eliminate x_γ by the constraints of $(P_\lambda(I))$ as

$$x_\gamma = M_{\gamma\alpha}x_\alpha.$$

Define:

$$\varphi(x_\alpha, x_\beta) = f(\dot{x}), \quad q(y_{I^c}) = p(\dot{y}),$$

where

$$\dot{x}_\alpha = x_\alpha, \quad \dot{x}_\beta = x_\beta, \quad \dot{x}_\gamma = M_{\gamma\alpha}x_\alpha,$$

and

$$\dot{y}_i = \begin{cases} y_i, & \text{if } i \in I^c, \\ 0 & \text{if } i \in I. \end{cases}$$

It is not difficult to realize that we can solve problem $(P_\lambda(I))$ via solving the following reduced optimization problem:

$$\begin{aligned} \min_{x_\alpha \in \mathbb{R}^{|\alpha|}, x_\beta \in \mathbb{R}^{|\beta|}, y_{I^c} \in \mathbb{R}^{|I^c|}} & \varphi(x_\alpha, x_\beta) + \lambda q(y_{I^c}) \\ \text{s.t.} & (B_{I^c\alpha} + B_{I^c\gamma}M_{\gamma\alpha})x_\alpha + B_{I^c\beta}x_\beta - y_{I^c} = 0. \end{aligned} \quad (RP_\lambda(I))$$

²Here, we abuse the notation a little bit to indicate the dependence of $M_{\gamma\alpha}$ on the index sets α and γ . The uniqueness is in the sense of a given partition.

The Lagrange function corresponding to $(RP_\lambda(I))$ is given by

$$l(x_\alpha, x_\beta, y_{I^c}, \xi) = \varphi(x_\alpha, x_\beta) + \lambda q(y_{I^c}) + \langle \xi, (B_{I^c\alpha} + B_{I^c\gamma}M_{\gamma\alpha})x_\alpha + B_{I^c\beta}x_\beta - y_{I^c} \rangle,$$

where $\xi \in \mathbb{R}^{|I^c|}$ is the Lagrange multiplier.

Now, if we solve $(RP_\lambda(I))$ and obtain a solution $(\hat{x}_\alpha, \hat{x}_\beta, \hat{y}_{I^c})$, then, there exists a $\hat{\xi}$ that satisfies the following KKT condition:

$$\begin{cases} (\nabla f(\hat{x}))_\alpha + M_{\gamma\alpha}^T(\nabla f(\hat{x}))_\gamma + (B_{I^c\alpha} + B_{I^c\gamma}M_{\gamma\alpha})^T\hat{\xi} = 0, \\ (\nabla f(\hat{x}))_\beta + B_{I^c\beta}^T\hat{\xi} = 0, \quad \hat{\xi} \in (\partial(\lambda p(\hat{y})))_{I^c}, \\ (B_{I^c\alpha} + B_{I^c\gamma}M_{\gamma\alpha})\hat{x}_\alpha + B_{I^c\beta}\hat{x}_\beta - \hat{y}_{I^c} = 0, \end{cases} \quad (6)$$

where \hat{x} and \hat{y} are defined as

$$\hat{x}_\alpha = \hat{x}_\alpha, \quad \hat{x}_\beta = \hat{x}_\beta, \quad \hat{x}_\gamma = M_{\gamma\alpha}\hat{x}_\alpha$$

and

$$\hat{y}_{I^c} = \hat{y}_{I^c}, \quad \hat{y}_I = 0,$$

respectively. Then (\bar{x}, \bar{y}) , which is constructed by

$$\begin{cases} \bar{x}_\alpha = \hat{x}_\alpha, \quad \bar{x}_\beta = \hat{x}_\beta, \quad \bar{x}_\gamma = M_{\gamma\alpha}\hat{x}_\alpha, \\ \bar{y}_{I^c} = \hat{y}_{I^c}, \quad \bar{y}_I = 0 \end{cases} \quad (7)$$

is a solution to problem $(P_\lambda(I))$. Thus, in order to obtain a solution to $(P_\lambda(I))$, we only need to solve a corresponding reduced problem $(RP_\lambda(I))$ whose dimension can be much smaller.

Remark 3. *We make some remarks to close this subsection.*

1. *We reduce the dimension of the problem from $\mathbb{R}^n \times \mathbb{R}^m$ to $\mathbb{R}^{n-|\gamma|} \times \mathbb{R}^{m-|I|}$, which can be a substantial reduction. For example, if the solution of (P_λ) is indeed sparse (this is an intrinsic property since we can obtain a sparse solution in general for large λ), then $|I|$ is close to m and $|\gamma|$ is close to n simultaneously.*
2. *In many real applications (e.g. convex clustering), we can identify the index set α , β , γ and construct the matrix $M_{\gamma\alpha}$ at a low cost. Also, since the linear map B is designed to encode some structures of the solution, it is usually very sparse.*

2.2 An Adaptive Sieving Technique for (P_λ) with a Fixed $\lambda > 0$

Now, we move on to present the details of the AS technique. We fix the parameter λ for now and we will generalize it to handle the case for a sequence of $\lambda > 0$ later. Also, for simplicity, we first present the idea with the assumption that we can solve $(RP_\lambda(I))$ exactly. The same idea will be generalized to the inexact setting without much difficulties later.

We first show how we can update the index set I if the current obtained solution (\bar{x}, \bar{y}) via solving $(P_\lambda(I))$ is not an optimal solution to (P_λ) . The key idea is to construct a corresponding dual variable pair $(\bar{u}, \bar{w}) \in \mathbb{R}^m \times \mathbb{R}^{|I|}$ which satisfies the following KKT condition for $(P_\lambda(I))$:

$$\begin{cases} (\nabla f(\bar{x}))_\alpha + B_{I\alpha}^T\bar{u}_I + B_{I^c\alpha}^T\bar{u}_{I^c} = 0, \\ (\nabla f(\bar{x}))_\beta + B_{I\beta}^T\bar{u}_I + B_{I^c\beta}^T\bar{u}_{I^c} = 0, \quad (\nabla f(\bar{x}))_\gamma + B_{I\gamma}^T\bar{u}_I + B_{I^c\gamma}^T\bar{u}_{I^c} = 0, \\ \bar{u}_I - \bar{w} \in \lambda(\partial p(\bar{y}))_I, \quad \bar{u}_{I^c} \in \lambda(\partial p(\bar{y}))_{I^c}, \\ B\bar{x} - \bar{y} = 0, \quad \bar{y}_I = 0. \end{cases} \quad (8)$$

Since $(\bar{x}, \bar{y}) = (\hat{x}, \hat{y})$ and $(\hat{x}, \hat{y}, \hat{\xi})$ is a solution to (6), we must have

$$B_{I^c\beta}^T \hat{\xi} = B_{I^c\beta}^T \bar{u}_{I^c}, \quad \hat{\xi} \in (\partial(\lambda p(\hat{y})))_{I^c}, \quad \bar{u}_{I^c} \in (\partial(\lambda p(\bar{y})))_{I^c}.$$

Aggressively, we construct \bar{u}_{I^c} as

$$\bar{u}_{I^c} = \hat{\xi}. \quad (9)$$

By the above construction of \bar{u}_{I^c} and the equation (5), the first equation of (8) is implied by the third equation of (8) and the first equation of (6). Thus, we can construct the pair (\bar{u}_I, \bar{w}) via solving the following equations for (u_I, w) :

$$\begin{cases} (\nabla f(\bar{x}))_\gamma + B_{I\gamma}^T u_I + B_{I^c\gamma}^T \bar{u}_{I^c} = 0, \\ u_I - w \in (\partial(\lambda p(\bar{y})))_I. \end{cases} \quad (10)$$

Since w is an unconstrained variable, for any \hat{u}_I satisfying the first equation of (10), there exists a \hat{w} such that the second one is satisfied. However, realizing the fact that if there exists a \tilde{u}_I such that $(\tilde{u}_I, 0)$ is a solution to (10), then the current solution pair (\bar{x}, \bar{y}) is an optimal solution to (P_λ) . Thus, we propose to construct the pair (\bar{u}_I, \bar{w}) such that \bar{w} has the minimum Euclidean norm. Since $B_{I\gamma}$ has full column rank, we can construct a particular solution to the first equation of (10) as

$$(\bar{u}_I)_0 = -B_{I\gamma}(B_{I\gamma}^T B_{I\gamma})^{-1}((\nabla f(\bar{x}))_\gamma + B_{I^c\gamma}^T \bar{u}_{I^c}). \quad (11)$$

Thus, all the solution to the first equation of (10) is given by

$$u_I = (\bar{u}_I)_0 + d,$$

where $d \in \text{Null}(B_{I\gamma}^T)$. In summary, we construct the solution pair (\bar{u}_I, \bar{w}) as follows:

$$\bar{u}_I = (\bar{u}_I)_0 + \bar{d}, \quad \bar{w} = \bar{u}_I - \Pi_{(\partial(\lambda p(\bar{y})))_I}(\bar{u}_I), \quad (12)$$

where \bar{d} is a solution to the following auxiliary optimization problem:

$$\begin{aligned} \min_{d \in \mathbb{R}^{|I|}} & \quad \frac{1}{2} \|((\bar{u}_I)_0 + d) - \Pi_{(\partial(\lambda p(\bar{y})))_I}((\bar{u}_I)_0 + d)\|^2 \\ \text{s.t.} & \quad d \in \text{Null}(B_{I\gamma}^T). \end{aligned} \quad (13)$$

Up to this point, we have completed the construction of a dual solution pair (\bar{u}, \bar{w}) . We show the nice properties of the constructed (\bar{u}, \bar{w}) in Theorem 1 and Theorem 2.

