

Gradient Projection Newton Algorithm for Sparse Collaborative Learning Using Synthetic and Real Datasets of Applications

Jun Sun, Lingchen Kong, and Shenglong Zhou

Department of Applied Mathematics, Beijing Jiaotong University, CN

Department of EEE, Imperial College London, UK

Abstract: Exploring the relationship among multiple sets of data from one same group enables practitioners to make better decisions in medical science and engineering. In this paper, we propose a sparse collaborative learning (SCL) model, an optimization with double-sparsity constraints, to process the problem with two sets of data and a shared response variable. It is capable of dealing with the classification problems or the regression problems dependent on the discreteness of the response variable as well as exploring the relationship between two datasets simultaneously. To solve SCL, we first present some necessary and sufficient optimality conditions and then design a gradient projection Newton algorithm which has proven to converge to a unique locally optimal solution globally with at least a quadratic convergence rate. Finally, the reported numerical experiments illustrate the efficiency of the proposed method.

Keywords: Sparse collaborative learning, double-sparsity, stationary point, gradient projection Newton, convergence analysis, numerical experiment

1 Introduction

There are many scenarios where datasets from the same group can be collected from various sources. Therefore, they differ but interact [14, 27, 30, 33]. For example, a researcher studying cancer outcomes may collect gene expression data and copy number data from a group of patients. The traditional approaches to do predictions are either merging two datasets or using two datasets separately. Both ways ignore the fact that they are from different sources with different meanings (e.g., gene expression and copy number). As stated in [26], exploring the relationship between sources allows for extracting informative biomarkers and improving clinical outcome predictions. Motivated by such practical applications, in this paper, we propose the following sparse collaborative learning (SCL) problem:

$$\begin{aligned} \min_{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2} \quad & \frac{1}{n} \left[a \cdot \ell(\boldsymbol{\theta}_1; X, \mathbf{y}) + b \cdot \ell(\boldsymbol{\theta}_2; Z, \mathbf{y}) + \frac{c}{2} \|X\boldsymbol{\theta}_1 - Z\boldsymbol{\theta}_2\|^2 \right] =: f(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \\ \text{s.t.} \quad & \|\boldsymbol{\theta}_1\|_0 \leq s_1, \quad \|\boldsymbol{\theta}_2\|_0 \leq s_2, \end{aligned} \tag{1.1}$$

where $\ell(\cdot)$ is a general loss function, $X \in \mathbb{R}^{n \times p_1}$, $Z \in \mathbb{R}^{n \times p_2}$ are two datasets from two different sources and $\mathbf{y} \in \mathbb{R}^n$ is the shared response, n is the sample/subject size, and p_1, p_2 represent the feature/variable sizes of two datasets. Here, $\|\boldsymbol{\theta}\|_0$ is the zero norm of $\boldsymbol{\theta}$, counting the number of its nonzero elements, $s_1 \ll p_1$, $s_2 \ll p_2$ are two integers representing the prior information on the upper bounds of the signal sparsity, a, b and c are positive parameters, and $\|\cdot\|$ represents the Euclidean norm. SCL models have been applied into many real-world applications, such as face recognition by using a mixture of synthetic and real images with dynamic weight [9], medical diagnosis including schizophrenia, Alzheimer's disease, or various neurocognitive phenotypes by using genetic and imaging data [11, 38, 12].

Two typical examples of ℓ will be investigated in this paper. When ℓ is the linear regression loss,

$$\ell_{lin}(\boldsymbol{\theta}; X, \mathbf{y}) := \frac{1}{2} \sum_{i=1}^n (y_i - \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle)^2,$$

SCL is called sparse collaborative regression (SCoRe [11]) usually working for the continuous response \mathbf{y} . Here, $\langle \mathbf{x}, \mathbf{z} \rangle$ is the inner product of two vectors \mathbf{x} and \mathbf{z} and \mathbf{x}_i is a column vector corresponding the i -th row of X . SCoRe is a combination of linear regression and canonical correlation analysis (CCA). The former makes predictions via employing two different types of datasets and the latter explores the relationship between them. Examples of employing ℓ_{lin} include CoRe [5], multi-task CoRe [38] and the models studied in [9, 12].

We note that the aforementioned models based on ℓ_{lin} aimed to process the continuous response \mathbf{y} . However, various real-world applications involve discrete responses, in particular for those in classification problems including the severity of the disease, whether or not to die and to name a few. Under such circumstances, linear regression-based models are unlikely to provide accurate predictions and hence it is necessary to consider the logistic regression loss defined by,

$$\ell_{log}(\boldsymbol{\theta}; X, \mathbf{y}) := \sum_{i=1}^n \left(\log(1 + \exp\langle \mathbf{x}_i, \boldsymbol{\theta} \rangle) - y_i \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle \right).$$

SCL with such a loss is called the sparse logistic collaborative regression (SLCoRe), which can be used to deal with datasets with discrete response \mathbf{y} . SLCoRe is a combination of logistic regression and CCA, aiming at classifying the samples in each of the two datasets while exploring the relationship between them. It is well-known that discrete responses are frequently involved in classification problems, while most of the existing classification methods including support vector machines [16, 34] and logistic regression [15, 22, 28, 29] only target one dataset. Very little work makes predictions for multiple sets of data and explores the relationship among them at the same time.

However, to accurately characterize the sparsity, it is suggested to impose the sparsity constraints directly instead of using the approximations/regularizations. For example, Beck and Eldar [4] thoroughly studied a general sparsity-constrained optimization model and developed the famous iterative hard thresholding algorithm, in the meanwhile, Bahmani et al. [2] and Plan et al. [22] investigated the logistic regression model with sparsity constraints. After which there is a vast body of work on developing optimization algorithms and understanding the properties of various sparse estimators for the sparsity constrained optimization [28, 20, 21, 37, ?]. We emphasize that all those work aimed at addressing applications with single datasets rather than multiple datasets.

In this paper, we study two typical examples of SCL: SLCoRe with $\ell = \ell_{log}$ and SCoRe with $\ell = \ell_{lin}$. All results to be established are based on these two models. The main contributions of the paper are summarized as follows:

- I) We propose a unified framework, SCL, for the problems with discrete or continuous response variables and two different sets of data. It can classify or predict the data in each dataset, and explore the relationship between the two datasets. New model (1.1) exploits the sparsity constraints directly, which enables to select a sufficiently small portion of informative features in each dataset provided that s_1 and s_2 are small enough.

- II) We investigate the first-order necessary and sufficient optimality conditions (see Theorem 3.1 and Theorem 3.2) for SCL as well as the existence and the uniqueness of its solution (see Theorem 3.3). One of the optimality conditions is associated with the α -stationary point seen Definition 3.1 that allows for algorithmic design conveniently.
- III) We develop a gradient projection Newton algorithm (GPNA) that combines the gradient projection motivated by the α -stationary point and the Newton step to accelerate the convergence. We prove that GPNA not only converges to a unique local minimizer of problem (1.1) globally (see Theorem 4.1) but also has a quadratic convergence rate for SLCoRe and termination within finite steps for SCoRe (see Theorem 4.2) under a mild assumption. These nice convergence properties indicate that our proposed algorithm should behave excellently in terms of high accuracy and speed, which is testified by its outstanding numerical performance.

We note that SCL problem (1.1) has a close link to the multi-model problem where multiple models based on the learned data distributions are used to make predictions [10, 23, 31, 36]. In contrast, SCL focuses on two groups of data not only for the prediction but also for exploring their inter-group relationships. To this end, if two groups of data in the dataset are known, then SCL with $c = 0$ (namely, no inter-group relationships are investigated) in problem (1.1) can be deemed as a special case of the multi-model problem.

To end this section, we present the organization of this paper. The next section describes the notation that will be employed through this paper and displays some properties of the objective function of problem (1.1). In section 3, we establish the first-order necessary and sufficient optimality conditions as well as the existence and the uniqueness of the solutions to problem (1.1). The algorithm GPNA and its convergence properties are provided in section 4. Numerical experiments on synthetic and real data are reported in section 5, and some concluding remarks are given in the section 6.

2 Preliminaries

Before giving the main results, we define some notations that will be employed throughout the paper. Let $[p] := \{1, 2, \dots, p\}$, $[n] := \{1, 2, \dots, n\}$. We denote sparse set Σ_s^p in \mathbb{R}^p by

$$\Sigma_s^p := \{\boldsymbol{\theta} \in \mathbb{R}^p : \|\boldsymbol{\theta}\|_0 \leq s\},$$

where $s \ll p$ is an integer. For a vector $\boldsymbol{\theta}$, denote its neighborhood with a radius δ by $N(\boldsymbol{\theta}, \delta) := \{\mathbf{u} \in \mathbb{R}^p : \|\boldsymbol{\theta} - \mathbf{u}\| < \delta\}$, and its support set by $\Gamma(\boldsymbol{\theta}) := \{i \in [p] : \theta_i \neq 0\}$. The complement set of Γ is written as $\bar{\Gamma}$. For a given set T , its spanned subspace of \mathbb{R}^p is denoted by $\mathbb{R}_T^p := \{\boldsymbol{\theta} \in \mathbb{R}^p : \Gamma(\boldsymbol{\theta}) \subseteq T\}$. Let $\boldsymbol{\theta}_\Gamma$ be the subvector of $\boldsymbol{\theta}$ indexed on Γ . We merge two vectors $\boldsymbol{\theta}_1$ and $\boldsymbol{\theta}_2$ as a single column vector via $(\boldsymbol{\theta}_1; \boldsymbol{\theta}_2) := (\boldsymbol{\theta}_1^\top \boldsymbol{\theta}_2^\top)^\top$. Finally, for a matrix $A \in \mathbb{R}^{n \times p}$, let $\lambda_{\max}(A)$ and $\lambda_{\min}(A)$ present its largest and smallest eigenvalue, respectively, and A_{TJ} denotes the sub-matrix containing rows indexed by T and columns indexed by J . In particular, $A_T := A_{T[p]}$ and $A_{:J} := A_{[n]J}$.

To characterize the projection of $\boldsymbol{\theta}$ onto Σ_s^p , we denote θ_i^\downarrow the i th largest element in magnitude of $\boldsymbol{\theta}$. Based on this, projection $\Pi_{\Sigma_s^p}(\boldsymbol{\theta})$ that is given by

$$\Pi_{\Sigma_s^p}(\boldsymbol{\theta}) := \underset{\mathbf{u} \in \Sigma_s^p}{\operatorname{argmin}} \|\boldsymbol{\theta} - \mathbf{u}\|$$

can be derived as follows: If $\theta_s^\downarrow = 0$ or $\theta_s^\downarrow > \theta_{s+1}^\downarrow$, then it is unique, i.e.,

$$(\Pi_{\Sigma_s^p}(\boldsymbol{\theta}))_i = \begin{cases} \theta_i, & |\theta_i| \geq \theta_s^\downarrow, \\ 0, & |\theta_i| < \theta_s^\downarrow. \end{cases}$$

If there are more than one equal to θ_s^\downarrow , we can choose any one of them and let the rest be 0. If $\theta_s^\downarrow = \theta_{s+1}^\downarrow \neq 0$, then

$$(\Pi_{\Sigma_s^p}(\boldsymbol{\theta}))_i = \begin{cases} \theta_i, & |\theta_i| > \theta_s^\downarrow, \\ \theta_i \text{ or } 0, & |\theta_i| = \theta_s^\downarrow, \\ 0, & |\theta_i| < \theta_s^\downarrow. \end{cases}$$

For example, for $\boldsymbol{\theta} = \{2, 4, 3, -3, 1\}$ and $\Sigma_2^5 = \{\mathbf{u} \in \mathbb{R}^5 : \|\mathbf{u}\|_0 \leq 2\}$, we have $\Pi_{\Sigma_2^5}(\boldsymbol{\theta}) = (0, 4, 3, 0, 0)^\top$ or $(0, 4, 0, -3, 0)^\top$.

Below are some concepts that will be used in this paper.

Definition 2.1 (*s*-regularity [4]). *A matrix $A \in \mathbb{R}^{n \times p}$ is called *s*-regular if its any *s* columns are linearly independent.*

Definition 2.2 (Strong smoothness [13]). *If function f is continuously differentiable, then for any $\boldsymbol{\theta}, \mathbf{d} \in \mathbb{R}^p$, we say that function f is strongly smooth on \mathbb{R}^p with a parameter $L_f > 0$ if it holds that*

$$f(\boldsymbol{\theta} + \mathbf{d}) \leq f(\boldsymbol{\theta}) + \langle \nabla f(\boldsymbol{\theta}), \mathbf{d} \rangle + (L_f/2)\|\mathbf{d}\|^2.$$

Definition 2.3 (Restricted strong convexity [2, 37, 1, 25]). *If function f is twice continuously differentiable, then for any $\boldsymbol{\theta}, \mathbf{d} \in \Sigma_r^p$ satisfying $\boldsymbol{\theta} + \mathbf{d} \in \Sigma_r^p$, we say that function f is *r*-restricted strongly convex on Σ_r^p with a parameter $l_f > 0$ if it holds that*

$$f(\boldsymbol{\theta} + \mathbf{d}) \geq f(\boldsymbol{\theta}) + \langle \nabla f(\boldsymbol{\theta}), \mathbf{d} \rangle + (l_f/2)\|\mathbf{d}\|^2 \quad \text{or} \quad \langle \mathbf{d}, \nabla^2 f(\boldsymbol{\theta})\mathbf{d} \rangle \geq (l_f/2)\|\mathbf{d}\|^2.$$

*If these conditions hold for $l_f = 0$, then f is called *r*-restricted convex on Σ_r^p .*

We now give some properties of f in problem (1.1), including the strong smoothness and restricted strong convexity as well as the Lipschitz continuity of its gradient and Hessian matrix.