Theorem 1. *Assume that $(\hat{x}_\alpha, \hat{x}_\beta, \hat{y}_{I^c})$ is an optimal solution to the following optimization problem:*

$$\begin{aligned} \min_{x_\alpha \in \mathbb{R}^{|\alpha|}, x_\beta \in \mathbb{R}^{|\beta|}, y_{I^c} \in \mathbb{R}^{|I^c|}} & \quad \varphi(x_\alpha, x_\beta) + \lambda q(y_{I^c}) + \langle x_\alpha, \hat{\delta}_1 \rangle + \langle x_\beta, \hat{\delta}_2 \rangle - \langle y_{I^c}, \hat{\delta}_3 \rangle \\ \text{s.t.} & \quad (B_{I^c\alpha} + B_{I^c\gamma} M_{\gamma\alpha})x_\alpha + B_{I^c\beta} x_\beta - y_{I^c} = 0, \end{aligned} \quad (14)$$

and $\hat{\xi}$ is the corresponding Lagrange multiplier. Here, $(\hat{\delta}_1, \hat{\delta}_2, \hat{\delta}_3) \in \mathbb{R}^{|\alpha|} \times \mathbb{R}^{|\beta|} \times \mathbb{R}^{|I^c|}$ are given error terms satisfying $\|\hat{\delta}_1\| + \|\hat{\delta}_2\| + \|\hat{\delta}_3\| \leq \epsilon$. Let $(\bar{x}, \bar{y}, \bar{u}_{I^c}, \bar{u}_I, \bar{w})$ be the solution that is constructed from (7), (9), and (12). Define $J(\lambda)$ as follows:

$$J(\lambda) := \{j \in I \mid \bar{u}_j \notin (\partial(\lambda p(\bar{y})))_j\}. \quad (15)$$

Then, $J(\lambda) \neq \emptyset$ if

$$\|R_\lambda(\bar{x}, \bar{y}, \bar{u})\| > \epsilon.$$

Proof. Since $(\hat{x}_\alpha, \hat{x}_\beta, \hat{y}_{I^c})$ is an optimal solution to (14) and $\hat{\xi}$ is the corresponding Lagrange multiplier, the following KKT system holds:

$$\begin{cases} (\nabla f(\hat{x}))_\alpha + M_{\gamma\alpha}^T(\nabla f(\hat{x}))_\gamma + (B_{I^c\alpha} + B_{I^c\gamma}M_{\gamma\alpha})^T\hat{\xi} + \hat{\delta}_1 = 0, \\ (\nabla f(\hat{x}))_\beta + B_{I^c\beta}^T\hat{\xi} + \hat{\delta}_2 = 0, \\ \hat{\xi} + \hat{\delta}_3 \in (\partial(\lambda p(\hat{y})))_{I^c}, \\ (B_{I^c\alpha} + B_{I^c\gamma}M_{\gamma\alpha})\hat{x}_\alpha + B_{I^c\beta}\hat{x}_\beta - \hat{y}_{I^c} = 0. \end{cases} \quad (16)$$

By construction, $(\bar{x}, \bar{y}, \bar{u}_{I^c}, \bar{u}_I, \bar{w})$ is a solution to:

$$\begin{cases} (\nabla f(\bar{x}))_\alpha + B_{I\alpha}^T\bar{u}_I + B_{I^c\alpha}^T\bar{u}_{I^c} + \hat{\delta}_1 = 0, \\ (\nabla f(\bar{x}))_\beta + B_{I\beta}^T\bar{u}_I + B_{I^c\beta}^T\bar{u}_{I^c} + \hat{\delta}_2 = 0, \\ (\nabla f(\bar{x}))_\gamma + B_{I\gamma}^T\bar{u}_I + B_{I^c\gamma}^T\bar{u}_{I^c} = 0, \\ \bar{u}_I - \bar{w} \in \lambda(\partial p(\bar{y}))_I, \\ \bar{u}_{I^c} + \hat{\delta}_3 \in \lambda(\partial p(\bar{y}))_{I^c}, \\ B\bar{x} - \bar{y} = 0, \quad \bar{y}_I = 0. \end{cases} \quad (17)$$

Now, we prove that $J(\lambda) \neq \emptyset$ provided $\|R_\lambda(\bar{x}, \bar{y}, \bar{u})\| > \epsilon$. We prove it by contradiction. Assume that

$$J(\lambda) = \emptyset.$$

Then we have

$$\bar{u} + \hat{\delta} \in \partial(\lambda p(\bar{y})),$$

where $\hat{\delta} = (\hat{\delta}_I, \hat{\delta}_{I^c}) = (0, \hat{\delta}_3)$. This implies that

$$\bar{y} - \text{Prox}_{\lambda p}(\bar{y} + (\bar{u} + \hat{\delta})) = 0.$$

Then,

$$\begin{aligned} \|R_\lambda(\bar{x}, \bar{y}, \bar{u})\| &= \|(\nabla f(\bar{x}) + B^T\bar{u}, \bar{y} - \text{Prox}_{\lambda p}(\bar{y} + \bar{u}), B\bar{x} - \bar{y})\| \\ &= \|((-\hat{\delta}_1, -\hat{\delta}_2, 0), \text{Prox}_{\lambda p}(\bar{y} + (\bar{u} + \hat{\delta})) - \text{Prox}_{\lambda p}(\bar{y} + \bar{u}), 0)\| \\ &\leq \|\hat{\delta}_1\| + \|\hat{\delta}_2\| + \|\hat{\delta}_3\| \\ &\leq \epsilon. \end{aligned} \quad (18)$$

Here, we used the property that the proximal mapping is Lipschitz continuous with modulus 1. This is a contradiction. Thus $J(\lambda) \neq \emptyset$ and we proved the statement in the theorem. \square

Remark 4. We do not need to specify a priori error terms $\hat{\delta}_1, \hat{\delta}_2, \hat{\delta}_3$ in Theorem 1. They should be interpreted as the errors incurred when we solve the problem $(RP_\lambda(I))$ inexactly with a given tolerance.

An important implication of Theorem 1 is that, if the current obtained solution pair (\bar{x}, \bar{y}) is not an inexact optimal solution to (P_λ) under the given tolerance, we can update the index set I by removing the identified violated index set $J(\lambda)$. This important implication motivates us to propose the adaptive sieving (AS) technique for (P_λ) with a given fixed $\lambda > 0$, which is presented in Algorithm 1.

Algorithm 1 Adaptive sieving for solving (P_λ) with a fixed $\lambda > 0$

1: **Input:** a given hyper-parameter $\lambda > 0$ and a given tolerance $\epsilon > 0$.
2: **Output:** $(x^*(\lambda), y^*(\lambda), z^*(\lambda))$.
3: **Initialization:** Generate an initial index set by a predefined initialization strategy: $I^0(\lambda) \subseteq [m]$.
4: **for** $i = 0, 1, 2, \dots$ **do**
5: **1.** For the given index set $I^i(\lambda)$, construct the index partition $\{\alpha^i, \beta^i, \gamma^i\}$ and the corresponding $M_{\gamma^i \alpha^i}$.
6: **2.** Apply any well designed algorithm to solve problem $(RP_\lambda(I))$ with $\{I^i(\lambda), \alpha^i, \beta^i, \gamma^i, M_{\gamma^i \alpha^i}\}$ and obtain an inexact solution $(\hat{x}_{\alpha^i}^i, \hat{x}_{\beta^i}^i, \hat{y}_{(I^i)^c(\lambda)}^i, \hat{\xi}^i)$ which satisfies the corresponding KKT system (16) with the latent error terms $(\hat{\delta}_1^i, \hat{\delta}_2^i, \hat{\delta}_3^i)$ such that $\|\hat{\delta}_1^i\| + \|\hat{\delta}_2^i\| + \|\hat{\delta}_3^i\| \leq \epsilon$.
7: **3.** Recover a solution $(\bar{x}^i, \bar{y}^i, \bar{u}^i, \bar{w}^i)$ by the construction of (7), (9) and (12), respectively.
8: **if** $\|R_\lambda(\bar{x}^i, \bar{y}^i, \bar{u}^i)\| \leq \epsilon$ **then**
9: Set $(x^*(\lambda), y^*(\lambda), z^*(\lambda)) = (\bar{x}^i, \bar{y}^i, \bar{u}^i)$.
10: **break.**
11: **else**
12: Create $J^i(\lambda)$:

$$J^i(\lambda) = \{j \in I^i(\lambda) \mid \bar{u}_j^i \notin \partial(\lambda p(\bar{y}^i))_j\}, \quad (19)$$

13: **if** $J^i(\lambda) \neq \emptyset$ **then**
14: Update $I^{i+1}(\lambda)$ as:

$$I^{i+1}(\lambda) \leftarrow I^i(\lambda) \setminus J^i(\lambda).$$

15: **else**
16: Set $(x^*(\lambda), y^*(\lambda), z^*(\lambda)) = (\bar{x}^i, \bar{y}^i, \bar{u}^i)$.
17: **break.**
18: **end if**
19: **end if**
20: **end for**
21: **return** $(x^*(\lambda), y^*(\lambda), z^*(\lambda))$.

Theorem 2. For a given $\epsilon > 0$, with any well designed algorithm which can solve the reduced subproblem $(RP_\lambda(I))$ to the given accuracy, Algorithm 1 is guaranteed to converge in finite number of iterations. Moreover, the obtained pair $(x^*(\lambda), y^*(\lambda), z^*(\lambda))$ is a solution to (P_λ) in the sense that

$$\|R_\lambda(x^*(\lambda), y^*(\lambda), z^*(\lambda))\| \leq \epsilon.$$

We omit the proof of Theorem 2 here as it is a byproduct of Theorem 1

Remark 5. We close this subsection by making some remarks here.

1. The proposed AS technique is a practical implementation of the aforementioned principal idea, which is solver independent and answers Q1, Q2, Q4 and Q5 simultaneously.
2. However, it may fail to answer the question Q3. The whole procedure described in Algorithm 1 is not guaranteed to certify the optimality of a given solution pair (\bar{x}, \bar{y}) , even if it is already optimal for (P_λ) . The constructed \bar{u} may not be the correct corresponding Lagrange multiplier. The main reason is because we have aggressively set $\bar{u}_{I^c} = \hat{\xi}$ in (9).

3. Although Algorithm 1 may fail to answer the question Q3 and it may need additional iterations to terminate the whole algorithm, the practical performance of Algorithm 1 is actually quite promising. Readers can find the numerical performance in the numerical experiments section.
4. In order to address the possible weakness of the construction of \bar{u} mentioned in item 2, we will propose an enhanced AS technique which can answer all the five questions simultaneously in the next subsection.

2.3 An Enhanced Adaptive Sieving Technique

Now, we introduce an enhanced adaptive sieving technique which can certify the optimality of the obtained pair (\bar{x}, \bar{y}) via solving the reduced subproblem $(RP_\lambda(I))$ if it is optimal to (P_λ) . With the enhanced AS technique, we can potentially reduce the number of sieving iterations of Algorithm 1.

The key idea is to deal with the issue we mentioned in Remark 5. Now, assume that (\bar{x}, \bar{y}) is an optimal solution to $(P_\lambda(I))$, which could be recovered by (7) with a solution of $(RP_\lambda(I))$. We can then define a new index set \tilde{I} as follows:

$$\tilde{I} := \{i \in [m] \mid \bar{y}_i = 0\}. \quad (20)$$

By the construction, we have $I \subseteq \tilde{I}$. It is not difficult to see that (\bar{x}, \bar{y}) is actually an optimal solution to the following constrained optimization problem:

$$\begin{aligned} \min_{x \in \mathbb{R}^n, y \in \mathbb{R}^m} \quad & f(x) + \lambda p(y) \\ \text{s.t.} \quad & Bx - y = 0, \\ & y_{\tilde{I}} = 0. \end{aligned} \quad (P_\lambda(\tilde{I}))$$

In a similar manner, we can define the index sets $\tilde{\alpha}$, $\tilde{\beta}$ and $\tilde{\gamma}$ with $|\tilde{\gamma}| = \tilde{r}$, which form a partition of $[n]$, such that $B_{\tilde{I}\tilde{\beta}} = 0$ and $B_{\tilde{I}\tilde{\gamma}}$ has full column rank. Again, we assume that $\tilde{\alpha} \neq \emptyset$. Thus, there exists a $M_{\tilde{\gamma}\tilde{\alpha}} \in \mathbb{R}^{|\tilde{\gamma}| \times |\tilde{\alpha}|}$ such that

$$B_{\tilde{I}\tilde{\alpha}} + B_{\tilde{I}\tilde{\gamma}}M_{\tilde{\gamma}\tilde{\alpha}} = 0.$$

Then, we can eliminate $x_{\tilde{\gamma}}$ by the constraints of $(P_\lambda(\tilde{I}))$ as

$$x_{\tilde{\gamma}} = M_{\tilde{\gamma}\tilde{\alpha}}x_{\tilde{\alpha}}.$$

The Lagrangian function corresponding to $(P_\lambda(\tilde{I}))$ is given by

$$l(x, y, v, s) = f(x) + \lambda p(y) + \langle v, Bx - y \rangle + \langle s, y_{\tilde{I}} \rangle,$$

where $v \in \mathbb{R}^m$ and $s \in \mathbb{R}^{|\tilde{I}|}$ are the Lagrange multipliers. For notational consistency, we denote $(\tilde{x}, \tilde{y}) = (\bar{x}, \bar{y})$. Since (\tilde{x}, \tilde{y}) is an optimal solution to $(P_\lambda(\tilde{I}))$, there exists (\tilde{v}, \tilde{s}) such that the following KKT condition for $(P_\lambda(\tilde{I}))$ is satisfied:

$$\begin{cases} (\nabla f(\tilde{x}))_{\tilde{\alpha}} + B_{\tilde{I}\tilde{\alpha}}^T \tilde{v}_{\tilde{I}} + B_{\tilde{I}c\tilde{\alpha}}^T \tilde{v}_{\tilde{I}c} = 0, & (\nabla f(\tilde{x}))_{\tilde{\beta}} + B_{\tilde{I}\tilde{\beta}}^T \tilde{v}_{\tilde{I}c} = 0, \\ (\nabla f(\tilde{x}))_{\tilde{\gamma}} + B_{\tilde{I}\tilde{\gamma}}^T \tilde{v}_{\tilde{I}} + B_{\tilde{I}c\tilde{\gamma}}^T \tilde{v}_{\tilde{I}c} = 0, \\ \tilde{v}_{\tilde{I}} - \tilde{s} \in \lambda(\partial p(\tilde{y}))_{\tilde{I}}, & \tilde{v}_{\tilde{I}c} \in \lambda(\partial p(\tilde{y}))_{\tilde{I}c}, \\ B\tilde{x} - \tilde{y} = 0, & \tilde{y}_{\tilde{I}} = 0. \end{cases} \quad (21)$$

On the other hand, we know that $(\tilde{x}_{\tilde{\alpha}}, \tilde{x}_{\tilde{\beta}}, \tilde{y}_{\tilde{I}^c})$ is an optimal solution to the following reduced problem corresponding to $(P_\lambda(\tilde{I}))$:

$$\begin{aligned} \min_{x_{\tilde{\alpha}} \in \mathbb{R}^{|\tilde{\alpha}|}, x_{\tilde{\beta}} \in \mathbb{R}^{|\tilde{\beta}|}, y_{\tilde{I}^c} \in \mathbb{R}^{|\tilde{I}^c|}} \quad & \tilde{\varphi}(x_{\tilde{\alpha}}, x_{\tilde{\beta}}) + \lambda \tilde{q}(y_{\tilde{I}^c}) \\ \text{s.t.} \quad & (B_{\tilde{I}^c \tilde{\alpha}} + B_{\tilde{I}^c \tilde{\gamma}} M_{\tilde{\gamma} \tilde{\alpha}}) x_{\tilde{\alpha}} + B_{\tilde{I}^c \tilde{\beta}} x_{\tilde{\beta}} - y_{\tilde{I}^c} = 0, \end{aligned} \quad (RP_\lambda(\tilde{I}))$$

where

$$\tilde{\varphi}(x_{\tilde{\alpha}}, x_{\tilde{\beta}}) = f(\tilde{x}), \quad \tilde{q}(y_{\tilde{I}^c}) = p(\tilde{y}).$$

Here

$$\dot{x}_{\tilde{\alpha}} = x_{\tilde{\alpha}}, \quad \dot{x}_{\tilde{\beta}} = x_{\tilde{\beta}}, \quad \dot{x}_{\tilde{\gamma}} = M_{\tilde{\gamma} \tilde{\alpha}} x_{\tilde{\alpha}},$$

and

$$\dot{y}_i = \begin{cases} y_i, & \text{if } i \in \tilde{I}^c, \\ 0 & \text{if } i \in \tilde{I}. \end{cases}$$

Since $(\tilde{x}_{\tilde{\alpha}}, \tilde{x}_{\tilde{\beta}}, \tilde{y}_{\tilde{I}^c})$ is an optimal solution to $(RP_\lambda(\tilde{I}))$, there exists a $\tilde{\theta} \in \mathbb{R}^{|\tilde{I}^c|}$ such that the following KKT condition is satisfied:

$$\begin{cases} (\nabla f(\tilde{x}))_{\tilde{\alpha}} + M_{\tilde{\gamma} \tilde{\alpha}}^T (\nabla f(\tilde{x}))_{\tilde{\gamma}} + (B_{\tilde{I}^c \tilde{\alpha}} + B_{\tilde{I}^c \tilde{\gamma}} M_{\tilde{\gamma} \tilde{\alpha}})^T \tilde{\theta} = 0, \\ (\nabla f(\tilde{x}))_{\tilde{\beta}} + B_{\tilde{I}^c \tilde{\beta}}^T \tilde{\theta} = 0, \quad \tilde{\theta} \in \lambda (\partial p(\tilde{y}))_{\tilde{I}^c}, \\ (B_{\tilde{I}^c \tilde{\alpha}} + B_{\tilde{I}^c \tilde{\gamma}} M_{\tilde{\gamma} \tilde{\alpha}}) \tilde{x}_{\tilde{\alpha}} + B_{\tilde{I}^c \tilde{\beta}} \tilde{x}_{\tilde{\beta}} - \tilde{y}_{\tilde{I}^c} = 0. \end{cases} \quad (22)$$

Again, the key is to construct a dual pair (\tilde{v}, \tilde{s}) from the KKT system (22) such that $(\tilde{x}, \tilde{y}, \tilde{v}, \tilde{s})$ is a solution to (21). Fortunately, by Assumption 2 and the fact $\tilde{I}^c = \{i \in [m] \mid \tilde{y}_i \neq 0\}$, we have

$$\tilde{v}_{\tilde{I}^c} = (\partial(\lambda p(\tilde{y})))_{\tilde{I}^c} = \tilde{\theta}. \quad (23)$$

Thus, by the uniqueness of $\tilde{v}_{\tilde{I}^c}$, the second equation of (21) must be satisfied.