Proposition 2.1. *Let $\boldsymbol{\theta} := (\boldsymbol{\theta}_1; \boldsymbol{\theta}_2)$ and $\ell = \ell_{\log}$. Objective function f in problem (1.1) has the following properties.*

- 1) *It is convex, twice continuously differentiable and strongly smooth with parameter L_f given by*

$$L_f := \lambda_{\max} \left(\frac{1}{n} \begin{bmatrix} (a/4 + c)X^\top X & -cX^\top Z \\ -cZ^\top X & (b/4 + c)Z^\top Z \end{bmatrix} \right),$$

which indicates that ∇f is Lipschitz continuous with parameter L_f for any $\boldsymbol{\theta}$ and $\boldsymbol{\theta}'$,

$$\|\nabla f(\boldsymbol{\theta}) - \nabla f(\boldsymbol{\theta}')\| \leq L_f \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|.$$

- 2) *Its Hessian matrix $\nabla^2 f(\boldsymbol{\theta})$ takes the form of*

$$\nabla^2 f(\boldsymbol{\theta}) = \frac{1}{n} \begin{bmatrix} X^\top (aD_1 + cI)X & -cX^\top Z \\ -cZ^\top X & Z^\top (bD_2 + cI)Z \end{bmatrix},$$

where I is the identity matrix, D_1 and D_2 are two diagonal matrices with

$$(D_1)_{ii} = \frac{\exp\langle \mathbf{x}_i, \boldsymbol{\theta}_1 \rangle}{(1 + \exp\langle \mathbf{x}_i, \boldsymbol{\theta}_1 \rangle)^2}, \quad i \in [n],$$

$$(D_2)_{ii} = \frac{\exp\langle \mathbf{z}_i, \boldsymbol{\theta}_2 \rangle}{(1 + \exp\langle \mathbf{z}_i, \boldsymbol{\theta}_2 \rangle)^2}, \quad i \in [n].$$

Moreover, $\nabla^2 f(\cdot)$ is Lipschitz continuous with constant C_f , namely,

$$\|\nabla^2 f(\boldsymbol{\theta}) - \nabla^2 f(\boldsymbol{\theta}')\| \leq C_f \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|, \quad (2.1)$$

for any $\boldsymbol{\theta}$ and $\boldsymbol{\theta}'$, where

$$C_f := \frac{3\sqrt{2}}{n} \max \left\{ a \max_{i \in [n]} \|\mathbf{x}_i\|_1 \lambda_{\max}(X^\top X), b \max_{i \in [n]} \|\mathbf{z}_i\|_1 \lambda_{\max}(Z^\top Z) \right\}.$$

3) If matrix $[X \ Z]$ is $(s_1 + s_2)$ -regular, then it is $(s_1 + s_2)$ -restricted strongly convex on $\Sigma_{s_1+s_2}^{p_1+p_2}$ with a positive parameter l_f given by

$$l_f := \min_{|T| \leq s_1+s_2} \lambda_{\min} \left(\frac{c}{n} \begin{bmatrix} X^\top X & -X^\top Z \\ -Z^\top X & Z^\top Z \end{bmatrix}_{TT} \right). \quad (2.2)$$

Proof. 1) It is easy to see that f is convex and twice continuously differentiable. Since $t/(1+t)^2 \leq 1/4$ for any $t \geq 0$, it follows $\lambda_{\max}(\nabla^2 f(\boldsymbol{\theta})) \leq L_f$ for any $\boldsymbol{\theta} \in \mathbb{R}^{p_1+p_2}$. This can show that the gradient of f is Lipschitz continuous with parameter L_f immediately.

2) It follows from [28, Lemma A.3] that $\nabla^2 \ell(\boldsymbol{\theta}_1; X)$ and both $\nabla^2 \ell(\boldsymbol{\theta}_2; Z)$ are Lipschitz continuous with constants

$$C_1 := (3/n) \max_{i \in [n]} \|\mathbf{x}_i\|_1 \lambda_{\max}(X^\top X), \quad C_2 := (3/n) \max_{i \in [n]} \|\mathbf{z}_i\|_1 \lambda_{\max}(Z^\top Z).$$

Then we have

$$\begin{aligned} \|\nabla^2 f(\boldsymbol{\theta}) - \nabla^2 f(\boldsymbol{\theta}')\| &= \|a\nabla^2 \ell(\boldsymbol{\theta}_1; X) + b\nabla^2 \ell(\boldsymbol{\theta}_2; Z) - a\nabla^2 \ell(\boldsymbol{\theta}'_1; X) - b\nabla^2 \ell(\boldsymbol{\theta}'_2; Z)\| \\ &\leq aC_1 \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}'_1\| + bC_2 \|\boldsymbol{\theta}_2 - \boldsymbol{\theta}'_2\| \\ &\leq \max\{aC_1, bC_2\} (\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}'_1\| + \|\boldsymbol{\theta}_2 - \boldsymbol{\theta}'_2\|) \\ &\leq \sqrt{2} \max\{aC_1, bC_2\} \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|. \end{aligned}$$

3) If matrix $[X \ Z]$ is $(s_1 + s_2)$ -regular, then so is matrix $[X \ -Z]$. Note that

$$\nabla^2 f(\boldsymbol{\theta}) = \frac{1}{n} \begin{bmatrix} aX^\top D_1 X & 0 \\ 0 & bZ^\top D_2 Z \end{bmatrix} + \frac{c}{n} \begin{bmatrix} X^\top X & -X^\top Z \\ -Z^\top X & Z^\top Z \end{bmatrix} =: A + B.$$

Clearly, both A and B are positive semi-definite. Moreover, $B = (c/n)[X \ -Z]^\top [X \ -Z]$. According to the $(s_1 + s_2)$ -regular of the matrix $[X \ Z]$, we can get B_{TT} is positive definite. Therefore, for any $\mathbf{d} := (\mathbf{d}_1; \mathbf{d}_2) \neq 0$ with $\|\mathbf{d}_1\|_0 \leq s_1$ and $\|\mathbf{d}_2\|_0 \leq s_2$, we have

$$\langle \mathbf{d}, \nabla^2 f(\boldsymbol{\theta}) \mathbf{d} \rangle = \langle \mathbf{d}, (A + B) \mathbf{d} \rangle \geq \langle \mathbf{d}, B \mathbf{d} \rangle \geq l_f \|\mathbf{d}\|^2 > 0.$$

This displays that the $(s_1 + s_2)$ -restricted strong convexity of $f(\boldsymbol{\theta})$ on $\Sigma_{s_1+s_2}^{p_1+p_2}$. The proof is complete. \square

\square

We note that the classical logistic regression which has been shown to be only strictly convex instead of being restricted strongly convex even though the assumption of the regularity of the sample matrix is imposed. However, the objective function of SLCoRe can be restricted strongly convex if the sample matrix is regular. In addition, if we only have one dataset, SLCoRe will degenerate into the classical sparse logistic regression. At this point, see the example in [28], similar results can be obtained. Similarly, for the objective function of SCoRe, we easily obtain the following results.

Proposition 2.2. *Let $\boldsymbol{\theta} := (\boldsymbol{\theta}_1; \boldsymbol{\theta}_2)$ and $\ell = \ell_{lin}$. Objective function f in (1.1) is convex, twice continuously differentiable and has Hessian matrix $\nabla^2 f(\boldsymbol{\theta})$ in the form of*

$$\nabla^2 f(\boldsymbol{\theta}) = \frac{1}{n} \begin{bmatrix} (a+c)X^\top X & -cX^\top Z \\ -cZ^\top X & (b+c)Z^\top Z \end{bmatrix} =: Q.$$

Moreover, it is strongly smooth with parameter $L_f := \lambda_{\max}(Q)$ and thus ∇f is Lipschitz continuous with parameter L_f . If $[X \ Z]$ is $(s_1 + s_2)$ -regular, then it is $(s_1 + s_2)$ -restricted strongly convex on $\Sigma_{s_1+s_2}^{p_1+p_2}$ with a positive parameter $l_f > 0$ given by

$$l_f := \min_{|T| \leq s_1+s_2} \lambda_{\min}(Q_{TT}). \quad (2.3)$$

It is worth mentioning that the main theorems in the sequel are established based on the assumption of s -regularity. So, to end this section, we would like to see which types of matrices $[X \ Z]$ could satisfies s -regularity. To proceed with that, we introduce the famous Restricted Isometry Property (RIP, [6]). A matrix $\Phi \in \mathbb{R}^{n \times p}$ is said to satisfy s -order RIP, if there exists a constant $\delta_s \in [0, 1)$ such that

$$(1 - \delta_s) \|\boldsymbol{\theta}\|^2 \leq \|\Phi \boldsymbol{\theta}\|^2 \leq (1 + \delta_s) \|\boldsymbol{\theta}\|^2$$

for all vectors $\boldsymbol{\theta} \in \Sigma_s^p$. This definition is equivalent to

$$(1 - \delta_s) \leq \lambda_{\min}(\Phi_{:T}^\top \Phi_{:T}) \leq \lambda_{\max}(\Phi_{:T}^\top \Phi_{:T}) \leq (1 + \delta_s), \quad \forall |T| \leq s.$$

Therefore, matrices satisfying s -order RIP must satisfy s -regularity. On the other hand, it has proven in [7, 3] that random Gaussian matrix, random binary matrix, and Fourier matrix satisfy s -order RIP with a high probability when s is small enough. Hence, these matrices also satisfy s -regularity.

3 Optimality Conditions

This section establishes the optimality conditions of SCL being useful for the algorithmic development, before which, for notational convenience, we define

$$\begin{aligned} \boldsymbol{\theta} &:= (\boldsymbol{\theta}_1; \boldsymbol{\theta}_2), \\ \nabla_i f(\boldsymbol{\theta}) &:= \nabla_{\boldsymbol{\theta}_i} f(\boldsymbol{\theta}), \quad i = 1, 2, \\ \Sigma_i &:= \Sigma_{s_i}^{p_i}, \quad i = 1, 2, \\ \Sigma &:= \{\boldsymbol{\theta} \in \mathbb{R}^{p_1+p_2} : \boldsymbol{\theta}_1 \in \Sigma_1, \boldsymbol{\theta}_2 \in \Sigma_2\}, \\ s &:= s_1 + s_2. \end{aligned} \quad (3.1)$$

Similar rules are also applied for $\boldsymbol{\theta}_1^*$ and $\boldsymbol{\theta}_2^*$. Based on these notation, we now establish the first-order necessary and sufficient optimality conditions for problem (1.1).

Theorem 3.1. *Let $\boldsymbol{\theta}^*$ be a point that satisfies*

$$\begin{aligned} (\nabla_j f(\boldsymbol{\theta}^*))_i &= 0, \quad i \in \Gamma(\boldsymbol{\theta}_j^*), \quad \text{if } \|\boldsymbol{\theta}_j^*\|_0 = s_j, \\ \nabla_j f(\boldsymbol{\theta}^*) &= 0, \quad \text{if } \|\boldsymbol{\theta}_j^*\|_0 < s_j, \end{aligned} \quad (3.2)$$

for $j = 1, 2$. Then $\boldsymbol{\theta}^*$ is a local minimizer of (1.1) if and only if it satisfies (3.2).

Proof. Necessity. Based on [24, Theorem 6.12], a local minimizer $\boldsymbol{\theta}^*$ of the problem (1.1) must satisfy that $-\nabla f(\boldsymbol{\theta}^*) \in \mathcal{N}_\Sigma(\boldsymbol{\theta}^*) = \mathcal{N}_{\Sigma_1}(\boldsymbol{\theta}_1^*) \times \mathcal{N}_{\Sigma_2}(\boldsymbol{\theta}_2^*)$, where $\mathcal{N}_\Sigma(\boldsymbol{\theta}^*)$ is the normal cone of Σ at $\boldsymbol{\theta}^*$ and the equality is by [24, Theorem 6.41]. Then the explicit expression (see [20, Table 1]) of normal cone $\mathcal{N}_{\Sigma_j}(\boldsymbol{\theta}_j^*)$ enable us to derive (3.2) immediately.

Sufficiency. Let $\boldsymbol{\theta}^*$ satisfy (3.2). The convexity of f leads to

$$f(\boldsymbol{\theta}) \geq f(\boldsymbol{\theta}^*) + \langle \nabla_1 f(\boldsymbol{\theta}^*), \boldsymbol{\theta}_1 - \boldsymbol{\theta}_1^* \rangle + \langle \nabla_2 f(\boldsymbol{\theta}^*), \boldsymbol{\theta}_2 - \boldsymbol{\theta}_2^* \rangle.$$

If there is a $\delta > 0$ such that, for any $\boldsymbol{\theta} \in \Sigma \cap N(\boldsymbol{\theta}^*, \delta)$,

$$\langle \nabla_1 f(\boldsymbol{\theta}^*), \boldsymbol{\theta}_1 - \boldsymbol{\theta}_1^* \rangle = \langle \nabla_2 f(\boldsymbol{\theta}^*), \boldsymbol{\theta}_2 - \boldsymbol{\theta}_2^* \rangle = 0, \quad (3.3)$$

then the conclusion can be made immediately. Therefore, we next to show (3.3). In fact, by (3.2), we note that $\nabla_j f(\boldsymbol{\theta}^*) = 0$ if $\|\boldsymbol{\theta}_j^*\|_0 < s_j$, which indicates it suffices to consider the worst case of $\|\boldsymbol{\theta}_j^*\|_0 = s_j, j = 1, 2$. Under such a case, we define

$$\delta := \min_{j=1,2} \min_{i \in \Gamma(\boldsymbol{\theta}_j^*)} |(\boldsymbol{\theta}_j^*)_i|.$$

Then for any $\boldsymbol{\theta} \in N(\boldsymbol{\theta}^*, \delta) \cap \Sigma$, we have

$$\begin{aligned} \forall i \in \Gamma(\boldsymbol{\theta}_j^*), \quad |(\boldsymbol{\theta}_j)_i| &= |(\boldsymbol{\theta}_j)_i^* - (\boldsymbol{\theta}_j)_i^* + (\boldsymbol{\theta}_j)_i| \\ &\geq |(\boldsymbol{\theta}_j)_i^*| - |(\boldsymbol{\theta}_j)_i^* - (\boldsymbol{\theta}_j)_i| \\ &\geq |(\boldsymbol{\theta}_j)_i^*| - \|\boldsymbol{\theta}_j^* - \boldsymbol{\theta}_j\| \\ &> |(\boldsymbol{\theta}_j)_i^*| - \delta \\ &\geq 0. \end{aligned}$$

The above relationship means that $i \in \Gamma(\boldsymbol{\theta}_j^*)$ (i.e. $(\boldsymbol{\theta}_j^*)_i \neq 0$) implies $|(\boldsymbol{\theta}_j)_i| > 0$ (i.e. $(\boldsymbol{\theta}_j)_i \neq 0$), which leads to $\Gamma(\boldsymbol{\theta}_j^*) \subseteq \Gamma(\boldsymbol{\theta}_j)$. This by $\|\boldsymbol{\theta}_j\|_0 \leq s_j = \|\boldsymbol{\theta}_j^*\|_0 = |\Gamma(\boldsymbol{\theta}_j^*)|$ allows us to yield that

$$\Gamma(\boldsymbol{\theta}_j^*) = \Gamma(\boldsymbol{\theta}_j), j = 1, 2, \forall \boldsymbol{\theta} \in N(\boldsymbol{\theta}^*, \delta) \cap \Sigma.$$

Using the above fact and (3.2) derives that

$$\langle \nabla_j f(\boldsymbol{\theta}^*), \boldsymbol{\theta}_j - \boldsymbol{\theta}_j^* \rangle = \langle (\nabla_j f(\boldsymbol{\theta}^*))_{\Gamma(\boldsymbol{\theta}_j^*)}, (\boldsymbol{\theta}_j - \boldsymbol{\theta}_j^*)_{\Gamma(\boldsymbol{\theta}_j^*)} \rangle = 0.$$

The proof is complete. \square

\square

Based on Theorem 3.1, however, the necessary and sufficient optimality conditions (3.2) mean that there is no useful information for the case $i \notin \Gamma(\boldsymbol{\theta}_j^*)$ when $\|\boldsymbol{\theta}_j^*\|_0 = s_j$. So, we introduce the concept of the α -stationary point of (1.1).

Definition 3.1. We say that $\boldsymbol{\theta}^*$ is an α -stationary point of problem (1.1) if there exists an $\alpha > 0$ such that

$$\boldsymbol{\theta}_1^* \in \Pi_{\Sigma_1}(\boldsymbol{\theta}_1^* - \alpha \nabla_1 f(\boldsymbol{\theta}^*)), \quad \boldsymbol{\theta}_2^* \in \Pi_{\Sigma_2}(\boldsymbol{\theta}_2^* - \alpha \nabla_2 f(\boldsymbol{\theta}^*)).$$

If there is only one variable, the definition of the α -stationary points is the same as that in [4, 20] which allows us to derive its explicit expression as follows.

Lemma 3.1. For a given $\alpha > 0$, the point $\boldsymbol{\theta}^*$ is an α -stationary point of problem (1.1) if and only if for $j = 1, 2$, it satisfies

$$\alpha(\nabla_j f(\boldsymbol{\theta}^*))_i \begin{cases} = 0, & i \in \Gamma(\boldsymbol{\theta}_j^*), \\ \leq (\boldsymbol{\theta}^*)_s^\downarrow, & i \in \bar{\Gamma}(\boldsymbol{\theta}_j^*), \end{cases} \quad \text{if } \|\boldsymbol{\theta}_j^*\|_0 = s_j, \quad (3.4)$$

$$\nabla_j f(\boldsymbol{\theta}^*) = 0, \quad \text{if } \|\boldsymbol{\theta}_j^*\|_0 < s_j.$$

Comparing conditions (3.2) and (3.4), the latter provides more information for the case of $i \in \bar{\Gamma}(\boldsymbol{\theta}_j^*)$. It can be clearly seen that the latter is a stronger condition and suffices to the former.

The following result reveals the relationships among the α -stationary point and the global/local minimizers of problem (1.1).

Theorem 3.2. Let $\boldsymbol{\theta}^*$ be an α -stationary point of problem (1.1), then it is a local minimizer. Furthermore, if $\|\boldsymbol{\theta}_1^*\|_0 < s_1, \|\boldsymbol{\theta}_2^*\|_0 < s_2$, then it is also a global minimizer. Conversely, if $\boldsymbol{\theta}^*$ is a global minimizer of problem (1.1), then it is an α -stationary point with $0 < \alpha < 1/L_f$.

Proof. Since condition (3.4) imply (3.2) and a point satisfying (3.2) is a local minimizer by Theorem 3.1, an α -stationary point of (1.1) is a local minimizer.

Conversely, suppose that a global minimizer $\boldsymbol{\theta}^*$ of problem (1.1) is not an α -stationary point with $0 < \alpha < 1/L_f$, that is, there exists $\boldsymbol{\eta}_1^* \neq \boldsymbol{\theta}_1^*$ or $\boldsymbol{\eta}_2^* \neq \boldsymbol{\theta}_2^*$ such that

$$\boldsymbol{\eta}_1^* \in \Pi_{\Sigma_1}(\boldsymbol{\theta}_1^* - \alpha \nabla_1 f(\boldsymbol{\theta}^*)) \quad \text{or} \quad \boldsymbol{\eta}_2^* \in \Pi_{\Sigma_2}(\boldsymbol{\theta}_2^* - \alpha \nabla_2 f(\boldsymbol{\theta}^*)).$$

Without loss of any generality, we have both of the above conditions. Then

$$\|\boldsymbol{\eta}_j^* - \boldsymbol{\theta}_j^* + \alpha \nabla_j f(\boldsymbol{\theta}^*)\|^2 \leq \|\boldsymbol{\theta}_j^* - \boldsymbol{\theta}_j^* + \alpha \nabla_j f(\boldsymbol{\theta}^*)\|^2, \quad j = 1, 2,$$

by the definition of projection $\Pi(\cdot)$, which implies

$$\langle \boldsymbol{\eta}_j^* - \boldsymbol{\theta}_j^*, \nabla_j f(\boldsymbol{\theta}^*) \rangle \leq -(1/2\alpha) \|\boldsymbol{\eta}_j^* - \boldsymbol{\theta}_j^*\|^2, \quad i = 1, 2.$$

Using this condition and the strong smoothness of f results in

$$\begin{aligned} f(\boldsymbol{\eta}^*) &\leq f(\boldsymbol{\theta}^*) + \langle \nabla f(\boldsymbol{\theta}^*), \boldsymbol{\eta}^* - \boldsymbol{\theta}^* \rangle + (L_f/2) \|\boldsymbol{\eta}^* - \boldsymbol{\theta}^*\|^2 \\ &\leq f(\boldsymbol{\theta}^*) + (L_f/2 - 1/(2\alpha)) \|\boldsymbol{\eta}^* - \boldsymbol{\theta}^*\|^2 < f(\boldsymbol{\theta}^*), \end{aligned}$$

where the last inequality is from $0 < \alpha < 1/L_f$. The above condition contradicts with the optimality of $\boldsymbol{\theta}^*$. So $\boldsymbol{\theta}^*$ is an α -stationary point with $0 < \alpha < 1/L_f$. The proof is complete. \square \square

To end this section, we would like to see the existence and uniqueness of solutions to problem (1.1), which is revealed by the following theorem.

Theorem 3.3. *If matrix $[X \ Z]$ is s -regular, then the global minimizer of problem (1.1) exists, and the local minimizers are finitely many and each of them is unique.*

Proof. Based on our notation $\mathbb{R}_{T_j}^{p_j} = \{\boldsymbol{\theta}_j \in \mathbb{R}^{p_j} : \Gamma(\boldsymbol{\theta}_j) \subseteq T_j\}$, we note that $\boldsymbol{\theta}_j \in \mathbb{R}_{T_j}^{p_j}$ implies $|\Gamma(\boldsymbol{\theta}_j)| \leq |T_j|, j = 1, 2$. Therefore, original problem (1.1) is equivalent to

$$\begin{aligned} & \min_{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2} f(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \\ & \text{s.t. } \boldsymbol{\theta}_1 \in \mathbb{R}_{T_1}^{p_1}, \quad \forall |T_1| = s_1, \\ & \quad \boldsymbol{\theta}_2 \in \mathbb{R}_{T_2}^{p_2}, \quad \forall |T_2| = s_2. \end{aligned}$$

This problem is clearly equivalent to

$$\min_{|T_1|=s_1, |T_2|=s_2} \left\{ \min_{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2} f(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2), \quad \text{s.t. } \boldsymbol{\theta}_1 \in \mathbb{R}_{T_1}^{p_1}, \boldsymbol{\theta}_2 \in \mathbb{R}_{T_2}^{p_2} \right\}. \quad (3.5)$$

If matrix $[X \ Z]$ is s -regular, then f is s -restricted strongly convex on $\Sigma_s^{p_1+p_2} = \Sigma_{s_1+s_2}^{p_1+p_2}$ by Proposition 2.1 or Proposition 2.2, and hence it is s -restricted strongly convex on $\mathbb{R}_{T_1}^{p_1} \times \mathbb{R}_{T_2}^{p_2}$ due to $\mathbb{R}_{T_1}^{p_1} \times \mathbb{R}_{T_2}^{p_2} \subseteq \Sigma_{s_1+s_2}^{p_1+p_2}$.

It follows from [19, Lemma 6] that the inner program admits a unique global minimizer denoted by $(\boldsymbol{\theta}_1^*(T_1), \boldsymbol{\theta}_2^*(T_2))$. Note that $T_1 \subseteq [p_1]$ and $T_2 \subseteq [p_2]$. Thus there are finitely many T_1 and T_2 such that $|T_1| = s_1$ and $|T_2| = s_2$, and so are the inner programs. This indicates that $(\boldsymbol{\theta}_1^*(T_1), \boldsymbol{\theta}_2^*(T_2))$ is finitely many. To derive the global minimizer of (3.5), we only pick one $(\boldsymbol{\theta}_1^*(T_1), \boldsymbol{\theta}_2^*(T_2))$ that makes the objective function value of (3.5) minimal. Therefore, the global minimizers exist.

We next show that any local minimizer $\boldsymbol{\theta}^*$ is unique. To proceed with that, denote $\delta := \min\{\delta_1, \delta_2\}$ where

$$\delta_j := \begin{cases} +\infty, & \boldsymbol{\theta}_j^* = 0, \\ \min_{i \in \Gamma(\boldsymbol{\theta}_j^*)} |(\boldsymbol{\theta}_j^*)_i|, & \boldsymbol{\theta}_j^* \neq 0, \end{cases} \quad j = 1, 2.$$

Clearly, $\delta_1, \delta_2 > 0$ and hence $\delta > 0$. Then, similar reasoning allows us to derive (3.3) for any $\boldsymbol{\theta} \in \Sigma \cap N(\boldsymbol{\theta}^*, \delta)$. This and f being s -restricted strongly convex lead to

$$f(\boldsymbol{\theta}) \geq f(\boldsymbol{\theta}^*) + (l_f/2)\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|^2.$$

The above condition indicates $\boldsymbol{\theta}^*$ is the unique global minimizer of problem $\min\{f(\boldsymbol{\theta}) : \boldsymbol{\theta} \in \Sigma \cap N(\boldsymbol{\theta}^*, \delta)\}$, namely, $\boldsymbol{\theta}^*$ is the unique local minimizer of problem (1.1). The proof is complete. $\square \square$

4 Gradient Projection Newton Algorithm

In this section, we propose the gradient projection Newton algorithm (GPNA) for problem (1.1). Again, for notational simplicity, we define some notations

$$\begin{aligned} \mathbf{u}^k & := (\mathbf{u}_1^k; \mathbf{u}_2^k), & \boldsymbol{\theta}^k & := (\boldsymbol{\theta}_1^k; \boldsymbol{\theta}_2^k), \\ \Gamma_k & := \Gamma(\mathbf{u}^k), & H^k & := \nabla^2 f(\mathbf{u}^k), \\ \boldsymbol{\theta}^k(\alpha) & := (\boldsymbol{\theta}_1^k(\alpha); \boldsymbol{\theta}_2^k(\alpha)), & \boldsymbol{\theta}_j^k(\alpha) & \in \Pi_{\Sigma_j}(\boldsymbol{\theta}_j^k - \alpha \nabla_j f(\boldsymbol{\theta}^k)), \quad j = 1, 2. \end{aligned}$$

Based on the notation in (3.1), we actually have

$$\boldsymbol{\theta}^k(\alpha) \in \Pi_{\Sigma}(\boldsymbol{\theta}^k - \alpha \nabla f(\boldsymbol{\theta}^k)). \quad (4.1)$$

The algorithmic framework of GPNA summarized in Algorithm 1 consists of two major components. The first one is based on the two projected gradient steps, which enforces two variables to satisfy the sparsity constraints. The second part adopts a Newton step to speed up the convergence. However, the Newton step is only performed when one of the following conditions is satisfied,

$$\begin{aligned} \text{Condition 1) } & \Gamma(\boldsymbol{\theta}_1^k) = \Gamma(\mathbf{u}_1^k), \quad \Gamma(\boldsymbol{\theta}_2^k) = \Gamma(\mathbf{u}_2^k), \\ \text{Condition 2) } & \|\nabla_1 f(\mathbf{u}^k)\| < \epsilon, \quad \Gamma(\boldsymbol{\theta}_2^k) = \Gamma(\mathbf{u}_2^k), \\ \text{Condition 3) } & \|\nabla_2 f(\mathbf{u}^k)\| < \epsilon, \quad \Gamma(\boldsymbol{\theta}_1^k) = \Gamma(\mathbf{u}_1^k), \\ \text{Condition 4) } & \|\nabla_1 f(\mathbf{u}^k)\| < \epsilon, \quad \|\nabla_2 f(\mathbf{u}^k)\| < \epsilon, \end{aligned} \quad (4.2)$$

where $\epsilon > 0$ is a given tolerance.

Algorithm 1 GPNA: Gradient Projection Newton Algorithm

Require: Initialize $\boldsymbol{\theta}^0$. Let $0 < \sigma, 0 < \epsilon, 0 < \alpha_0 \leq 1, 0 < \gamma < 1, 0 < \varepsilon < \text{tol}_0$ and set $k \Leftarrow 0$.

1: **while** $\text{tol}_k > \varepsilon$ **do**

2: Gradient projection: Find the smallest integer $q_k = 0, 1, \dots$ such that

$$f(\boldsymbol{\theta}^k(\alpha_0 \gamma^{q_k})) \leq f(\boldsymbol{\theta}^k) - (\sigma/2) \|\boldsymbol{\theta}^k(\alpha_0 \gamma^{q_k}) - \boldsymbol{\theta}^k\|^2.$$

4: Set $\alpha_k = \alpha_0 \gamma^{q_k}$, $\mathbf{u}^k = \boldsymbol{\theta}^k(\alpha_k)$ and $\boldsymbol{\theta}^{k+1} = \mathbf{u}^k$.

5: **if** one of the conditions in (4.2) is satisfied **then**

6: Newton step: If the following equations are solvable

$$H_{\Gamma_k \Gamma_k}^k (\mathbf{v}_{\Gamma_k}^k - \mathbf{u}_{\Gamma_k}^k) = -(\nabla f(\mathbf{u}^k))_{\Gamma_k}, \quad \mathbf{v}_{\bar{\Gamma}_k}^k = 0, \quad (4.3)$$

8: and the solution \mathbf{v}^k satisfies

$$f(\mathbf{v}^k) \leq f(\mathbf{u}^k) - (\sigma/2) \|\mathbf{v}^k - \mathbf{u}^k\|^2, \quad (4.4)$$

10: then set $\boldsymbol{\theta}^{k+1} = \mathbf{v}^k$.

11: **end if**

12: Compute $\text{tol}_k := \|(\nabla f(\boldsymbol{\theta}^{k+1}))_{\Gamma_k}\|$ and set $k := k + 1$.

13: **end while**

14: Output the solution $\boldsymbol{\theta}^k$.

Remark 4.1. We have some comments on the halting condition and computational complexity for GPNA in Algorithm 1.