Similarly, we construct $(\tilde{v}_{\tilde{I}}, \tilde{s})$ as follows:

$$\tilde{v}_{\tilde{I}} = (\tilde{v}_{\tilde{I}})_0 + \tilde{d}, \quad \tilde{s} = \tilde{v}_{\tilde{I}} - \Pi_{(\partial(\lambda p(\tilde{y})))_{\tilde{I}}}(\tilde{v}_{\tilde{I}}), \quad (24)$$

where

$$(\tilde{v}_{\tilde{I}})_0 = -B_{\tilde{I} \tilde{\gamma}} (B_{\tilde{I} \tilde{\gamma}}^T B_{\tilde{I} \tilde{\gamma}})^{-1} ((\nabla f(\tilde{x}))_{\tilde{\gamma}} + B_{\tilde{I}^c \tilde{\gamma}}^T \tilde{v}_{\tilde{I}^c})$$

and \tilde{d} is an optimal solution to the following auxiliary optimization problem:

$$\begin{aligned} \min_{d \in \mathbb{R}^{|\tilde{I}|}} \quad & \frac{1}{2} \|((\tilde{v}_{\tilde{I}})_0 + d) - \Pi_{(\partial(\lambda p(\tilde{y})))_{\tilde{I}}}((\tilde{v}_{\tilde{I}})_0 + d)\|^2 \\ \text{s.t.} \quad & d \in \text{Null}(B_{\tilde{I} \tilde{\gamma}}^T). \end{aligned} \quad (25)$$

For the above constructed (\tilde{v}, \tilde{s}) , it has a nice property to be summarized in the following theorem. It shows that the constructed dual variable \tilde{v} can certify the optimality of \tilde{x} .

Theorem 3. *For a given $\epsilon > 0$, if the current obtained solution \tilde{x} by solving $(RP_\lambda(I))$ is an optimal solution to the following perturbed optimization problem*

$$\min_{x \in \mathbb{R}^n} \quad f(x) + \lambda p(Bx) + \langle x, \tilde{\delta} \rangle, \quad (26)$$

where $\tilde{\delta} \in \mathbb{R}^n$ is a latent error vector such that $\|\tilde{\delta}\| \leq \frac{\epsilon}{1+2L_{\tilde{\gamma}}}$, with $L_{\tilde{\gamma}} = \|B_{\tilde{I} \tilde{\gamma}} (B_{\tilde{I} \tilde{\gamma}}^T B_{\tilde{I} \tilde{\gamma}})^{-1}\|$. Then, we must have

$$\|R_\lambda(\tilde{x}, \tilde{y}, \tilde{v})\| \leq \epsilon,$$

where $\tilde{y} = B\tilde{x}$ and \tilde{v} is constructed in (23), (24) and (25). Thus we can certify the optimality of \tilde{x} .

Proof. If \tilde{x} is an optimal solution to (26), then (\tilde{x}, \tilde{y}) is an optimal solution to

$$\begin{aligned} \min_{x \in \mathbb{R}^n, y \in \mathbb{R}^m} \quad & f(x) + \lambda p(y) + \langle x, \tilde{\delta} \rangle \\ \text{s.t.} \quad & Bx - y = 0. \end{aligned} \quad (27)$$

Then, there exists a $\tilde{z} \in \mathbb{R}^m$ which satisfies the following KKT system:

$$\begin{cases} (\nabla f(\tilde{x}))_{\tilde{\alpha}} + B_{\tilde{I}\tilde{\alpha}}^T \tilde{z}_{\tilde{I}} + B_{\tilde{I}^c\tilde{\alpha}}^T \tilde{z}_{\tilde{I}^c} + \tilde{\delta}_{\tilde{\alpha}} = 0, \\ (\nabla f(\tilde{x}))_{\tilde{\beta}} + B_{\tilde{I}^c\tilde{\beta}}^T \tilde{z}_{\tilde{I}^c} + \tilde{\delta}_{\tilde{\beta}} = 0, \\ (\nabla f(\tilde{x}))_{\tilde{\gamma}} + B_{\tilde{I}\tilde{\gamma}}^T \tilde{z}_{\tilde{I}} + B_{\tilde{I}^c\tilde{\gamma}}^T \tilde{z}_{\tilde{I}^c} + \tilde{\delta}_{\tilde{\gamma}} = 0, \\ \tilde{z} \in \partial(\lambda p(\tilde{y})), \\ B\tilde{x} - \tilde{y} = 0. \end{cases} \quad (28)$$

By Assumption 2 and the fact $\tilde{I}^c = \{i \in [m] \mid \tilde{y}_i \neq 0\}$, $(\partial(\lambda p(\tilde{y})))_{\tilde{I}^c}$ is a singleton. Thus we must have

$$\tilde{z}_{\tilde{I}^c} = \tilde{v}_{\tilde{I}^c}.$$

Therefore, $\tilde{v}_{\tilde{I}} = \tilde{z}_{\tilde{I}} - B_{\tilde{I}\tilde{\gamma}}^T (B_{\tilde{I}\tilde{\gamma}}^T B_{\tilde{I}\tilde{\gamma}})^{-1} \tilde{\delta}_{\tilde{\gamma}}$ is a solution to

$$(\nabla f(\tilde{x}))_{\tilde{\gamma}} + B_{\tilde{I}\tilde{\gamma}}^T \tilde{v}_{\tilde{I}} + B_{\tilde{I}^c\tilde{\gamma}}^T \tilde{v}_{\tilde{I}^c} = 0.$$

Since $\tilde{v}_{\tilde{I}}$ is a solution to (25), we have

$$\begin{aligned} \|\tilde{s}\| &= \|\tilde{v}_{\tilde{I}} - \Pi_{(\partial(\lambda p(\tilde{y})))_{\tilde{I}}}(\tilde{v}_{\tilde{I}})\| \\ &\leq \|\tilde{v}_{\tilde{I}} - \Pi_{(\partial(\lambda p(\tilde{y})))_{\tilde{I}}}(\tilde{v}_{\tilde{I}})\| \\ &= \|(\tilde{z}_{\tilde{I}} - B_{\tilde{I}\tilde{\gamma}}^T (B_{\tilde{I}\tilde{\gamma}}^T B_{\tilde{I}\tilde{\gamma}})^{-1} \tilde{\delta}_{\tilde{\gamma}}) - \Pi_{(\partial(\lambda p(\tilde{y})))_{\tilde{I}}}(\tilde{z}_{\tilde{I}} - B_{\tilde{I}\tilde{\gamma}}^T (B_{\tilde{I}\tilde{\gamma}}^T B_{\tilde{I}\tilde{\gamma}})^{-1} \tilde{\delta}_{\tilde{\gamma}})\| \\ &= \|-B_{\tilde{I}\tilde{\gamma}}^T (B_{\tilde{I}\tilde{\gamma}}^T B_{\tilde{I}\tilde{\gamma}})^{-1} \tilde{\delta}_{\tilde{\gamma}} + (\Pi_{(\partial(\lambda p(\tilde{y})))_{\tilde{I}}}(\tilde{z}_{\tilde{I}}) - \Pi_{(\partial(\lambda p(\tilde{y})))_{\tilde{I}}}(\tilde{z}_{\tilde{I}} - B_{\tilde{I}\tilde{\gamma}}^T (B_{\tilde{I}\tilde{\gamma}}^T B_{\tilde{I}\tilde{\gamma}})^{-1} \tilde{\delta}_{\tilde{\gamma}}))\| \\ &\leq 2\|B_{\tilde{I}\tilde{\gamma}}^T (B_{\tilde{I}\tilde{\gamma}}^T B_{\tilde{I}\tilde{\gamma}})^{-1} \tilde{\delta}_{\tilde{\gamma}}\| \\ &\leq 2L_{\tilde{\gamma}} \|\tilde{\delta}_{\tilde{\gamma}}\|. \end{aligned}$$

On the other hand, by the construction, we know that $(\tilde{x}, \tilde{y}, \tilde{v}, \tilde{s})$ satisfies the following KKT system

$$\begin{cases} (\nabla f(\tilde{x}))_{\tilde{\alpha}} + B_{\tilde{I}\tilde{\alpha}}^T \tilde{v}_{\tilde{I}} + B_{\tilde{I}^c\tilde{\alpha}}^T \tilde{v}_{\tilde{I}^c} + \tilde{\delta}_{\tilde{\alpha}} = 0, \\ (\nabla f(\tilde{x}))_{\tilde{\beta}} + B_{\tilde{I}^c\tilde{\beta}}^T \tilde{v}_{\tilde{I}^c} + \tilde{\delta}_{\tilde{\beta}} = 0, \\ (\nabla f(\tilde{x}))_{\tilde{\gamma}} + B_{\tilde{I}\tilde{\gamma}}^T \tilde{v}_{\tilde{I}} + B_{\tilde{I}^c\tilde{\gamma}}^T \tilde{v}_{\tilde{I}^c} = 0, \\ \tilde{v}_{\tilde{I}} - \tilde{s} \in (\partial(\lambda p(\tilde{y})))_{\tilde{I}}, \quad \tilde{v}_{\tilde{I}^c} \in (\partial(\lambda p(\tilde{y})))_{\tilde{I}^c}. \\ B\tilde{x} - \tilde{y} = 0. \end{cases}$$

Then, we have

$$\begin{aligned} \|R_{\lambda}(\tilde{x}, \tilde{y}, \tilde{v})\| &= \|(\nabla f(\tilde{x}) + B^T \tilde{v}, \quad \tilde{y} - \text{Prox}_{\lambda p}(\tilde{y} + \tilde{v}), \quad B\tilde{x} - \tilde{y})\| \\ &\leq \|(\tilde{\delta}_{\tilde{\alpha}}, \tilde{\delta}_{\tilde{\beta}}, 0)\| + \|\tilde{s}\| \\ &\leq \|\tilde{\delta}\| + 2L_{\tilde{\gamma}} \|\tilde{\delta}\| \\ &\leq \epsilon. \end{aligned}$$

This completes the proof of the theorem. \square

Algorithm 2 An enhanced adaptive sieving for solving (P_λ) with a fixed $\lambda > 0$

1: **Input:** a given hyper-parameter $\lambda > 0$ and a given tolerance $\epsilon > 0$.
2: **Output:** $(x^*(\lambda), y^*(\lambda), z^*(\lambda))$.
3: **Initialization:** Generate an initial index set by a predefined initialization strategy: $I^0(\lambda) \subseteq [m]$.
4: **for** $i = 0, 1, 2, \dots$ **do**
5: **1.** For the given index set $I^i(\lambda)$, construct the index partition $\{\alpha^i, \beta^i, \gamma^i\}$ and the corresponding $M_{\gamma^i \alpha^i}$.
6: **2.** Apply any well designed algorithm to solve problem $(RP_\lambda(I))$ with $\{I^i(\lambda), \alpha^i, \beta^i, \gamma^i, M_{\gamma^i \alpha^i}\}$ and obtain an inexact solution $(\hat{x}_{\alpha^i}^i, \hat{x}_{\beta^i}^i, \hat{y}_{(I^i)^c(\lambda)}^i, \hat{\xi}^i)$ which satisfies the corresponding KKT system (16) with the latent error terms $(\hat{\delta}_1^i, \hat{\delta}_2^i, \hat{\delta}_3^i)$ such that $\|\hat{\delta}_1^i\| + \|\hat{\delta}_2^i\| + \|\hat{\delta}_3^i\| \leq \epsilon$.
7: **3.** Recover a solution (\bar{x}^i, \bar{y}^i) by (7).
8: **if** $i > 1$ and $|F_\lambda(\bar{x}^i) - F_\lambda(\bar{x}^{i-1})| \leq \epsilon$ **then**
9: Define $(\tilde{x}^i, \tilde{y}^i) = (\bar{x}^i, B\bar{x}^i)$ and $\tilde{I}^i = \{i \in [m] \mid \tilde{y}_i = 0\}$. Construct $\{\tilde{\alpha}^i, \tilde{\beta}^i, \tilde{\gamma}^i, M_{\tilde{\gamma}^i \tilde{\alpha}^i}\}$.
10: Construct $(\tilde{v}^i, \tilde{s}^i)$ by (23), (24) and (25).
11: **if** $\|R_\lambda(\tilde{x}^i, \tilde{y}^i, \tilde{v})\| \leq \epsilon$ **then**
12: Set $(x^*(\lambda), y^*(\lambda), z^*(\lambda)) = (\tilde{x}^i, \tilde{y}^i, \tilde{v}^i)$.
13: **break.**
14: **end if**
15: **end if**
16: **4.** Recover a pair (\bar{u}^i, \bar{w}^i) by (9) and (12), respectively.
17: **if** $\|R_\lambda(\bar{x}^i, \bar{y}^i, \bar{u}^i)\| \leq \epsilon$ **then**
18: Set $(x^*(\lambda), y^*(\lambda), z^*(\lambda)) = (\bar{x}^i, \bar{y}^i, \bar{u}^i)$.
19: **break.**
20: **else**
21: Create $J^i(\lambda)$:

$$J^i(\lambda) = \{j \in I^i(\lambda) \mid \bar{u}_j^i \notin \partial(\lambda p(\bar{y}^i))_j\}. \quad (29)$$

22: **if** $J^i(\lambda) \neq \emptyset$ **then**
23: Update $I^{i+1}(\lambda)$ as:

$$I^{i+1}(\lambda) \leftarrow I^i(\lambda) \setminus J^i(\lambda).$$

24: **else**
25: Set $(x^*(\lambda), y^*(\lambda), z^*(\lambda)) = (\bar{x}^i, \bar{y}^i, \bar{u}^i)$.
26: **break.**
27: **end if**
28: **end if**
29: **end for**
30: **return** $(x^*(\lambda), y^*(\lambda), z^*(\lambda))$.

Now, we present the enhanced AS technique in Algorithm 2. As a byproduct of Theorem 1, Theorem 2 and Theorem 3, we have the following property.

Theorem 4. *For a given $\epsilon > 0$, Algorithm 2 is guaranteed to converge in finite number of iterations. The number of sieving iterations of Algorithm 2 is no more than that of Algorithm 1. Moreover, the obtained pair $(x^*(\lambda), y^*(\lambda), z^*(\lambda))$ is a solution to (P_λ) in the sense that*

$$\|R_\lambda(x^*(\lambda), y^*(\lambda), z^*(\lambda))\| \leq \epsilon.$$

Remark 6. *We close this subsection by making some remarks.*

1. *The enhanced adaptive sieving technique described in Algorithm 2 is a rigorous implementation of the aforementioned principal idea which simultaneously answers all the five questions we asked earlier in Section 2.*
2. *A natural question is, why should we still perform the sieving based on \bar{u} instead of \tilde{v} directly? Now we explain the reason. If we define*

$$\tilde{J}(\lambda) = \{j \in \tilde{I} \mid \tilde{v}_j \notin (\partial(\lambda p(\tilde{y})))_j\},$$

and assuming that $\tilde{J}(\lambda) \neq \emptyset$, we cannot guarantee that $\tilde{J}(\lambda) \cap I \neq \emptyset$, which is required to update the index set I .

3. *The main idea for the enhanced algorithm is to certify the optimality of the current solution if it is an optimal solution of (P_λ) . Then we can stop the sieving procedure earlier, comparing to Algorithm 1. It is a natural idea that we only try to certify the optimality of the current obtained solution if it is the solution to (P_λ) with high probability. This is implied by the condition $|F_\lambda(\bar{x}^i) - F_\lambda(\bar{x}^{i-1})| < \epsilon$, which is used in Algorithm 2. The reason we use the difference of the consecutive function values instead of the solution vectors is because the optimal solutions of (P_λ) may not be unique, but they all have the same objective function value.*
4. *In practice, the AS technique is sometimes better than the enhanced AS technique in terms of running time although the enhanced AS could potentially reduce the number of AS iterations. But of course, the enhanced AS technique is the one with a better theoretical guarantee. Detailed empirical comparison of these techniques can be found in the numerical experiments.*

2.4 An Accelerated Proximal Gradient Algorithm for Dual Variables Recovery

As aforementioned, the key step to recover the dual variables and applying the AS technique is to recover \bar{u} (or \tilde{v}) via solving the optimization problem (13) (or (25)). In this paper, we adopt the accelerated proximal gradient (APG) algorithm [1, 14] to solve it. Since the optimization problem (25) has the same form as (13), we use the problem (13) as an example.

First of all, we could rewrite the constrained optimization problem (13) equivalently as

$$\min_d h(d) + \delta_{\text{Null}(B_{I_\gamma}^T)}(d), \quad (30)$$

where $h(d) = \frac{1}{2} \|((\bar{u}_I)_0 + d) - \Pi_{\partial(\lambda p(\bar{y}))_I}((\bar{u}_I)_0 + d)\|^2$ and $\delta_{\text{Null}(B_{I_\gamma}^T)}(\cdot)$ is the indicator function of the Null space of $B_{I_\gamma}^T$.

In order to apply the APG algorithm, we need to derive the proximal mapping of the indicator function $\delta_{\text{Null}(B_{I_\gamma}^T)}(\cdot)$, which is the projection operator onto the null space of $B_{I_\gamma}^T$. Since $B_{I_\gamma}^T$ is of full row rank, the projection of a given vector $a \in \mathbb{R}^{|I|}$ onto the null space of $B_{I_\gamma}^T$ is computed by

$$\Pi_{\text{Null}(B_{I_\gamma}^T)}(a) = (I - B_{I_\gamma}(B_{I_\gamma}^T B_{I_\gamma})^{-1} B_{I_\gamma}^T) a.$$

On the other hand, the function $h(\cdot)$ is continuously differentiable and the gradient of $h(\cdot)$ is

$$\nabla h(d) = ((\bar{u}_I)_0 + d) - \Pi_{\partial(\lambda p(\bar{y}))_I}((\bar{u}_I)_0 + d) = \Pi_{(\partial(\lambda p(\bar{y}))_I)^\circ}((\bar{u}_I)_0 + d).$$

Here, $(\partial(\lambda p(\bar{y}))_I)^\circ$ is the polar of the closed convex set $\partial(\lambda p(\bar{y}))_I$ and the second equality comes from the Moreau identity [13]. Thus, $\nabla h(\cdot)$ is Lipschitz continuous with modulus 1 [30]. The APG algorithm for solving the optimization problem (30) is shown in Algorithm 3.

Algorithm 3 Accelerated proximal gradient algorithm for (30)

Input: $\epsilon > 0$ and maxiter.
Output: \bar{d} .
Initialization: $L = 1$, $d^0 = 0$, $\hat{d}^1 = d^0$, $k = 0$ and $t_1 = 1$.
while $k < \text{maxiter}$ **do**
 $k = k + 1$,
 $d^k = \Pi_{\text{Null}(B_{I\gamma}^T)}(\hat{d}^k - \frac{1}{L}\nabla h(\hat{d}^k))$,
 if $\max(\|d^k - d^{k-1}\|, \|((\bar{u}_I)_0 + d^k) - \Pi_{\partial(\lambda p(\bar{y}))_I}((\bar{u}_I)_0 + d^k)\|) \leq \epsilon$ **then**
 break,
 end if
 $t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}$,
 $\hat{d}^{k+1} = d^k + (\frac{t_k - 1}{t_{k+1}})(d^k - d^{k-1})$.
end while
 $\bar{d} = d^k$,
return \bar{d} .