- One can discern that if $\boldsymbol{\theta}^{k+1} = \mathbf{u}^k$, then $\boldsymbol{\theta}_{\Gamma_k}^{k+1} = 0$. If $\boldsymbol{\theta}^{k+1} = \mathbf{v}^k$, then the updating rule (4.3) for \mathbf{v}^k indicates that

$$\Gamma(\mathbf{v}^k) \subseteq \Gamma_k = \Gamma(\mathbf{u}^k), \quad (4.5)$$

which also implies $\boldsymbol{\theta}_{\Gamma_k}^{k+1} = 0$. Now suppose $\text{tol}_k = 0$, i.e., $(\nabla f(\boldsymbol{\theta}^{k+1}))_{\Gamma_k} = 0$. Then $\boldsymbol{\theta}^{k+1}$ satisfies (3.2) and thus is a local minimizer of problem (1.1). Therefore, it makes sense to terminate the algorithm when $\text{tol}_k < \varepsilon$.

- We note that the calculations of $\Pi_{\Sigma_1}, \Pi_{\Sigma_2}$ and gradient ∇f dominate the computation for the gradient projection step. And these three terms are easy to calculate and their total computational complexity is about $O(n(p_1 + p_2))$. For the Newton step, if matrix $[X \ Z]$ is s -regular, then the inverse of $H_{\Gamma_k \Gamma_k}^k$ exists due to $|\Gamma_k| \leq s$, which means that every Newton step is well defined. Moreover, the worst-case computational complexity of deriving \mathbf{v}^k is about $O(s^3 + ns^2)$. Overall, the entire computational complexity of the k th iteration of Algorithm 1 is $O(s^3 + ns^2 + q_k n(p_1 + p_2))$. We prove that α_k is bounded by upper and lower bounds. If we know the strong smooth parameter L_f of the objective function f , then q_k may be taken as 1 or a small positive integer.

4.1 Global convergence

Before establishing the main convergence results, we define a constant $\underline{\alpha}$ by

$$\underline{\alpha} := \min \{1, \gamma(\sigma + L_f)^{-1}\},$$

which is a positive scalar. We first need the following lemma.

Lemma 4.1. *Let $\{\boldsymbol{\theta}^k\}$ be the sequence generated by GPNA. The following statements are true.*

- 1) For any $0 < \alpha \leq 1/(\sigma + L_f)$, it holds that

$$f(\boldsymbol{\theta}^k(\alpha)) \leq f(\boldsymbol{\theta}^k) - (\sigma/2)\|\boldsymbol{\theta}^k(\alpha) - \boldsymbol{\theta}^k\|^2, \quad (4.6)$$

and thus $\inf_{k \geq 0} \{\alpha_k\} \geq \underline{\alpha} > 0$.

- 2) $\{f(\boldsymbol{\theta}^k)\}$ is a non-increasing sequence and

$$\lim_{k \rightarrow \infty} \|\mathbf{u}^k - \boldsymbol{\theta}^k\| = \lim_{k \rightarrow \infty} \|\boldsymbol{\theta}^{k+1} - \boldsymbol{\theta}^k\| = 0.$$

- 3) Any accumulating point of sequence $\{\boldsymbol{\theta}^k\}$ is an α -stationary point with $0 < \alpha \leq \underline{\alpha}$ of problem (1.1).

Proof. 1) It follows from (4.1) that $\boldsymbol{\theta}^k(\alpha) \in \Pi_{\Sigma}(\boldsymbol{\theta}^k - \alpha \nabla f(\boldsymbol{\theta}^k))$ and thus

$$\|\boldsymbol{\theta}^k(\alpha) - (\boldsymbol{\theta}^k - \alpha \nabla f(\boldsymbol{\theta}^k))\|^2 \leq \|\boldsymbol{\theta}^k - (\boldsymbol{\theta}^k - \alpha \nabla f(\boldsymbol{\theta}^k))\|^2,$$

which results in

$$2\alpha \langle \nabla f(\boldsymbol{\theta}^k), \boldsymbol{\theta}^k(\alpha) - \boldsymbol{\theta}^k \rangle \leq -\|\boldsymbol{\theta}^k(\alpha) - \boldsymbol{\theta}^k\|^2.$$

This and the strong smoothness of f with constant L_f derive that

$$\begin{aligned} f(\boldsymbol{\theta}^k(\alpha)) &\leq f(\boldsymbol{\theta}^k) + \langle \nabla f(\boldsymbol{\theta}^k), \boldsymbol{\theta}^k(\alpha) - \boldsymbol{\theta}^k \rangle + (L_f/2)\|\boldsymbol{\theta}^k(\alpha) - \boldsymbol{\theta}^k\|^2 \\ &\leq f(\boldsymbol{\theta}^k) - (1/(2\alpha) - (L_f/2))\|\boldsymbol{\theta}^k(\alpha) - \boldsymbol{\theta}^k\|^2 \\ &\leq f(\boldsymbol{\theta}^k) - (\sigma/2)\|\boldsymbol{\theta}^k(\alpha) - \boldsymbol{\theta}^k\|^2, \end{aligned}$$

where the last inequality is from $0 < \alpha \leq 1/(\sigma + L_f)$. Invoking the Armijo-type step size rule, one has $\alpha_k \geq \gamma/(\sigma + L_f)$, which by $\alpha_k \leq 1$ proves the desired assertion.

2) By (4.6) and $\mathbf{u}^k = \boldsymbol{\theta}^k(\alpha_k)$, we have

$$f(\mathbf{u}^k) \leq f(\boldsymbol{\theta}^k) - (\sigma/2)\|\mathbf{u}^k - \boldsymbol{\theta}^k\|^2. \quad (4.7)$$

By the framework of Algorithm 1, if $\boldsymbol{\theta}^{k+1} = \mathbf{u}^k$, then the above condition implies,

$$f(\boldsymbol{\theta}^{k+1}) \leq f(\boldsymbol{\theta}^k) - (\sigma/2)\|\boldsymbol{\theta}^{k+1} - \boldsymbol{\theta}^k\|^2.$$

If $\boldsymbol{\theta}^{k+1} = \mathbf{v}^k$, then we obtain

$$\begin{aligned} f(\boldsymbol{\theta}^{k+1}) &= f(\mathbf{v}^k) \leq f(\mathbf{u}^k) - (\sigma/2)\|\boldsymbol{\theta}^{k+1} - \mathbf{u}^k\|^2 \\ &\leq f(\boldsymbol{\theta}^k) - (\sigma/2)\|\mathbf{u}^k - \boldsymbol{\theta}^k\|^2 - (\sigma/2)\|\boldsymbol{\theta}^{k+1} - \mathbf{u}^k\|^2 \\ &\leq f(\boldsymbol{\theta}^k) - (\sigma/4)\|\boldsymbol{\theta}^{k+1} - \boldsymbol{\theta}^k\|^2, \end{aligned}$$

where the second and last inequalities used (4.7) and a fact $\|\mathbf{a} + \mathbf{b}\|^2 \leq 2\|\mathbf{a}\|^2 + 2\|\mathbf{b}\|^2$ for all vectors \mathbf{a} and \mathbf{b} . Both cases lead to

$$\begin{aligned} f(\boldsymbol{\theta}^{k+1}) &\leq f(\boldsymbol{\theta}^k) - (\sigma/4)\|\boldsymbol{\theta}^{k+1} - \boldsymbol{\theta}^k\|^2, \\ f(\boldsymbol{\theta}^{k+1}) &\leq f(\boldsymbol{\theta}^k) - (\sigma/2)\|\mathbf{u}^k - \boldsymbol{\theta}^k\|^2. \end{aligned} \quad (4.8)$$

Therefore, $\{f(\boldsymbol{\theta}^k)\}$ is non-increasing, which by (4.8) and $f \geq 0$ yields

$$\begin{aligned} &\sum_{k \geq 0} \max\{(\sigma/4)\|\boldsymbol{\theta}^{k+1} - \boldsymbol{\theta}^k\|^2, (\sigma/2)\|\mathbf{u}^k - \boldsymbol{\theta}^k\|^2\} \\ &\leq \sum_{k \geq 0} [f(\boldsymbol{\theta}^k) - f(\boldsymbol{\theta}^{k+1})] = f(\boldsymbol{\theta}^0) - \lim_{k \rightarrow \infty} f(\boldsymbol{\theta}^{k+1}) \leq f(\boldsymbol{\theta}^0). \end{aligned}$$

The above condition suffices to $\lim_{k \rightarrow \infty} \|\boldsymbol{\theta}^{k+1} - \boldsymbol{\theta}^k\| = \lim_{k \rightarrow \infty} \|\mathbf{u}^k - \boldsymbol{\theta}^k\| = 0$.

3) Let $\boldsymbol{\theta}^*$ be any accumulating point of $\{\boldsymbol{\theta}^k\}$. Then there exists a subset M of $\{0, 1, 2, \dots\}$ such that $\lim_{k(\in M) \rightarrow \infty} \boldsymbol{\theta}^k = \boldsymbol{\theta}^*$. This further implies $\lim_{k(\in M) \rightarrow \infty} \mathbf{u}^k = \boldsymbol{\theta}^*$ by applying 2). In addition, as stated in 1), we have $\{\alpha_k\} \subseteq [\underline{\alpha}, 1]$, which indicates that one can find a subsequence K of M and a scalar $\alpha_* \in [\underline{\alpha}, 1]$ such that $\{\alpha_k : k \in K\} \rightarrow \alpha_*$. Overall, we have

$$\lim_{k(\in K) \rightarrow \infty} \boldsymbol{\theta}^k = \lim_{k(\in K) \rightarrow \infty} \mathbf{u}^k = \boldsymbol{\theta}^*, \quad \lim_{k(\in K) \rightarrow \infty} \alpha_k = \alpha_* \in [\underline{\alpha}, 1]. \quad (4.9)$$

Let $\boldsymbol{\eta}^k := \boldsymbol{\theta}^k - \alpha_k \nabla f(\boldsymbol{\theta}^k)$. The framework of Algorithm 1 implies

$$\mathbf{u}^k \in \Pi_{\Sigma}(\boldsymbol{\eta}^k), \quad \lim_{k(\in K) \rightarrow \infty} \boldsymbol{\eta}^k = \boldsymbol{\theta}^* - \alpha_* \nabla f(\boldsymbol{\theta}^*) =: \boldsymbol{\eta}^*. \quad (4.10)$$

The first condition means $\mathbf{u}^k \in \Sigma$ for any $k \geq 1$. Note that Σ is closed and $\boldsymbol{\theta}^*$ is the accumulating point of $\{\mathbf{u}^k\}$ by (4.9). Therefore, $\boldsymbol{\theta}^* \in \Sigma$, which results in

$$\min_{\boldsymbol{\theta} \in \Sigma} \|\boldsymbol{\theta} - \boldsymbol{\eta}^*\| \leq \|\boldsymbol{\theta}^* - \boldsymbol{\eta}^*\|. \quad (4.11)$$

If ‘<’ holds in the above condition, then there is an $\varepsilon_0 > 0$ such that

$$\begin{aligned} \|\boldsymbol{\theta}^* - \boldsymbol{\eta}^*\| - \varepsilon_0 &= \min_{\boldsymbol{\theta} \in \Sigma} \|\boldsymbol{\theta} - \boldsymbol{\eta}^*\| \\ &\geq \min_{\boldsymbol{\theta} \in \Sigma} (\|\boldsymbol{\theta} - \boldsymbol{\eta}^k\| - \|\boldsymbol{\eta}^k - \boldsymbol{\eta}^*\|) \\ &= \|\mathbf{u}^k - \boldsymbol{\eta}^k\| - \|\boldsymbol{\eta}^k - \boldsymbol{\eta}^*\|, \end{aligned}$$

where the last equality is from (4.10). Taking the limit of both sides of the above condition along $k \in K \rightarrow \infty$ yields $\|\boldsymbol{\theta}^* - \boldsymbol{\eta}^*\| - \varepsilon_0 \geq \|\boldsymbol{\theta}^* - \boldsymbol{\eta}^*\|$ by (4.9) and (4.10), a contradiction with $\varepsilon_0 > 0$. Therefore, we must have the equality holds in (4.11), showing that

$$\boldsymbol{\theta}^* \in \Pi_{\Sigma}(\boldsymbol{\eta}^*) = \Pi_{\Sigma}(\boldsymbol{\theta}^* - \alpha_* \nabla f(\boldsymbol{\theta}^*)).$$

The above relation means the conditions in (3.4) hold for $\alpha = \alpha_*$, then these conditions must hold for any $0 < \alpha \leq \underline{\alpha}$ due to $\underline{\alpha} \leq \alpha_*$ from (4.9), namely,

$$\boldsymbol{\theta}^* \in \Pi_{\Sigma}(\boldsymbol{\theta}^* - \alpha \nabla f(\boldsymbol{\theta}^*)),$$

displaying that $\boldsymbol{\theta}^*$ is an α -stationary point of problem (1.1), as desired. The proof is complete. \square

\square

The above lemma allows us to conclude that the whole sequence converges.

Theorem 4.1. *Let $\{\boldsymbol{\theta}^k\}$ be the sequence generated by GPNA. Then the whole sequence converges to a unique local minimizer of (1.1) if $[X \ Z]$ is s -regular.*

Proof. As shown in Lemma 4.1, $\{\boldsymbol{\theta}^k\} \subseteq \{\boldsymbol{\theta} : f(\boldsymbol{\theta}) \leq f(\boldsymbol{\theta}^0), \boldsymbol{\theta} \in \Sigma\}$ is a bounded set due to s -restricted strong convexity of f from the s -regularity of $[X \ Z]$. Therefore, one can find a subsequence of $\{\boldsymbol{\theta}^k\}$ that converges to α -stationary point $\boldsymbol{\theta}^*$ with $0 < \alpha \leq \underline{\alpha}$ of problem (1.1). Recall that an α -stationary point $\boldsymbol{\theta}^*$ is also a local minimizer by Theorem 3.2, which by Theorem 3.3 indicates that $\boldsymbol{\theta}^*$ is unique if $[X \ Z]$ is s -regular. In other words, $\boldsymbol{\theta}^*$ is an isolated local minimizer of problem (1.1). Finally, it follows from $\boldsymbol{\theta}^*$ being isolated, [18, Lemma 4.10] and $\lim_{k \rightarrow \infty} \|\boldsymbol{\theta}^{k+1} - \boldsymbol{\theta}^k\| = 0$ by Lemma 4.1 that the whole sequence converges to the unique local minimizer, $\boldsymbol{\theta}^*$. The proof is complete. \square

\square

4.2 Convergence rate

This part aims to establish the convergence rate of GPNA when the sequence falls into a local area of its limiting point. Before the main result, we claim the following facts.

Lemma 4.2. *Suppose $[X \ Z]$ is s -regular. Let $\{\boldsymbol{\theta}^k\}$ be the sequence generated by GPNA and $\boldsymbol{\theta}^*$ be its limit. The following results hold for sufficiently large k .*

1) *The support set of $\boldsymbol{\theta}^*$ can be identified by*

$$\Gamma(\boldsymbol{\theta}_j^*) \begin{cases} \subseteq (\Gamma(\boldsymbol{\theta}_j^k) \cap \Gamma(\mathbf{u}_j^k)), & \text{if } \|\boldsymbol{\theta}_j^*\|_0 < s_j, \\ \equiv \Gamma(\boldsymbol{\theta}_j^k) \equiv \Gamma(\mathbf{u}_j^k), & \text{if } \|\boldsymbol{\theta}_j^*\|_0 = s_j, \end{cases} \quad j = 1, 2. \quad (4.12)$$

2) The Newton step is always admitted if we set $\sigma \in (0, l_f/2)$.