It is well known that the sequence $\{d^k\}$ generated by the APG algorithm have the following $O(1/k^2)$ complexity [1, 14].

Theorem 5. *Let $\{d^k\}$ and $\{y^k\}$ be the sequences generated by Algorithm 3. Then for any $k \geq 1$, we have*

$$h(d^k) - h(d^*) \leq \frac{2\|d^*\|^2}{(k+1)^2}, \quad (31)$$

where d^* is any optimal solution to (30).

Remark 7. *Although we need to solve an additional optimization problem (30) in order to apply the AS technique, the computational cost is affordable. Now, we explain the key insights behind. In the enhanced AS technique, if we do obtain an optimal solution of (P_λ) via solving the current subproblem $(P_\lambda(I))$, then we must have $h(d^*) < \frac{\epsilon^2}{2}$ by Theorem 3. Moreover, $\|d^*\|$ must be relatively small. By the above complexity result, we could obtain an inexact solution to the problem (30) in several cheap iterations. On the other hand, if the objective function value of (30) is still large after several iterations (say 10 iterations), we can terminate the algorithm since this phenomenon indicates that we have not yet obtained an optimal solution to the problem (P_λ) . In other words, the current index set I is incorrect and we need to update it by removing violating indices. In short, although we need to solve an additional optimization problem, we only need to run APG for several iterations.*

The main computational cost for each iteration of APG is from two projections. For most of the commonly used regularizers p (for example, ℓ_1 norm, ℓ_2 norm), the projection of a given vector onto the subdifferential set is very cheap. On the other hand, in order to compute the projection onto the null space of $B_{I\gamma}^T$, the main computational cost is from computing $(B_{I\gamma}^T B_{I\gamma})^{-1}$. However, as we mentioned earlier, the matrix B is usually very sparse in many applications, the sparse Cholesky decomposition is not costly. Thus, the computational cost for one iteration of APG is affordable,

even for large scale problems. This is also one of the main reason for us to adopt APG to solve the optimization problem (30).

2.5 An Adaptive Sieving Technique for Solution Path

It is not difficult for us to generalize Algorithm 2 to obtain a solution path for problem (P_λ) with a sequence of parameters $\lambda_1 > \lambda_2 > \dots > \lambda_l > 0$. The key idea is that, if we obtain a solution $(x^*(\lambda_i), y^*(\lambda_i), z^*(\lambda_i))$ for (P_λ) with $\lambda = \lambda_i$, then, we can initialize the index set I_{i+1}^0 in Algorithm 2 for $\lambda = \lambda_{i+1}$ as

$$I_{i+1}^0 := \{k \in [m] \mid |(Bx^*(\lambda_i))_k| < \hat{\epsilon}\}, \quad (32)$$

where $\hat{\epsilon} > 0$ is a given tolerance. The algorithm for applying the AS technique (or the EAS technique) to generate a solution path is shown in Algorithm 4.

Algorithm 4 Generate solution path for (P_λ) with the AS technique (or the EAS technique)

Input: $\epsilon > 0, \hat{\epsilon} > 0$ and a sequence $\lambda_1 > \lambda_2 > \dots > \lambda_l > 0$.

Output: A solution path for (P_λ) : $\{(x^*(\lambda_1), y^*(\lambda_1), z^*(\lambda_1)), \dots, (x^*(\lambda_l), y^*(\lambda_l), z^*(\lambda_l))\}$.

Initialization: Initialize index set $I^0(\lambda_1) \subseteq [m]$ by a predefined initialization strategy.

for $k = 1, 2, \dots, l$ **do**

Step 1. Obtain $(x^*(\lambda_k), y^*(\lambda_k), z^*(\lambda_k))$ by calling Algorithm 1 (or Algorithm 2) with $\{\lambda, \epsilon, I^0(\lambda)\} = \{\lambda_k, \epsilon, I^0(\lambda_k)\}$.

if $k < l$ **then**

Step 2. Define

$$I^0(\lambda_{k+1}) := \{j \in [m] \mid |(Bx^*(\lambda_k))_j| < \hat{\epsilon}\}.$$

end if

end for

return $\{(x^*(\lambda_1), y^*(\lambda_1), z^*(\lambda_1)), \dots, (x^*(\lambda_l), y^*(\lambda_l), z^*(\lambda_l))\}$.

3 Adaptive Sieving and Enhanced Adaptive Sieving Technique for Convex Clustering

In this section, we will show how to apply the AS technique and the EAS technique on the convex clustering model (1).

Denote $\mathcal{E} := \{(i, j) \mid w_{ij} > 0, 1 \leq i < j \leq n\}$. Then $\mathcal{G} = ([n], \mathcal{E})$ forms an undirected graph and the weighted convex clustering model (1) is equivalent to:

$$\min_{X \in \mathbb{R}^{d \times N}} \frac{1}{2} \sum_{i=1}^N \|X_{:i} - A_{:i}\|_2^2 + \lambda \sum_{(i,j) \in \mathcal{E}} w_{ij} \|X_{:i} - X_{:j}\|_p. \quad (33)$$

We enumerate the index pairs in \mathcal{E} by the lexicographic order and denote by $l(i, j)$ for the pair (i, j) . Define the linear map $\mathcal{B} : \mathbb{R}^{d \times N} \rightarrow \mathbb{R}^{d \times |\mathcal{E}|}$ by

$$(\mathcal{B}(X))_{:,l(i,j)} = X_{:i} - X_{:j},$$

and the node-arc incidence matrix $J \in \mathbb{R}^{N \times |\mathcal{E}|}$ as

$$J_{k,l(i,j)} = \begin{cases} 1, & \text{if } k = i, \\ -1, & \text{if } k = j, \\ 0, & \text{otherwise.} \end{cases} \quad (34)$$

Then, for any given $X \in \mathbb{R}^{d \times N}$ and $Z \in \mathbb{R}^{d \times |\mathcal{E}|}$, we have

$$\mathcal{B}(X) = XJ, \quad \mathcal{B}^*(Z) = ZJ^T. \quad (35)$$

It is not difficult to see that the convex clustering model (33) is a special case of (2).

3.1 A Construction of the Reduced Problem

The main step for constructing the reduced subproblem is to construct the index sets α, β, γ and the corresponding matrix $M_{\gamma\alpha}$. For a given index set

$$I := \{l(i, j)\} \subseteq \{1, 2, \dots, |\mathcal{E}|\}, \quad (36)$$

we can construct a subgraph $\hat{\mathcal{G}} \subseteq \mathcal{G}$ with edges $\hat{\mathcal{E}} := \{(i, j) \mid l(i, j) \in I\}$ and all the corresponding nodes. Then, we can decompose the graph $\hat{\mathcal{G}}$ as

$$\hat{\mathcal{G}} = \hat{\mathcal{G}}_1 \cup \hat{\mathcal{G}}_2 \cup \dots \cup \hat{\mathcal{G}}_s,$$

where $\hat{\mathcal{G}}_i$ are disjoint connected subgraph of $\hat{\mathcal{G}}$. Denote the node index set of $\hat{\mathcal{G}}_i$ as $\hat{\mathcal{N}}_i$ and we define

$$\alpha_i = \min\{k \mid k \in \hat{\mathcal{N}}_i\}, \quad i = 1, 2, \dots, s.$$

Then, we can uniquely determine the index sets α, β and γ as

$$\alpha = \{\alpha_1, \dots, \alpha_s\}, \quad \beta = [N] \setminus (\hat{\mathcal{N}}_1 \cup \dots \cup \hat{\mathcal{N}}_s), \quad \text{and} \quad \gamma = (\hat{\mathcal{N}}_1 \cup \dots \cup \hat{\mathcal{N}}_s) \setminus \alpha.$$

The index sets α, β and γ have clear meanings in the convex clustering model. For a given index set $I \subseteq [|\mathcal{E}|]$ and the generated graph $\hat{\mathcal{G}}$, α_i is the index of the selected representative point for the i -th cluster identified by the connected component $\hat{\mathcal{G}}_i$. On the other hand, β is the collection of the indices of the isolated clusters which contain only a singleton.

Furthermore, we could have an explicit formula for $M_{\gamma\alpha} \in \mathbb{R}^{|\alpha| \times |\gamma|}$, which is given by

$$(M_{\gamma\alpha})_{ij} = \begin{cases} 1, & \text{if } j \in \hat{\mathcal{N}}_i, \\ 0, & \text{otherwise.} \end{cases}$$

Then

$$X_{:\gamma} = X_{:\alpha} M_{\gamma\alpha},$$

which actually maps the data points indexed by γ to the corresponding centroids with indices in the set α .

4 Numerical Experiments

In this section, we demonstrate the efficiency of the proposed solver independent AS technique and EAS technique via the important convex clustering model (33) (with $p = 2$). In this paper, we mainly focus on the numerical efficiency of our proposed techniques, readers can refer to [4, 10, 22] and the references therein for the performance of clustering by the convex clustering model (33). We test the AS technique with AMA [2], ADMM [2] and SSNAL [29], which are the three of the most popular algorithms for solving (33). Due to the limited length of the paper, we omit the details of these three algorithms but refer the readers to consult the aforementioned references. In our experiments, by default, we will generate the clustering path with $\lambda =: [10 : -0.2 : 1]$. The weights w_{ij} will be defined by the following Gaussian kernel with k -nearest neighbors (we choose $k = 10$ in our experiments):

$$w_{ij} = \begin{cases} \exp(-\frac{1}{2}\|A_{\cdot i} - A_{\cdot j}\|^2), & \text{if } (i, j) \in \mathcal{E}, \\ 0, & \text{if } (i, j) \notin \mathcal{E}, \end{cases} \quad (37)$$

where $\mathcal{E} = \{(i, j) \mid A_{\cdot i} \text{ is among } A_{\cdot j}\text{'s } k \text{ nearest neighbors}\}$.