Proof. 1) If $\|\boldsymbol{\theta}_j^*\|_0 = s_j$, then by $\boldsymbol{\theta}_j^k \rightarrow \boldsymbol{\theta}_j^*$, $\mathbf{u}_j^k \rightarrow \boldsymbol{\theta}_j^*$ and $\|\boldsymbol{\theta}_j^k\|_0 \leq s_j$, $\|\mathbf{u}_j^k\|_0 \leq s_j$, we must have $\Gamma(\boldsymbol{\theta}_j^*) \equiv \Gamma(\boldsymbol{\theta}_j^k) \equiv \Gamma(\mathbf{u}_j^k)$ for sufficiently large k . If $\|\boldsymbol{\theta}_j^*\|_0 < s_j$, similar reasoning allows for deriving $\Gamma(\boldsymbol{\theta}_j^*) \subseteq \Gamma(\boldsymbol{\theta}_j^k)$ and $\Gamma(\boldsymbol{\theta}_j^*) \subseteq \Gamma(\mathbf{u}_j^k)$.

2) By Theorem 4.1, the limiting point, $\boldsymbol{\theta}^*$, is a local minimizer of problem (1.1). Therefore, it satisfies (3.2) from Theorem 3.1. We first conclude that for sufficiently large k , one of the four conditions in (4.2) must be satisfied. In fact, there are four cases for $\boldsymbol{\theta}^*$ and each case can imply one condition in (4.2) as follows:

$$\text{Case 1) } \|\boldsymbol{\theta}_1^*\|_0 = s_1, \|\boldsymbol{\theta}_2^*\|_0 = s_2 \implies \text{Condition 1),}$$

$$\text{Case 2) } \|\boldsymbol{\theta}_1^*\|_0 < s_1, \|\boldsymbol{\theta}_2^*\|_0 = s_2 \implies \text{Condition 2),}$$

$$\text{Case 3) } \|\boldsymbol{\theta}_1^*\|_0 = s_1, \|\boldsymbol{\theta}_2^*\|_0 < s_2 \implies \text{Condition 3),}$$

$$\text{Case 4) } \|\boldsymbol{\theta}_1^*\|_0 < s_1, \|\boldsymbol{\theta}_2^*\|_0 < s_2 \implies \text{Condition 4).}$$

We now show them one by one. The Lipschitz continuity of ∇f indicates that

$$\begin{aligned} & \max\{\|\nabla_j f(\mathbf{u}^k) - \nabla_j f(\boldsymbol{\theta}^*)\|, \|(\nabla f(\mathbf{u}^k))_{\Gamma_k} - (\nabla f(\boldsymbol{\theta}^*))_{\Gamma_k}\| \} \\ & \leq \|\nabla f(\mathbf{u}^k) - \nabla f(\boldsymbol{\theta}^*)\| \leq L_f \|\mathbf{u}^k - \boldsymbol{\theta}^*\|. \end{aligned} \quad (4.13)$$

The relation of Case 1) \Rightarrow Condition 1) can be derived by (4.12) immediately. For Case 2), we have $\Gamma(\boldsymbol{\theta}_2^k) \equiv \Gamma(\mathbf{u}_2^k)$ by (4.12) and

$$\begin{aligned} \|\nabla_1 f(\mathbf{u}^k)\| &= \|\nabla_1 f(\mathbf{u}^k) - \nabla_1 f(\boldsymbol{\theta}^*)\| \quad (\text{by (3.2)}) \\ &\leq L_f \|\mathbf{u}^k - \boldsymbol{\theta}^*\| \quad (\text{by (4.13)}) \\ &\leq \epsilon. \quad (\text{by } \mathbf{u}^k \rightarrow \boldsymbol{\theta}^*) \end{aligned}$$

Therefore, Case 2) \Rightarrow Condition 2). Similarly, we can show the last two relations.

Next, since $[X \ Z]$ is s -regular, $H_{\Gamma_k \Gamma_k}^k$ is non-singular, which means that the equations (4.3) are solvable. Finally, we show the inequality (4.4) is true when $\sigma \in (0, l_f/2)$. In fact, the conditions (4.12) and (3.2) enable to derive

$$(\nabla f(\boldsymbol{\theta}^*))_{\Gamma_k} = 0, \quad (4.14)$$

for sufficiently large k . Then it follows from (4.3) that

$$\begin{aligned} \|\mathbf{v}^k - \mathbf{u}^k\| &= \|\mathbf{v}_{\Gamma_k}^k - \mathbf{u}_{\Gamma_k}^k\| \quad (\text{by (4.5)}) \\ &= \|(H_{\Gamma_k \Gamma_k}^k)^{-1} (\nabla f(\mathbf{u}^k))_{\Gamma_k}\| \quad (\text{by (4.3)}) \\ &\leq (1/l_f) \|(\nabla f(\mathbf{u}^k))_{\Gamma_k}\| \quad (\text{by (2.2) or (2.3)}) \\ &= (1/l_f) \|(\nabla f(\mathbf{u}^k))_{\Gamma_k} - (\nabla f(\boldsymbol{\theta}^*))_{\Gamma_k}\| \quad (\text{by (4.14)}) \\ &\leq (L_f/l_f) \|\mathbf{u}^k - \boldsymbol{\theta}^*\| \rightarrow 0. \quad (\text{by (4.13)}) \end{aligned}$$

The above condition indicates that $\|\mathbf{v}^k - \mathbf{u}^k\| \rightarrow 0$, resulting in

$$o(\|\mathbf{v}^k - \mathbf{u}^k\|^2) \leq (l_f/4)\|\mathbf{v}^k - \mathbf{u}^k\|^2, \quad (4.15)$$

for sufficiently large k . Now, we have the following chain of inequalities,

$$\begin{aligned} 2f(\mathbf{v}^k) - 2f(\mathbf{u}^k) &= 2\langle \nabla f(\mathbf{u}^k), \mathbf{v}^k - \mathbf{u}^k \rangle + 2o(\|\mathbf{v}^k - \mathbf{u}^k\|^2) \\ &\quad + \langle \nabla^2 f(\mathbf{u}^k)(\mathbf{v}^k - \mathbf{u}^k), \mathbf{v}^k - \mathbf{u}^k \rangle \quad (\text{by Taylor expansion}) \\ &= 2\langle (\nabla f(\mathbf{u}^k))_{\Gamma_k}, (\mathbf{v}^k - \mathbf{u}^k)_{\Gamma_k} \rangle + 2o(\|\mathbf{v}^k - \mathbf{u}^k\|^2) \\ &\quad + \langle H_{\Gamma_k \Gamma_k}^k (\mathbf{v}^k - \mathbf{u}^k)_{\Gamma_k}, (\mathbf{v}^k - \mathbf{u}^k)_{\Gamma_k} \rangle \quad (\text{by (4.5)}) \\ &= -\langle H_{\Gamma_k \Gamma_k}^k (\mathbf{v}^k - \mathbf{u}^k)_{\Gamma_k}, (\mathbf{v}^k - \mathbf{u}^k)_{\Gamma_k} \rangle + 2o(\|\mathbf{v}^k - \mathbf{u}^k\|^2) \quad (\text{by (4.3)}) \\ &\leq -l_f \|(\mathbf{v}^k - \mathbf{u}^k)_{\Gamma_k}\|^2 + o(\|\mathbf{v}^k - \mathbf{u}^k\|^2) \quad (\text{by (2.2) or (2.3)}) \\ &= -l_f \|\mathbf{v}^k - \mathbf{u}^k\|^2 + 2o(\|\mathbf{v}^k - \mathbf{u}^k\|^2) \quad (\text{by (4.5)}) \\ &\leq -(l_f/2)\|\mathbf{v}^k - \mathbf{u}^k\|^2 \quad (\text{by (4.15)}) \\ &\leq -\sigma \|\mathbf{v}^k - \mathbf{u}^k\|^2. \quad (\text{by } \sigma \in (0, l_f/2)) \end{aligned}$$

Overall, the Newton step is always admitted for sufficiently large k . The proof is complete. \square \square

Finally, we conclude that GPNA can converge quadratically for SLCoRe and terminate within finite steps for SCoRe by the following theorem.

Theorem 4.2. *Suppose $[X \ Z]$ is s -regular. Then the sequence generated by GPNA with $\sigma \in (0, l_f/2)$ eventually converges to its limit quadratically for SLCoRe or within finitely many steps for SCoRe, namely, for sufficiently large k ,*

$$\begin{aligned} \|\boldsymbol{\theta}^{k+1} - \boldsymbol{\theta}^*\| &\leq \frac{(1+L_f)^2 C_f}{2l_f} \|\boldsymbol{\theta}^k - \boldsymbol{\theta}^*\|^2, \quad \text{if } \ell = \ell_{\log}, \\ \boldsymbol{\theta}^{k+1} &= \boldsymbol{\theta}^*, \quad \text{if } \ell = \ell_{\text{lin}}. \end{aligned}$$

Proof. We first estimate $\|\mathbf{u}^k - \boldsymbol{\theta}^*\|$. Recalling (4.1) that

$$\mathbf{u}^k = \boldsymbol{\theta}^k(\alpha_k) \in \Pi_{\Sigma}(\boldsymbol{\theta}^k - \alpha_k \nabla f(\boldsymbol{\theta}^k))$$

and $\Gamma_k = \Gamma(\mathbf{u}^k)$, we have

$$\mathbf{u}_{\Gamma_k}^k = \boldsymbol{\theta}_{\Gamma_k}^k - \alpha_k (\nabla f(\boldsymbol{\theta}^k))_{\Gamma_k}, \quad \mathbf{u}_{\bar{\Gamma}_k}^k = 0.$$

This enables us to deliver that

$$\begin{aligned} \|\mathbf{u}^k - \boldsymbol{\theta}^*\| &= \|\boldsymbol{\theta}_{\Gamma_k}^k - \alpha_k (\nabla f(\boldsymbol{\theta}^k))_{\Gamma_k} - \boldsymbol{\theta}_{\Gamma_k}^*\| \quad (\text{by } \mathbf{u}_{\bar{\Gamma}_k}^k = \boldsymbol{\theta}_{\bar{\Gamma}_k}^* = 0 \text{ from (4.12)}) \\ &= \|\boldsymbol{\theta}_{\Gamma_k}^k - \alpha_k (\nabla f(\boldsymbol{\theta}^k))_{\Gamma_k} - \boldsymbol{\theta}_{\Gamma_k}^* - \alpha_k (\nabla f(\boldsymbol{\theta}^*))_{\Gamma_k}\| \quad (\text{by (4.14)}) \\ &\leq \|\boldsymbol{\theta}_{\Gamma_k}^k - \boldsymbol{\theta}_{\Gamma_k}^*\| + \alpha_k \|(\nabla f(\boldsymbol{\theta}^k))_{\Gamma_k} - (\nabla f(\boldsymbol{\theta}^*))_{\Gamma_k}\| \\ &\leq (1 + L_f) \|\boldsymbol{\theta}^k - \boldsymbol{\theta}^*\|. \quad (\text{by } 0 < \alpha_k \leq 1 \text{ and (4.13)}) \end{aligned} \quad (4.16)$$

By Lemma 4.2 2), the Newton step is always admitted for sufficiently large k . Then direct calculations lead the following chain of inequalities,

$$\begin{aligned}
& \|\boldsymbol{\theta}^{k+1} - \boldsymbol{\theta}^*\| = \|\mathbf{v}^k - \boldsymbol{\theta}^*\| \\
& = \|\mathbf{v}_{\Gamma_k}^k - \boldsymbol{\theta}_{\Gamma_k}^*\| \quad (\text{by } \mathbf{v}_{\Gamma_k}^k = \boldsymbol{\theta}_{\Gamma_k}^* = 0 \text{ from (4.12)}) \\
& = \|\mathbf{u}_{\Gamma_k}^k - \boldsymbol{\theta}_{\Gamma_k}^* - (H_{\Gamma_k \Gamma_k}^k)^{-1}(\nabla f(\mathbf{u}^k))_{\Gamma_k}\| \quad (\text{by (4.3)}) \\
& = \|\mathbf{u}_{\Gamma_k}^k - \boldsymbol{\theta}_{\Gamma_k}^* - (H_{\Gamma_k \Gamma_k}^k)^{-1}((\nabla f(\mathbf{u}^k))_{\Gamma_k} - (\nabla f(\boldsymbol{\theta}^*))_{\Gamma_k})\| \quad (\text{by (4.14)}) \\
& \leq (1/l_f) \|H_{\Gamma_k \Gamma_k}^k(\mathbf{u}_{\Gamma_k}^k - \boldsymbol{\theta}_{\Gamma_k}^*) - ((\nabla f(\mathbf{u}^k))_{\Gamma_k} - (\nabla f(\boldsymbol{\theta}^*))_{\Gamma_k})\| \quad (\text{by (2.2) or (2.3)}) \\
& \leq (1/l_f) \|\nabla^2 f(\mathbf{u}^k)(\mathbf{u}^k - \boldsymbol{\theta}^*) - (\nabla f(\mathbf{u}^k) - \nabla f(\boldsymbol{\theta}^*))\| \\
& = (1/l_f) \left\| \int_0^1 (\nabla^2 f(\mathbf{u}^* + t(\mathbf{u}^k - \boldsymbol{\theta}^*)) - \nabla^2 f(\mathbf{u}^k))(\mathbf{u}^k - \boldsymbol{\theta}^*) dt \right\| \\
& \leq (1/l_f) \int_0^1 \|\nabla^2 f(\mathbf{u}^* + t(\mathbf{u}^k - \boldsymbol{\theta}^*)) - \nabla^2 f(\mathbf{u}^k)\| \|\mathbf{u}^k - \boldsymbol{\theta}^*\| dt.
\end{aligned}$$

Note that if $\ell = \ell_{lin}$, then $\nabla^2 f(\mathbf{u}^* + t(\mathbf{u}^k - \boldsymbol{\theta}^*)) = \nabla^2 f(\mathbf{u}^k) = Q$. The above condition implies $\|\boldsymbol{\theta}^{k+1} - \boldsymbol{\theta}^*\| \leq 0$, namely, $\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^*$. If $\ell = \ell_{log}$, then above condition implies

$$\begin{aligned}
\|\boldsymbol{\theta}^{k+1} - \boldsymbol{\theta}^*\| & \leq (1/l_f) \int_0^1 C_f \|\mathbf{u}^* + t(\mathbf{u}^k - \boldsymbol{\theta}^*) - \mathbf{u}^k\| \|\mathbf{u}^k - \boldsymbol{\theta}^*\| dt \quad (\text{by (2.1)}) \\
& \leq (C_f/l_f) \|\mathbf{u}^k - \boldsymbol{\theta}^*\|^2 \int_0^1 (1-t) dt \\
& = (C_f/(2l_f)) \|\mathbf{u}^k - \boldsymbol{\theta}^*\|^2.
\end{aligned}$$

which combining (4.16) can make the conclusion immediately. The proof is complete. \square \square

5 Numerical experiments

This section implements GPNA to solve SCL with synthetic datasets and real datasets. All numerical experiments are conducted by running MATLAB (R2018b) on an ideapad with CPU @2.30GHz 2.40GHz and 4GB memory. Apart from the stopping criterion outlined in the algorithm, we also set the maximum number of iterations to 1000. We set $S = \varepsilon = 0.0001, \epsilon = 0.001, \alpha_0 = 1$ and $\gamma = 0.5$. The initial point is chosen as $\boldsymbol{\theta}^0 = 0$.

5.1 SLCoRe model for discrete response variables

In this subsection, we solve SCL with $\ell = \ell_{log}$, namely, SLCoRe. This model usually works well for the data with discrete response variables. In the sequel, we first present two testing examples, followed by the parameters' tuning for GPNA and its numerical comparisons with some benchmark methods on synthetic and real datasets.

5.1.1 Test examples

Synthetic and real data are tested for SLCoRe.

Example 5.1 (Synthetic data). Similar to [2], each sample $\mathbf{x}_i, i \in [n]$ in $X \in \mathbb{R}^{n \times p_1}$ is independently generated by an autoregressive process

$$x_{i(j+1)} = \theta x_{ij} + \sqrt{1 - \theta^2} c_j \quad \text{for all } j \in [p_1 - 1],$$

with $x_{i1} \in \mathcal{N}(0, 1), c_j \in \mathcal{N}(0, 1)$ and $\theta \in [0, 1)$ being the correlation parameter. Note that the larger θ is, the more correlated two columns are. Let $Z = X + 0.01 \cdot \Lambda$ with $\Lambda_{ij} \in \mathcal{N}(0, 1)$. Therefore, for such an example, $p_1 = p_2 =: p$. The sparse parameters $\boldsymbol{\theta}_1 \in \mathbb{R}^p$ and $\boldsymbol{\theta}_2 \in \mathbb{R}^p$ have s_1 and s_2 nonzero entries that are drawn independently from the standard Gaussian distribution, respectively. Finally, response $\mathbf{y} \in \{0, 1\}^n$ is randomly generated from the Bernoulli distribution with

$$\text{Prob}\{y_i = 0 \mid \mathbf{x}_i, \mathbf{z}_i\} = \frac{1}{2} \left[\frac{1}{1 + \exp(-\langle \mathbf{x}_i, \boldsymbol{\theta}_1 \rangle)} + \frac{1}{1 + \exp(-\langle \mathbf{z}_i, \boldsymbol{\theta}_2 \rangle)} \right].$$

Example 5.2 (Real data). Two real datasets are taken into account. They are the alcohol dependence data with $n = 46, p_1 = 500$ and $p_2 = 300$ [33]¹ and Diffuse large B-cell lymphoma (DLBCL) data with $n = 203, p_1 = 17350$ and $p_2 = 386165$ [14]². All datasets are feature-wisely scaled to $[-1, 1]$.

To evaluate the performance of one method, we report the CPU time (in seconds), the classification error rate (CER) [28] and the canonical correlation value (CCV) defined by

$$\text{CER} := \frac{\|\text{sign}(X\boldsymbol{\theta}_1) - \mathbf{y}\|_0 + \|\text{sign}(Z\boldsymbol{\theta}_2) - \mathbf{y}\|_0}{n}, \quad \text{CCV} := \frac{\|X\boldsymbol{\theta}_1 - Z\boldsymbol{\theta}_2\|}{n},$$

where $\boldsymbol{\theta} = (\boldsymbol{\theta}_1; \boldsymbol{\theta}_2)$ is the solution obtained by one method and $(\text{sign}(\mathbf{x}))_i = 1$ if $x_i > 0$ and $(\text{sign}(\mathbf{x}))_i = 0$ otherwise for $i \in [n]$. Note that the smaller CER (or the smaller CCV or the shorter CPU time) the better performance.

5.1.2 Sensitivity analysis

We now implement GPNA to see its performance under different choices of (a, b, c, s_1, s_2) .

(a) Effect of (a, b, c) . Recall that there are three parameters (a, b, c) involved in problem (1.1). We fix $a = 1, c = 0.01$ but vary $b \in [0.01, 10]$ to see the effect of b and fix $a = 1, b = 1$ but change $c \in [0.0001, 1]$ to see the effect of c . The average results over 100 instances for Example 5.1 are presented in Fig. 1, where $n = 200, p = 2000$ and $s_1 = s_2 = 20$.

When a and c are fixed, from the three above sub-figures in Fig. 1, one can observe that CER is declining steadily when $b \in [0.01, 1)$ but dramatically when $b \in [1, 10]$. However, the best choice of b for CCV and CUP time is $b = 1 = a$. Therefore, for Example 5.1 with $p_1 = p_2$ and $s_1 = s_2$, the best option to set a and b should be $a = b$.

When a and b are fixed, from the three bottom sub-figures in Fig. 1, it can be clearly seen that the larger values of c , the smaller CCV and longer CPU time. One can observe that the variance of $c \in [0.0001, 0.01]$ do not influence CER significantly.

We test some other choices and find the following options for (a, b, c) that allows GPNA to render desirable overall performance:

$$a = \frac{s_1}{s_1 + s_2}, \quad b = \frac{s_2}{s_1 + s_2}, \quad c = \frac{1}{s_1 + s_2}.$$

¹Available at <https://github.com/cran/CVR/blob/master/data/alcohol.rda>

²Available at <http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE11318>

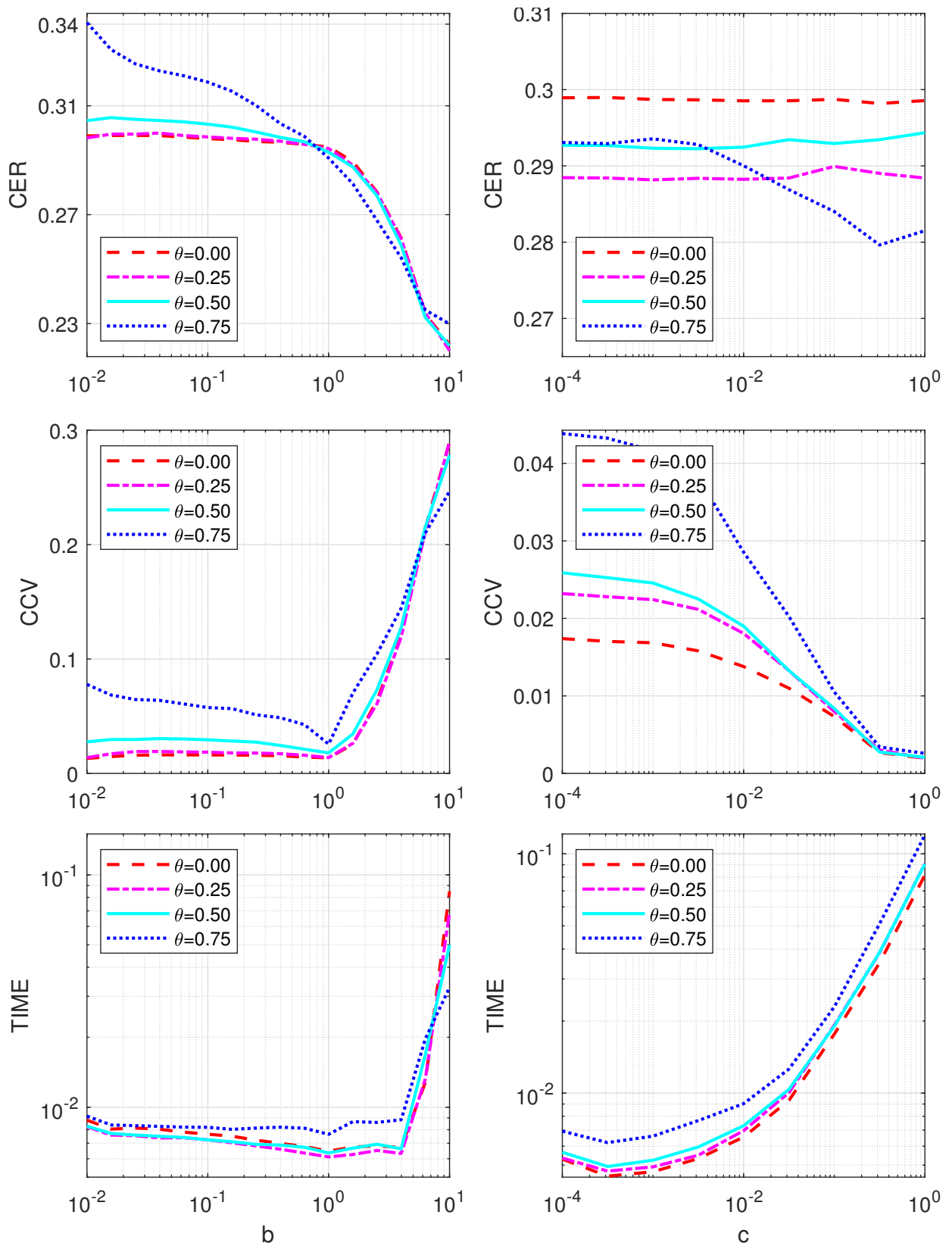


Figure 1: Effect of b and c for Example 5.1.

Therefore, in the following numerical experiments, we fix a, b, c as above choices if no additional information is provided.

(b) **Effect of (s_1, s_2) .** To see the effect of s_1 and s_2 , we choose both s_1 and s_2 from $\{5, 10, \dots, 40\}$. The average results of GPNA for Example 5.1 are shown in Fig. 2 where $n =$

200, $p = 2000$, $\theta = 0.5$. The figure demonstrates that the larger s_1 or s_2 the higher values of CER, leading to better performance. Moreover, the closer between s_1 and s_2 is, the smaller CCV is.

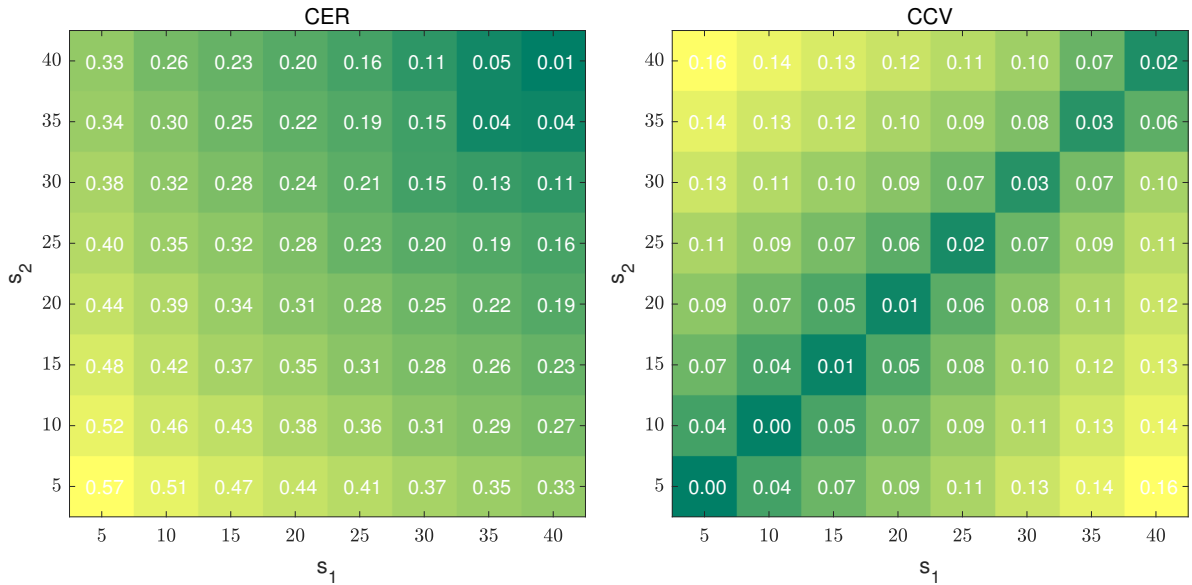


Figure 2: Effect of s_1 and s_2 for Example 5.1.

5.1.3 Effectiveness

To illustrate the effectiveness of our proposed model SCoRe as well as the method GPNA, several alternative approaches are selected. They are SCoRe [11], GPGN [28], GraSP [2], IIHT [21] and NTGP [32]. The first one is used to solve the SCoRe, which can be used to illustrate that SCoRe is a better model than SCoRe for the discrete response variables. GPGN, GraSP, IIHT and NTGP solve the sparse logistic regression that merges two datasets into a single one, which can be used to highlight the advantage of the model SCoRe for two interrelated datasets.

(c) Comparison for Example 5.1. For simplicity, we fix $n = 1000$, $p = 10000$ while choose $\theta \in \{0, 0.5, 0.8\}$ and $s_1, s_2 \in \{200, 300, 500\}$. For each case of (θ, s_1, s_2) , we test 100 instances and report the average results of GPNA, SCoRe, GPGN, IIHT, GraSP and NTGP. Some comments on the reported data in Tables 1 and 2 can be made.

Regarding CER, GPNA achieves the minimum values compared with other methods regardless of the sparsity and correlation how to change. The error rate of the other five methods is more than 10% for the case of two data sets. Moreover, CERs obtained by GPNA, GPGN, IIHT, GraSP and NTGP are smaller than SCoRe, which indicates that ℓ_{log} is more advantageous than ℓ_{lin} for the discrete responses.

Regarding CCV, GPNA delivers tiny values, which shows that there is a high correlation between the two datasets. Although SCoRe can also reveal the relationship between two datasets, the result is not as good as GPNA. Nevertheless, they both perform smaller CCVs than GPGN, IIHT, GraSP and NTGP since the latter four methods solve the model that ignores the relationship between two datasets.

Regarding CPU time, it is obvious that GPNA is the fastest and the calculations take less than a second for all scenarios. By contrast, the other methods need much longer time, especially for larger sparsity, with SCoRe taking 28 seconds, which is 47 times longer than GPNA.