For a fair comparison with the fast AMA algorithm, in this paper, we terminate all the algorithms based on the relative duality gap:

$$\eta = \frac{F_\lambda(X) - D_\lambda(Z)}{1 + |F_\lambda(X)| + |D_\lambda(Z)|} < \epsilon. \quad (38)$$

Here, $\epsilon > 0$ is a given tolerance, $F_\lambda(X)$ and $D_\lambda(Z)$ are the objective function value of the primal problem (2) and the dual problem (D_λ), respectively. We set $\epsilon = 10^{-6}$ in (38) and $\hat{\epsilon} = 2e-16$ in (32) by default in this paper. All our computational results are obtained by running MATLAB on a windows workstation (Intel Xeon E5-2680 @ 2.50GHz).

4.1 Simulated Data Sets

In this subsection, we provide some numerical results on the simulated two half-moon data, which is one of the most popular data sets for clustering.

First, we revisit the performance of fast AMA [2], ADMM [2] and SSNAL [29] for generating the clustering path directly. We implemented the three algorithms in MATLAB and tried our best to optimize the computations for a fair comparison³. As shown in Figure 1(a), SSNAL is the best among the three algorithms on this data set. However, unlike the statements in [2] stating that fast AMA is much better than ADMM, we actually observe some discrepancies in the performance. Fast AMA could not achieve the accuracy we set for most of the cases when n is relatively large (Figure 1(b)). For a fairer comparison, we revisit the performance of the three algorithms under the relatively low accuracy setting with $\epsilon = 10^{-4}$ (Figure 1(c), 1(d)), our numerical results show that ADMM is still better than fast AMA even in the low accuracy setting. Since the fast AMA has difficulty solving (33) to high accuracy, we focus on applying the AS technique with ADMM and SSNAL.

Now, we move on to present the numerical performance of the proposed AS technique. The details could be found in Figure 2. Our numerical results on the two half-moon data set show that

³Readers can find the implementations at: <https://blog.nus.edu.sg/mattohkc/software/convexclustering/>

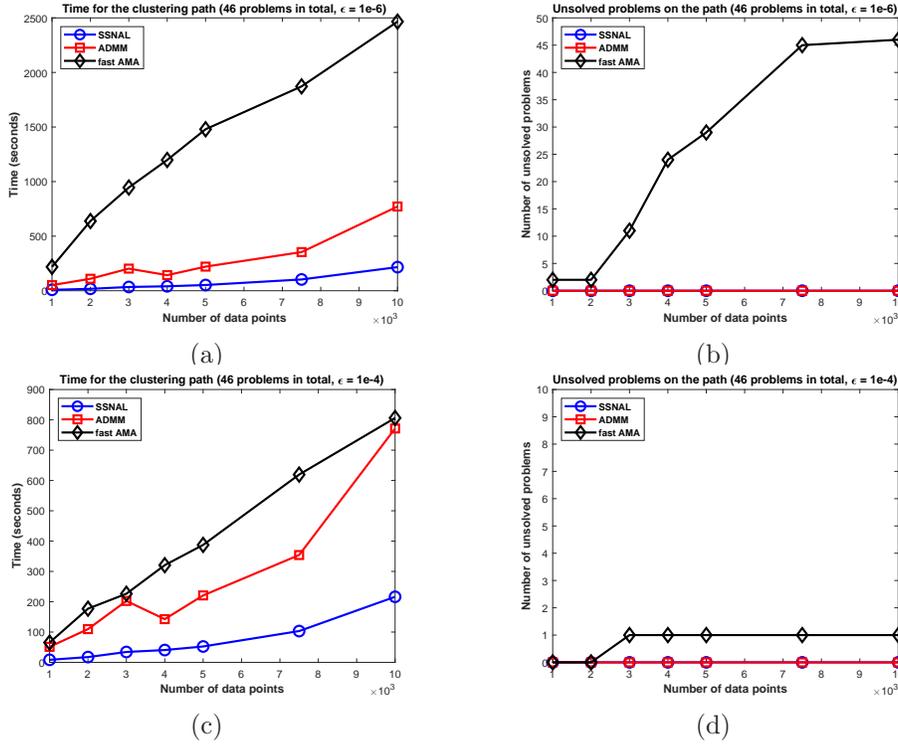


Figure 1: Numerical performance revisit on the two half moon data set.

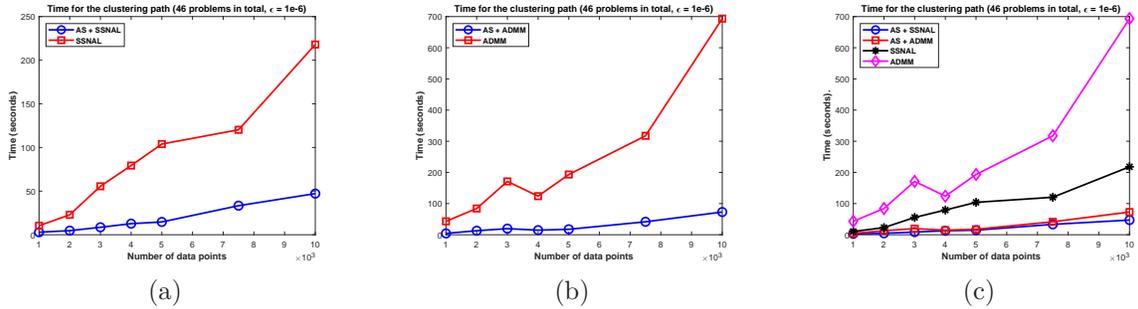


Figure 2: Numerical performance on the two half moon data set with $k = 10$.

the AS technique could accelerate the SSNAL and the ADMM by up to **4.8** times (Figure 2(a)) and **12.8** times (Figure 2(b)), respectively. With the help of the AS technique, AS+ADMM could even be comparable to AS+SSNAL (Figure 2(c)), which demonstrates the power of the AS technique for capturing the intrinsic structured sparsity of the convex clustering model. Since the AS technique can take advantage of the sparse structure to substantially reduce the dimension of the problem, we can apply the sparse Cholesky decomposition to solve the linear system involved in ADMM in a highly efficient way. This also partially demonstrates that ADMM is efficient to solve small scale convex clustering problems.

Next, we move on to present the empirical comparison between the AS technique and the EAS technique on the two half-moon data set. The results could be found in Figure 3. As shown in

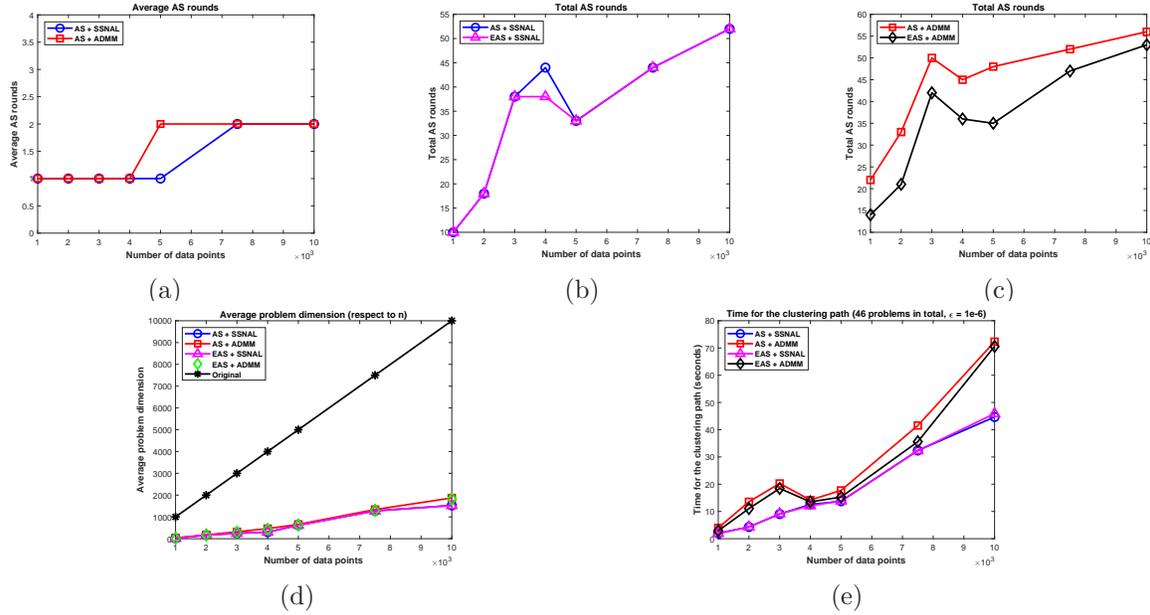


Figure 3: Numerical performance on the two half moon data set with $k = 10$.

Figure 3(a), the AS technique performs very well and the average AS rounds are very small, even for large scale problems. Furthermore, as shown in Figure 3(d), the sizes of the reduced problems are much smaller than those of the original problems. These are the main reasons why the AS technique can accelerate the algorithms. On the other hand, as one may imagine, the EAS could potentially early-stop the AS procedure. Thus the EAS technique could reduce the AS rounds and further accelerate the algorithms. This phenomenon is indeed observed in numerical experiments. As shown in Figure 3(b) and Figure 3(c), the EAS can have fewer AS rounds. The running time comparison could be found in Figure 3(e), which is consistent with our expectation.