Table 1: Comparison of the results for Example 5.1.

s_1	Algs.	CER			CCV			TIME		
		s_2			s_2			s_2		
		200	300	500	200	300	500	200	300	500
$\theta = 0$										
200	GPNA	0.013	0.016	0.022	0.040	0.074	0.086	00.5	00.6	00.5
	SCoRe	0.410	0.452	0.398	0.231	0.332	0.354	15.5	17.1	18.1
	IIHT	0.396	0.259	0.263	0.286	0.512	0.519	01.5	01.7	01.5
	GraSP	0.224	0.363	0.245	0.384	0.467	0.596	02.9	02.2	01.6
	GPGN	0.117	0.144	0.131	0.486	0.643	0.678	00.7	00.7	00.8
	NTGP	0.128	0.137	0.143	0.573	0.586	0.697	01.1	01.2	01.2
300	GPNA	0.012	0.000	0.000	0.084	0.032	0.062	00.4	00.5	00.5
	SCoRe	0.391	0.384	0.423	0.382	0.518	0.521	20.3	24.2	23.6
	IIHT	0.241	0.256	0.407	0.476	0.561	0.869	01.5	01.7	01.5
	GraSP	0.237	0.266	0.304	0.561	0.627	0.922	03.8	02.1	01.6
	GPGN	0.134	0.107	0.096	0.558	0.734	0.877	00.7	00.7	00.8
	NTGP	0.118	0.142	0.153	0.621	0.727	0.973	01.4	01.5	01.4
500	GPNA	0.024	0.000	0.000	0.087	0.061	0.014	00.5	00.7	00.6
	SCoRe	0.425	0.459	0.480	0.231	0.242	0.318	20.3	24.6	26.4
	IIHT	0.323	0.413	0.328	0.396	0.461	0.469	01.5	01.3	01.5
	GraSP	0.243	0.261	0.252	0.853	0.886	0.877	01.6	01.3	01.1
	GPGN	0.126	0.109	0.115	0.878	0.974	0.963	00.8	00.9	00.9
	NTGP	0.135	0.183	0.167	0.931	0.924	0.987	01.4	01.5	01.6
$\theta = 0.5$										
200	GPNA	0.014	0.021	0.023	0.050	0.086	0.087	00.4	00.4	00.5
	SCoRe	0.423	0.398	0.451	0.213	0.252	0.385	16.7	18.9	21.1
	IIHT	0.264	0.253	0.246	0.319	0.478	0.491	02.1	01.5	02.0
	GraSP	0.243	0.258	0.251	0.343	0.437	0.553	04.7	02.6	01.5
	GPGN	0.134	0.118	0.156	0.461	0.586	0.672	00.6	00.7	00.7
	NTGP	0.144	0.152	0.153	0.429	0.536	0.543	01.5	01.6	01.8
300	GPNA	0.017	0.000	0.000	0.086	0.032	0.043	00.4	00.6	00.7
	SCoRe	0.423	0.384	0.366	0.247	0.342	0.425	21.5	23.7	24.8
	IIHT	0.246	0.273	0.282	0.324	0.363	0.513	01.7	01.9	01.5
	GraSP	0.257	0.253	0.271	0.337	0.472	0.438	03.2	01.9	01.4
	GPGN	0.145	0.157	0.138	0.512	0.466	0.539	00.7	00.8	00.8
	NTGP	0.148	0.157	0.143	0.384	0.473	0.614	01.6	01.7	01.9
500	GPNA	0.018	0.000	0.000	0.089	0.071	0.020	00.5	00.7	00.7
	SCoRe	0.443	0.483	0.456	0.343	0.462	0.437	18.4	22.9	28.4
	IIHT	0.239	0.252	0.399	0.478	0.526	0.854	01.8	01.6	01.7
	GraSP	0.258	0.264	0.285	0.523	0.528	0.694	01.5	01.3	01.1
	GPGN	0.157	0.127	0.148	0.633	0.579	0.715	00.7	00.7	00.8
	NTGP	0.128	0.137	0.153	0.584	0.162	0.849	01.7	01.6	01.7

Table 2: Comparison of the results for Example 5.1.

s_1	Algs.	CER			CCV			TIME		
		s_2			s_2			s_2		
		200	300	500	200	300	500	200	300	500
$\theta = 0.8$										
200	GPNA	0.055	0.051	0.060	0.085	0.104	0.105	00.4	00.5	00.5
	SCoRe	0.491	0.473	0.432	0.252	0.344	0.335	14.7	19.4	21.6
	IIHT	0.282	0.260	0.251	0.274	0.417	0.423	02.0	02.1	01.9
	GraSP	0.301	0.253	0.268	0.284	0.349	0.464	06.4	04.1	03.5
	GPGN	0.167	0.137	0.142	0.433	0.512	0.641	00.6	00.7	00.7
	NTGP	0.211	0.176	0.189	0.343	0.487	0.622	01.7	01.7	01.8
300	GPNA	0.052	0.000	0.000	0.105	0.059	0.093	00.4	00.5	00.5
	SCoRe	0.435	0.412	0.397	0.334	0.338	0.396	16.4	20.3	23.7
	IIHT	0.268	0.271	0.257	0.434	0.475	0.587	01.6	01.7	01.4
	GraSP	0.276	0.284	0.245	0.376	0.433	0.639	03.9	03.2	01.8
	GPGN	0.154	0.138	0.165	0.533	0.537	0.626	00.7	00.7	00.8
	NTGP	0.204	0.225	0.188	0.526	0.491	0.654	01.6	01.7	01.9
500	GPNA	0.061	0.000	0.000	0.108	0.085	0.031	00.5	00.5	00.6
	SCoRe	0.457	0.423	0.382	0.324	0.356	0.431	18.4	22.8	28.7
	IIHT	0.255	0.258	0.266	0.527	0.529	0.912	01.6	01.4	02.0
	GraSP	0.239	0.254	0.263	0.518	0.538	0.883	02.4	01.8	01.4
	GPGN	0.154	0.162	0.166	0.634	0.568	0.942	00.7	00.8	00.8
	NTGP	0.173	0.213	0.188	0.638	0.652	0.875	01.7	01.9	01.8

(d) **Comparison for Example 5.2.** This part reports the numerical comparisons of GPNA, SCoRe, GPGN, IIHT, GraSP and NTGP for analysing two real datasets.

We first apply our method to jointly analyze methylation and gene expression data in an alcohol dependence study [33]. SCoRe can be used to identify the canonical variates from DNA methylation (corresponding to X) and gene expression (corresponding to Z) supervised by the phenotypical information, e.g., alcohol use disorder (AUD), which is observed as a binary indicator variable \mathbf{y} . In this study, genome-wide DNA methylation levels and genome-wide expression levels of genes are quantified for $n = 46$ European Australians. Similar to [17], we choose top $p_1 = 500$ CpG sites and $p_2 = 300$ genes associated with AUD.

We use a random splitting procedure to compare the six methods. At each split, 10 observations are randomly chosen as the testing data and the remaining 36 observations are the training data. The random splitting is repeated 100 times. We choose different sparsity and the average results are reported in Table 3 and show the better behaviour of GPNA since it obtains lower CER (meaning better predictions), smaller CCV and runs much faster.

We next deal with a higher dimensional real dataset DLBCL [14]. It comprises of $n = 203$ patients, each of which has $p_1 = 17350$ gene expression and $p_2 = 386165$ copy numbers. We fixate on the case where \mathbf{y} is a binary variable indicating the survival or death or the cancer subtype.

Table 3: Comparison of the results for Example 5.2.

	s_1	s_2	Training			Testing	
			CER	CCV	TIME(s)	CER	CCV
AUD							
SCoRe			0.617	0.200	001.6	0.582	0.278
GPNA	20	10	0.025	0.004	000.2	0.004	0.005
	20	20	0.020	0.009	000.2	0.002	0.007
	35	20	0.018	0.008	000.2	0.002	0.012
	35	35	0.017	0.007	000.3	0.000	0.005
IIHT	20	10	0.525	0.248	001.6	0.480	0.890
	20	20	0.472	0.251	001.6	0.530	0.893
	35	20	0.455	0.251	001.6	0.463	0.889
	35	35	0.466	0.253	001.7	0.428	0.871
GraSP	20	10	0.528	0.338	001.5	0.410	0.932
	20	20	0.443	0.336	001.3	0.500	0.919
	35	20	0.487	0.328	001.7	0.422	0.922
	35	35	0.482	0.334	001.7	0.338	0.916
GPGN	20	10	0.243	0.365	000.2	0.334	0.974
	20	20	0.284	0.378	000.3	0.347	0.868
	35	20	0.233	0.469	000.3	0.346	0.884
	35	35	0.215	0.478	000.3	0.317	0.967
NTGP	20	10	0.556	0.595	000.3	0.420	0.863
	20	20	0.524	0.553	000.3	0.376	0.761
	35	20	0.488	0.528	000.3	0.397	0.837
	35	35	0.472	0.557	000.3	0.385	0.868
DLBCL							
SCoRe			0.753	0.592	070.4	0.682	0.634
GPNA	50	50	0.054	0.036	000.3	0.024	0.029
	50	100	0.034	0.067	000.3	0.039	0.085
	100	100	0.000	0.024	000.3	0.001	0.022
	100	150	0.000	0.017	000.3	0.002	0.014
IIHT	50	50	0.471	0.796	042.5	0.464	0.732
	50	100	0.458	0.763	043.6	0.483	0.746
	100	100	0.488	0.743	046.3	0.462	0.737
	100	150	0.482	0.737	047.6	0.455	0.739
GraSP	50	50	0.456	0.854	235.4	0.472	0.861
	50	100	0.458	0.846	254.7	0.483	0.867
	100	100	0.432	0.852	228.6	0.457	0.854
	100	150	0.427	0.848	233.2	0.463	0.851
GPGN	50	50	0.408	0.973	000.9	0.487	0.832
	50	100	0.384	0.972	000.9	0.453	0.731
	100	100	0.387	0.881	001.1	0.473	0.848
	100	150	0.395	0.956	001.1	0.434	0.907
NTGP	50	50	0.421	0.834	038.2	0.428	0.911
	50	100	0.452	0.786	040.3	0.478	0.841
	100	100	0.478	0.879	041.8	0.503	0.812
	100	150	0.474	0.934	042.4	0.433	0.865

Again, the 203 samples are split into 153 ones as the training set and 50 ones as the testing set. The random splitting is repeated 100 times. Similar phenomenon to AUD data can be observed for DLBCL in Table 3, showing the better performance of GPNA.

5.2 SCoRe model for continuous response variables

In the subsequent numerical experiments, we focus on SCL with $\ell = \ell_{lin}$, namely, SCoRe. This model is proper for the data with continuous response variables. For such a model, we also do parameters' tuning for GPNA and get similar observations to that for SLCoRe. Therefore, we keep the same setting of parameters as previous examples for GPNA.

5.2.1 Test examples

Again, synthetic and real data are tested for SCoRe.

Example 5.3 (Synthetic data). *The sample data X and Z as well as the sparse parameters $\theta_1 \in \mathbb{R}^p$ and $\theta_2 \in \mathbb{R}^p$ are generated the same as Example 5.1, while the response \mathbf{y} is generated by $\mathbf{y} = (X\theta_1 + Z\theta_2)/2$.*

Example 5.4 (Real data). *Two real datasets are taken into consideration. They are the body mass index (BMI) of mouse data with $n = 294$, $p_1 = 163$ and $p_2 = 215$ [30]³ and DLBCL data. All datasets are feature-wisely scaled to $[-1, 1]$.*

To evaluate the performance of one method, we report the CPU time (in seconds), the mean square error (MSE) and CCV defined by

$$\text{MSE} := \frac{\|\mathbf{y} - X\theta_1\| + \|\mathbf{y} - Z\theta_2\|}{n}, \quad \text{CCV} := \frac{\|X\theta_1 - Z\theta_2\|}{n},$$

where $\theta = (\theta_1; \theta_2)$ is the solution obtained by one method.

5.2.2 Effectiveness

Besides three aforementioned methods SCoRe, GraSP, IIHT, we also select two additional methods SP [8] and LNA [35] for comparisons. Again, GraSP, IIHT, SP and LNA are solving the problem without consider the interrelationship between two datasets.

(e) Comparison for Example 5.3. We first compare GPNA with the other five methods for Example 5.3. For simplicity, we fix $n = 2000, p = 6000$ while choose $\theta \in \{0, 0.5\}$ and $s_1, s_2 \in \{100, 200, 500\}$. For each case of (θ, s_1, s_2) , we test 100 instances and report the average results of GPNA, SCoRe, IIHT, GraSP, SP and LNA. Some comments on the data in Table 4 can be made.

Regarding MSE, GPNA achieves the smallest values in comparison with the other methods regardless of how the sparsity levels s_1, s_2 and correlation parameter θ change. Once again, GPNA produces relatively small CCVs, which indicates that there is a high correlation between the two datasets. By contrast, since IIHT, GraSP, SP and LNA do not take the correlation into account, their generated CCVs are higher than these by GPNA and SCoRe. It can be clearly seen that GPNA runs the fastest, such as 0.6 seconds consumed when $s_1 = s_2 = 500, \theta = 0$ v.s. 23.7, 16.4, 87.6, 56.1 and 2.6 seconds by the other five methods.