Remark 8. We close this subsection by making some remarks.

1. One may be curious about the phenomenon where AS+SSNAL could have fewer AS rounds than AS+ADMM (Figure 3(a)). Now we try to give a plausible explanation. Although we set the same tolerance for terminating SSNAL and ADMM, the real accuracy achieved by the two algorithms are different. In our experiments, we observe that the SSNAL achieves higher accuracy than ADMM due to its faster convergence rate. This may be the main reason for the phenomenon shown in Figure 3(a).
2. As shown in Figure 3(e), EAS could further accelerate ADMM but may not be so for SSNAL, although it may early-terminate the AS procedure. This mainly because we need to solve additional auxiliary optimization problems in Algorithm 2 by the APG algorithm and it may spend more time than solving a few more reduced problems with SSNAL, because SSNAL is very efficient on this data set.
3. One may naturally agree that the AS technique and the EAS technique could be very powerful when λ is relatively large, since many data points are assigned to only a few clusters in this

case. However, since we generate the whole clustering path, some problems on the clustering path may not have this nice property when the parameter λ is small (which affects the efficiency of the AS technique and the EAS technique). But we still observe the distinctive advantages of them.

4.2 Real Data Sets

In this subsection, we will present the performance of the AS technique and the EAS technique for generating the clustering path on the MNIST dataset [7]. We adopt the preprocessing method described in [12], which applies a one hidden layer linear neural network to preprocess the raw images. Then, we apply the convex clustering model (33) on the preprocessed data. Our experiments is on the testing set of MNIST data and the dimension of the preprocessed data is 10×10000 . The details could be found in Table 1. From the results, we observe that the AS technique could

	SSNAL			ADMM		
	direct	with AS	with EAS	direct	with AS	with EAS
Time (seconds)	1207.7	156.3	157.3	1823.8	128.5	132.2
Total AS round	0	45	45	0	47	47
Average problem dimension	10000	1377	1377	10000	1389	1389

Table 1: Numerical performance on the MNIST data set.

accelerate the ADMM by up to **14.2** times and the SSNAL by up to **7.7** times. It is understandable that AS could be more attractive for ADMM, since the second-order sparsity embedded in the algorithm SSNAL has partially captured the structured sparsity already. Moreover, since the EAS technique does not reduce the sieving iterations on this data set comparing to the AS technique, the EAS technique will spend more time than the AS technique.

5 Conclusion

In this paper, we propose an AS technique and an enhanced AS technique, which are solver independent, for convex optimization problems with structured sparsity. The proposed techniques can accelerate various optimization algorithms by substantially reducing the dimension of the problems that need to be solved. Numerical performance on the convex clustering model has demonstrated the high efficiency of the proposed dimension reduction techniques. We also established a finite convergence property of the AS and enhanced AS techniques in this paper. However, we should note that in the worst-case, the AS technique may sieve all the indices. But based on our empirical evaluation, one can say that the AS technique works very well in practice. Thus, as a future research topic, we will make efforts to analyze the average-case complexity of the AS and the enhanced AS technique.

References

- [1] A. BECK AND M. TEBoulLE, *A fast iterative shrinkage-thresholding algorithm for linear inverse problems*, SIAM journal on imaging sciences, 2 (2009), pp. 183–202.

- [2] E. C. CHI AND K. LANGE, *Splitting methods for convex clustering*, Journal of Computational and Graphical Statistics, 24 (2015), pp. 994–1013.
- [3] E. C. CHI AND S. STEINERBERGER, *Recovering trees with convex clustering*, SIAM Journal on Mathematics of Data Science, 1 (2019), pp. 383–407.
- [4] T. D. HOCKING, A. JOULIN, F. BACH, AND J.-P. VERT, *Clusterpath an algorithm for clustering using convex fusion penalties*, in 28th international conference on machine learning, 2011, p. 1.
- [5] T. JIANG, *Sum-of-norms clustering: theoretical guarantee and post-processing*, Master’s thesis, University of Waterloo, 2020.
- [6] T. JIANG, S. VAVASIS, AND C. W. ZHAI, *Recovery of a mixture of Gaussians by sum-of-norms clustering*, Journal of Machine Learning Research, 21 (2020), pp. 1–16.
- [7] Y. LECUN, *The MNIST database of handwritten digits*, <http://yann.lecun.com/exdb/mnist/>, (1998).
- [8] M. LIN, Y. YUAN, D. SUN, AND K.-C. TOH, *Adaptive sieving with PPDNA: Generating solution paths of exclusive lasso models*, arXiv preprint arXiv:2009.08719, (2020).
- [9] Y. LIN AND S. CHEN, *A centroid auto-fused hierarchical fuzzy c-means clustering*, IEEE Transactions on Fuzzy Systems, (2020).
- [10] F. LINDSTEN, H. OHLSSON, AND L. LJUNG, *Clustering using sum-of-norms regularization: With application to particle filter output computation*, in 2011 IEEE Statistical Signal Processing Workshop (SSP), IEEE, 2011, pp. 201–204.
- [11] S. LLOYD, *Least squares quantization in pcm*, IEEE transactions on information theory, 28 (1982), pp. 129–137.
- [12] D. G. MIXON, S. VILLAR, AND R. WARD, *Clustering subgaussian mixtures by semidefinite programming*, Information and Inference: A Journal of the IMA, 6 (2017), pp. 389–415.
- [13] J.-J. MOREAU, *Proximité et dualité dans un espace Hilbertien*, Bulletin de la Société mathématique de France, 93 (1965), pp. 273–299.
- [14] Y. E. NESTEROV, *A method for solving the convex programming problem with convergence rate $o(1/k^2)$* , in Dokl. akad. nauk Sssr, vol. 269, 1983, pp. 543–547.
- [15] A. Y. NG, M. JORDAN, Y. WEISS, ET AL., *On spectral clustering: analysis and an algorithm*, Proceedings of IEEE Neural Information Processing Systems (NIPS), (2002).
- [16] A. PANAHİ, D. DUBHASHI, F. D. JOHANSSON, AND C. BHATTACHARYYA, *Clustering by sum of norms: Stochastic incremental algorithm, convergence and cluster recovery*, in International conference on machine learning, PMLR, 2017, pp. 2769–2777.
- [17] K. PELCKMANS, J. DE BRABANTER, J. A. SUYKENS, AND B. DE MOOR, *Convex clustering shrinkage*, in PASCAL Workshop on Statistics and Optimization of Clustering Workshop, 2005.

- [18] S. A. SHAH AND V. KOLTUN, *Robust continuous clustering*, Proceedings of the National Academy of Sciences, 114 (2017), pp. 9814–9819.
- [19] Y. SHE, *Sparse regression with exact clustering*, Electronic Journal of Statistics, 4 (2010), pp. 1055 – 1096.
- [20] J. SHI AND J. MALIK, *Normalized cuts and image segmentation*, IEEE Transactions on pattern analysis and machine intelligence, 22 (2000), pp. 888–905.
- [21] M. SOLTANOLKOTABI, E. ELHAMIFAR, E. J. CANDÉS, ET AL., *Robust subspace clustering*, Annals of Statistics, 42 (2014), pp. 669–699.
- [22] D. SUN, K.-C. TOH, AND Y. YUAN, *Convex clustering: model, theoretical guarantee and efficient algorithm*, Journal of Machine Learning Research, 22 (2021), pp. 1–32.
- [23] K. M. TAN AND D. WITTEN, *Statistical properties of convex clustering*, Electronic journal of statistics, 9 (2015), pp. 2324–2347.
- [24] R. TIBSHIRANI, *Regression shrinkage and selection via the lasso*, Journal of the Royal Statistical Society: Series B (Methodological), 58 (1996), pp. 267–288.
- [25] R. TIBSHIRANI, M. SAUNDERS, S. ROSSET, J. ZHU, AND K. KNIGHT, *Sparsity and smoothness via the fused lasso*, Journal of the Royal Statistical Society: Series B (Statistical Methodology), 67 (2005), pp. 91–108.
- [26] S. VASSILVITSKII AND D. ARTHUR, *k-means++: The advantages of careful seeding*, in Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms, 2006, pp. 1027–1035.
- [27] R. VIDAL, *Subspace clustering*, IEEE Signal Processing Magazine, 28 (2011), pp. 52–68.
- [28] M. YUAN AND Y. LIN, *Model selection and estimation in regression with grouped variables*, Journal of the Royal Statistical Society: Series B (Statistical Methodology), 68 (2006), pp. 49–67.
- [29] Y. YUAN, D. SUN, AND K.-C. TOH, *An efficient semismooth Newton based algorithm for convex clustering*, in International Conference on Machine Learning, PMLR, 2018, pp. 5718–5726.
- [30] E. H. ZARANTONELLO, *Projections on convex sets in Hilbert space and spectral theory: Part i. projections on convex sets: Part ii. spectral theory*, in Contributions to nonlinear functional analysis, Elsevier, 1971, pp. 237–424.
- [31] Y. ZHOU, R. JIN, AND S. C.-H. HOI, *Exclusive lasso for multi-task feature selection*, in Proceedings of the thirteenth international conference on artificial intelligence and statistics, 2010, pp. 988–995.
- [32] C. ZHU, H. XU, C. LENG, AND S. YAN, *Convex optimization procedure for clustering: Theoretical revisit*, Advances in Neural Information Processing Systems, 27 (2014), pp. 1619–1627.