³Available at <https://github.com/cran/CVR/blob/master/data/mouse.rda>

Table 4: Comparison of the results for Example 5.3.

s_1	Algs.	MSE			CCV			TIME		
		s_2			s_2			s_2		
		100	200	500	100	200	500	100	200	500
$\theta = 0$										
100	GPNA	0.083	0.097	0.171	0.015	0.081	0.169	00.3	00.4	00.4
	SCoRe	0.243	0.268	0.284	0.031	0.092	0.201	16.5	17.6	19.3
	IIHT	0.157	0.189	0.267	0.155	0.189	0.245	01.3	01.9	07.3
	GraSP	0.232	0.173	0.322	0.163	0.211	0.258	14.4	18.5	33.6
	SP	0.226	0.255	0.364	0.160	0.185	0.263	01.4	01.9	21.1
	LNA	0.167	0.173	0.247	0.174	0.177	0.223	00.7	00.8	01.4
200	GPNA	0.097	0.115	0.161	0.081	0.010	0.139	00.4	00.5	00.5
	SCoRe	0.214	0.277	0.286	0.082	0.093	0.175	19.3	19.2	21.7
	IIHT	0.207	0.227	0.287	0.195	0.223	0.240	01.9	03.2	09.8
	GraSP	0.211	0.252	0.324	0.243	0.245	0.251	21.7	27.2	47.4
	SP	0.243	0.309	0.394	0.174	0.223	0.271	01.8	03.1	24.6
	LNA	0.216	0.236	0.317	0.225	0.244	0.238	00.8	00.9	01.6
500	GPNA	0.165	0.148	0.182	0.163	0.142	0.063	00.4	00.4	00.6
	SCoRe	0.291	0.318	0.285	0.184	0.147	0.185	17.4	21.4	23.7
	IIHT	0.288	0.286	0.367	0.243	0.235	0.282	08.6	11.7	16.4
	GraSP	0.317	0.264	0.334	0.259	0.221	0.273	35.6	50.2	87.6
	SP	0.351	0.411	0.476	0.253	0.288	0.252	16.5	39.8	56.1
	LNA	0.256	0.342	0.329	0.285	0.283	0.264	01.5	01.7	02.6
$\theta = 0.5$										
100	GPNA	0.094	0.096	0.188	0.015	0.082	0.172	00.3	00.6	00.9
	SCoRe	0.242	0.265	0.267	0.144	0.167	0.184	16.7	19.4	22.3
	IIHT	0.172	0.196	0.285	0.163	0.182	0.241	01.8	02.5	11.5
	GraSP	0.187	0.224	0.273	0.152	0.187	0.252	18.3	24.2	32.6
	SP	0.224	0.252	0.377	0.159	0.183	0.269	01.3	01.9	27.4
	LNA	0.213	0.189	0.254	0.171	0.176	0.253	00.6	00.8	01.5
200	GPNA	0.095	0.117	0.168	0.088	0.031	0.137	00.4	00.6	00.8
	SCoRe	0.212	0.224	0.281	0.158	0.145	0.174	18.8	20.5	23.1
	IIHT	0.196	0.233	0.317	0.203	0.218	0.263	02.6	03.9	11.9
	GraSP	0.248	0.236	0.339	0.199	0.253	0.256	21.7	35.3	47.8
	SP	0.289	0.322	0.368	0.214	0.217	0.282	01.8	03.5	29.7
	LNA	0.226	0.265	0.287	0.248	0.255	0.252	00.9	00.9	01.7
500	GPNA	0.187	0.167	0.182	0.183	0.129	0.056	00.6	00.7	00.7
	SCoRe	0.287	0.245	0.272	0.195	0.146	0.122	20.4	21.8	23.3
	IIHT	0.275	0.315	0.346	0.242	0.264	0.144	14.4	21.7	24.3
	GraSP	0.244	0.337	0.386	0.296	0.302	0.221	42.2	51.7	70.4
	SP	0.358	0.402	0.461	0.263	0.278	0.266	21.8	32.3	57.4
	LNA	0.253	0.296	0.306	0.314	0.267	0.278	01.6	01.9	02.9

Table 5: Comparison of the results for Example 5.4.

	s_1	s_2	Training			Testing	
			MSE	CCV	TIME(s)	MSE	CCV
Mouse							
SCoRe			0.423	0.183	001.2	0.386	0.172
GPNA	20	10	0.196	0.121	000.1	0.256	0.151
	20	20	0.174	0.120	000.1	0.223	0.136
	40	20	0.141	0.103	000.1	0.184	0.126
	40	40	0.125	0.088	000.1	0.167	0.103
IIHT	20	10	0.325	0.228	000.5	0.315	0.233
	20	20	0.323	0.197	000.6	0.301	0.198
	40	20	0.319	0.159	000.6	0.305	0.227
	40	40	0.318	0.162	000.7	0.312	0.173
GraSP	20	10	0.324	0.265	000.7	0.336	0.302
	20	20	0.312	0.263	000.7	0.328	0.273
	40	20	0.286	0.237	000.8	0.313	0.262
	40	40	0.294	0.258	000.9	0.327	0.235
SP	20	10	0.337	0.169	000.4	0.344	0.183
	20	20	0.335	0.158	000.4	0.342	0.129
	40	20	0.338	0.146	000.6	0.353	0.187
	40	40	0.334	0.138	000.9	0.355	0.159
LNA	20	10	0.266	0.282	000.3	0.378	0.237
	20	20	0.275	0.269	000.3	0.346	0.245
	40	20	0.254	0.257	000.4	0.361	0.235
	40	40	0.253	0.263	000.5	0.334	0.239
DLBCL							
SCoRe			0.533	0.315	083.4	0.546	0.307
GPNA	50	50	0.267	0.166	000.6	0.313	0.213
	50	100	0.243	0.167	000.6	0.339	0.225
	100	100	0.234	0.159	000.6	0.324	0.212
	100	150	0.233	0.158	000.7	0.311	0.215
IIHT	50	50	0.417	0.352	039.2	0.445	0.326
	50	100	0.412	0.346	042.7	0.437	0.317
	100	100	0.403	0.337	045.6	0.431	0.314
	100	150	0.408	0.346	046.7	0.438	0.324
GraSP	50	50	0.456	0.434	235.5	0.441	0.362
	50	100	0.458	0.457	254.8	0.451	0.353
	100	100	0.432	0.442	228.6	0.439	0.351
	100	150	0.422	0.446	232.7	0.440	0.363
SP	50	50	0.426	0.423	013.5	0.435	0.334
	50	100	0.428	0.437	014.8	0.443	0.341
	100	100	0.416	0.425	015.6	0.432	0.331
	100	150	0.419	0.432	017.4	0.426	0.337
LNA	50	50	0.398	0.451	000.8	0.425	0.366
	50	100	0.414	0.417	000.9	0.407	0.351
	100	100	0.386	0.422	001.1	0.396	0.348
	100	150	0.381	0.436	001.2	0.413	0.359

(f) **Comparison for Example 5.4.** Finally, we report results of five methods for analysing two real datasets: Mouse data and DLBCL. For mouse gene expression data, similar to [17], we choose $p_1 = 163$ single nucleotide polymorphisms (SNPs corresponding to X) and $p_2 = 215$ genes (corresponding to Z) of $n = 294$ for analysis. Again random splitting procedure is employed. At each split, 140 observations are randomly chosen as the testing data and the remaining 154 observations are the training data. The random splitting is repeated 100 times. We choose different sparsity and the average results are reported in Table 5 and display the better behaviour of GPNA since it runs much faster and obtains lower MSE (meaning better predictions), smaller CCV. For DLBCL, results present in Table 5, where the random splitting procedure being same as Example 5.2. Similarly, GPNA obtains lower MSE (meaning better predictions), smaller CCV and runs the fastest, such as 0.6 seconds consumed when $s_1 = s_2 = 50$ v.s. 83.4, 39.2, 235.5, 13.5 and 0.9 seconds by the other five methods, which demonstrate better performance of GPNA.

6 Conclusions and Future work

The SCL model proposed in this paper not only fulfils the tasks of classification or regression for each dataset but also explores the relationship between two datasets. The usage of the double sparsity constraints makes it more efficient for feature selections. To solve the SCL problem, the optimality conditions have been investigated, leading to a gradient projection strategy in the algorithm. To accelerate the convergence, we employed a Newton step when the iteration met some conditions. The final developed gradient projection Newton algorithm has proven to be global and at least quadratic convergent and possessed an excellent numerical performance. We feel that the proposed method is capable of addressing some other general sparsity constrained optimization problems.

As pointed out by our referee, it is an interesting topic to apply the developed techniques and method into dealing with the multi-model problems in particular for some practical applications, such as regional climate prediction. We leave this as future research.

Acknowledgments

The authors would like to thank the Principal Editor and the anonymous referee for their helpful suggestions.

References

- [1] A. Agarwal, S. Negahban, and M. J. Wainwright. Fast global convergence of gradient methods for high-dimensional statistical recovery. *Advances in Neural Information Processing Systems*, 23:2452–2482, 2010.
- [2] S. Bahmani, P. Boufounos, and B. Raj. Greedy sparsity-constrained optimization. *J. Mach. Learn. Res.*, 14(1):807–841, 2013.
- [3] R. Baraniuk, M. Davenport, R. DeVore, and M. Wakin. A simple proof of the restricted isometry property for random matrices. *Constr. Approx.*, 28(3):253–263, 2008.

- [4] A. Beck and Y. C. Eldar. Sparsity constrained nonlinear optimization: optimality conditions and algorithms. *SIAM J. Optimiz.*, 23(3):1480–1509, 2013.
- [5] U. Brefeld, T. Gärtner, T. Scheffer, and S. Wrobel. Efficient co-regularised least squares regression. In *Proceedings of the 23rd international Conference on Machine learning*, pages 137–144, 2006.
- [6] E. J. Candes and T. Tao. Decoding by linear programming. *IEEE T. Inform. Theory*, 51(12):4203–4215, 2005.
- [7] E. J. Candes and T. Tao. Near-optimal signal recovery from random projections: Universal encoding strategies? *IEEE T. Inform. Theory*, 52(12):5406–5425, 2006.
- [8] W. Dai and O. Milenkovic. Subspace pursuit for compressive sensing signal reconstruction. *IEEE T. Inform. Theory*, 55(5):2230–2249, 2009.
- [9] Z. Feng, G. Hu, J. Kittler, W. Christmas, and X. Wu. Cascaded collaborative regression for robust facial landmark detection trained using a mixture of synthetic and real images with dynamic weighting. *IEEE T. Image Process.*, 24(11):3425–3440, 2015.
- [10] S. M. Goldfeld and R. E. Quandt. A Markov model for switching regression. *J. Econometrics*, 1(1):3–15, 1973.
- [11] S. M. Gross and R. Tibshirani. Collaborative regression. *Biostatistics*, 16(2):326–338, 2015.
- [12] W. Hu, B. Cai, A. Zhang, V. D. Calhoun, and Y. Wang. Deep collaborative learning with application to multimodal brain development study. *IEEE T. Bio-Med. Eng.*, 66(12):3346–3359, 2019.
- [13] A. Jalali, C. Johnson, and P. Ravikumar. On learning discrete graphical models using greedy methods. *Advances in Neural Information Processing Systems*, 24:1935–1943, 2011.
- [14] G. Lenz, G. Wright, N. Emre, H. Kohlhammer, S. Dave, R. Davis, S. Carty, L. Lam, A. Shaer, W. Xiao, J. Powell, A. Rosenwald, G. Ott, H. Muller, R. Gascoyne, J. Connors, E. Campo, E. Jae, J. Delabie, E. Smeland, L. Rimsza, R. Fisher, D. Weisenburger, W. Chan, and L. Staudt. Molecular subtypes of diffuse large b-cell lymphoma arise by distinct genetic pathways. *Proceedings of the National Academy of Sciences*, 105(36):13520–13525, 2008.
- [15] F. Liu, X. Huang, C. Gong, J. Yang, and J. Suykens. Indefinite kernel logistic regression with concave-inexact-convex procedure. *IEEE T. Neur. Net. Lear.*, 30(3):1–12, 2018.
- [16] X. Liu, B. Zhao, and W. He. Simultaneous feature selection and classification for data-adaptive kernel-penalized SVM. *Mathematics*, 8(10):1846, 2020.
- [17] C. Luo, J. Liu, D. Dey, and K. Chen. Canonical variate regression. *Biostatistics*, 17(3):468–483, 2017.
- [18] J. J. Moré and D. C. Sorensen. Computing a trust region step. *SIAM J. Sci. Stat. Comput.*, 4(3):553–572, 1983.

- [19] Y. Nesterov. Primal-dual subgradient methods for convex problems. *Math. Program.*, 120(1):221–259, 2009.
- [20] L. Pan, N. Xiu, and S. Zhou. On solutions of sparsity constrained optimization. *J. Oper. Res. Soc. China*, 3(4):421–439, 2017.
- [21] L. Pan, S. Zhou, N. Xiu, and H. Qi. A convergent iterative hard thresholding for sparsity and nonnegativity constrained optimization. *Pac. J. Optim.*, 13(2):325–353, 2017.
- [22] Y. Plan and R. Vershynin. Robust 1-bit compressed sensing and sparse logistic regression: a convex programming approach. *IEEE T. Inform. Theory*, 59(1):482–494, 2013.
- [23] R. Redner and H. Walker. Mixture densities, maximum likelihood and the EM algorithm. *SIAM Rev.*, 26(2):195–239, 1984.
- [24] R. Rockafellar and R. Wets. *Variational analysis*. Springer Science and Business Media, 2009.
- [25] S. Shalev-Shwartz, N. Srebro, and T. Zhang. Trading accuracy for sparsity in optimization problems with sparsity constraints. *SIAM J. Optimiz.*, 20(6):2807–2832, 2010.
- [26] P. M. Thompson, N. G. Martin, and M. J. Wright. Imaging genomics. *Curr. Opin. Neurol.*, 23(4):368–373, 2010.
- [27] P. M. Visscher, M. A. Brown, M. I. McCarthy, and J. Yang. Five years of GWAS discovery. *Am. J. Hum. Genet.*, 90(1):7–24, 2012.
- [28] R. Wang, N. Xiu, and C. Zhang. Greedy projected gradient-Newton method for sparse logistic regression. *IEEE T. Neur. Net. Lear.*, 31(2):527–538, 2020.
- [29] R. Wang, N. Xiu, and S. Zhou. An extended Newton-type algorithm for ℓ_2 -regularized sparse logistic regression and its efficiency for classifying large-scale datasets. *J. Comput. Appl. Math.*, 397:113656, 2021.
- [30] S. Wang, N. Yehya, E. E. Schadt, H. Wang, T. A. Drake, and A. J. Lusis. Genetic and genomic analysis of a fat mass trait with complex inheritance reveals marked sex specificity. *Plos Genet.*, 2(2):148–159, 2006.
- [31] Y. Xiao, J. Wu, Z. Lin, and X. Zhao. A deep learning-based multi-model ensemble method for cancer prediction. *Comput. Meth. Prog. Bio.*, 153:1–9, 2018.
- [32] X. Yuan and Q. Liu. Newton-type greedy selection methods for l_0 -constrained minimization. *IEEE T. Pattern Anal.*, 39(12):2437–2450, 2017.
- [33] H. Zhang, F. Wang, H. Xu, Y. Liu, J. Liu, H. Zhao, and J. Gelernter. Differentially co-expressed genes in postmortem prefrontal cortex of individuals with alcohol use disorders: influence on alcohol metabolism-related pathways. *Hum. Genet.*, 133(2):1383–1394, 2014.
- [34] X. Zhang, Y. Wu, L. Wang, and R. Li. Variable selection for support vector machines in moderately high dimensions. *J. R. Stat. Soc. B*, 78(1):53–76, 2016.

- [35] C. Zhao, N. Xiu, H. Qi, and Z. Luo. A Lagrange–Newton algorithm for sparse nonlinear programming. *Math. Program.*, pages <https://doi.org/10.1007/s10107-021-01719-x>, 2021.
- [36] L. Zhao, K. Oleson, E. Bou-Zeid, S. Krayenhoff, A. Bray, Q. Zhu, Z. Zheng, C. Chen, and M. Oppenheimer. Global multi-model projections of local urban climates. *Nat, Clim. Change*, 11(1):152–157, 2021.
- [37] S. Zhou, N. Xiu, and H. Qi. Global and quadratic convergence of Newton hard-thresholding pursuit. *J. Mach. Learn. Res.*, 22(12):1–45, 2021.
- [38] P. Zille, V. D. Calhoun, and Y. Wang. Enforcing co-expression within a brain-imaging genomics regression framework. *IEEE T. Med. Imaging*, 37(12):2561–2571, 2018.