

# A primal dual semismooth Newton algorithms framework for convex composite optimization

We develop a semismooth Newton based algorithm to solve a class of convex composite optimization problem. Different from many SSN based solvers which are designed case by case or only able to solve two blocks problems. The solver proposed in this paper is able to solve multi-block problems such as conic programming, Lasso type problems and quadratic programming. By dealing internally with the intrinsic nonsmooth property of  $p(\mathbf{x})$ , all the different problems are solved in a unified saddle point framework induced from augment Lagrangian strong duality which lowers the entrance barrier and hence user-friendly. Such structured constraints appear pervasively in image processing and machine learning such as robust PCA clustering problems.

Additional Key Words and Phrases: convex composite optimization, semismooth Newton method, Matlab software package.

## ACM Reference Format:

. 2018. A primal dual semismooth Newton algorithms framework for convex composite optimization. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 20 pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 Introduction

In this paper, we aim to develop a algorithm framework for the following convex composite optimization problem:

$$\begin{aligned} \min_{\mathbf{x}} \quad & \langle \mathbf{c}, \mathbf{x} \rangle + \frac{1}{2} \langle \mathbf{x}, \mathbf{Q}(\mathbf{x}) \rangle + f(\mathcal{B}(\mathbf{x})) + p(\mathbf{x}), \\ \text{s.t.} \quad & \mathbf{x} \in \mathcal{P}_1, \mathcal{A}(\mathbf{x}) \in \mathcal{P}_2. \end{aligned} \tag{1.1} \quad \{\text{general}\}$$

where  $\mathbf{c} \in \mathbb{R}^n$ ,  $\mathcal{A} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $\mathcal{B} : \mathbb{R}^n \rightarrow \mathbb{R}^p$  are linear operators,  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  is a convex function,  $\mathcal{P}_1 = \{\mathbf{x} \in \mathbb{R}^n | \mathbf{l} \leq \mathbf{x} \leq \mathbf{u}\}$  and  $\mathcal{P}_2 = \{\mathbf{x} \in \mathbb{R}^m | \mathbf{l}_b \leq \mathbf{x} \leq \mathbf{u}_b\}$ ,  $\mathbf{Q} \in \mathbb{S}_+^n$  is a positive semidefinite matrix or operator,  $p(\mathbf{x})$  is a convex and nonsmooth function. The choices of  $p(\mathbf{x})$  provide flexibility to handle many kinds of problems.

While the proposed model (1.1) focus on one single variables  $\mathbf{x}$ , our solver is capable of solving the following more general problem with  $N$  block of of variables with shift term:

$$\begin{aligned} \min_{\mathbf{x}} \quad & \sum_{i=1}^N \langle \mathbf{c}_i, \mathbf{x}_i \rangle + \sum_{i=1}^N \frac{1}{2} \langle \mathbf{x}_i, \mathbf{Q}_i(\mathbf{x}_i) \rangle + \sum_{i=1}^N f_i(\mathcal{B}_i(\mathbf{x}) - \mathbf{b}_{2,i}) + \sum_{i=1}^N p_i(\mathbf{x}_i - \mathbf{b}_{1,i}), \\ \text{s.t.} \quad & \mathbf{x}_i \in \mathcal{P}_{1,i}, \sum_{i=1}^N \mathcal{A}_i(\mathbf{x}_i) \in \mathcal{P}_{2,i}, \quad i = 1, \dots, N. \end{aligned} \tag{1.2} \quad \{\text{general-block}\}$$

---

Author's Contact Information:

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

Manuscript submitted to ACM

Manuscript submitted to ACM

In model (1.1),  $p(\mathbf{x})$  usually serves as a regularizer or constraint, which covers a wide range of nonsmooth problems and applications [5, 6]. Compared to the regularized least squares form, the least squares constrained formulation:

$$\min_{\mathbf{x}} \{p(\mathbf{x}) \mid \|\mathcal{A}\mathbf{x} - \mathbf{b}\|_{\infty} \leq \lambda\} \quad (1.3) \quad \{\text{comp2}\}$$

is widely believed to be computationally more challenging because of the complicated geometry of the feasible region. Yet, formulation (1.3) is sometimes preferred in real-data modeling since one can always control the noise level of the model by tuning the acceptable tolerance-the parameter. A potentially feasible approach for solving problem (1.3) is the recently developed level-set method while the level-set method is based on iteratively solving the regularized least squares problem with different  $\mu$ . Therefore, it is of great significance to design efficient algorithms for composite optimization problem (1.1).

However, the existing algorithms can roughly be divided into two categories. The first category consists of algorithms that use only the gradient information. For example, proximal gradient, ADMM, primal-dual, etc. The advantages of first order methods are they are easy to implement, converge fast to a solution with moderate accuracy. However, most of them has slow tail convergence and may fail on slightly more challenging problems. More importantly, the sparsity and low-rank property of the algorithms are not implemented in the iteration process. The other category is the second order methods such as Newton type methods. The advantages of these methods are they usually converge fast to a solution with high accuracy but they are not easy to implement and may encounter numerical difficulties. Moreover, nearly all of these second order information based solvers rely on certain nondegeneracy assumptions to guarantee the nonsingularity of the corresponding inner linear systems.

When  $p(\mathbf{x}) = \|\mathbf{x}\|_1$ , the standard way to solve (1.3) is via linear programming techniques [16] since it is well known that it can be recast as a linear program [17]. Some typical modern solvers rely on interior-point methods (IPM) which are somewhat problematic for large scale problems, since they do not scale well with size. There are some solvers that can handle (1.3) using proximal operator. However, they are designed case by case, it is not user-friendly if the users want to change the model such as adding linear, quadratic term, or shift term  $p_i(\mathbf{x}_i - \mathbf{b}_i)$ .

## 1.1 contribution

The contribution of this paper is listed as follows:

- Generalization: The software packages designed by [2, 8] cannot handle problems such as quadratic programming (QP) and semidefinite programming (SDP) with nonnegative constraints.
- Unification: Although the algorithms proposed in [2, 8] can handle many types of problems, their implementations are designed on a case-by-case basis.
- Second-order algorithms: To the best of our knowledge, we are the first to develop a MATLAB software package incorporating second-order algorithms to solve a wide range of optimization problems.

## 2 Examples

The classical models that (1.1) includes are classified as follows. Consider the following convex optimization problems with constraint:

$$\begin{aligned} \min_{\mathbf{x}} \quad & \|\mathbf{x}\|_1, \\ \text{s.t.} \quad & \|\mathcal{B}\mathbf{x} - \mathbf{b}\|_{\infty} < \lambda_1. \end{aligned} \quad (2.1)$$

For the  $\ell_1$  constrained problem:

$$\begin{aligned} \min_{\mathbf{x}} \quad & \|\mathcal{A}\mathbf{x} - \mathbf{b}\|_2, \\ \text{s.t.} \quad & \|\mathbf{x}\|_1 \leq \lambda_1, \end{aligned} \quad (2.2) \quad \{\text{pro-11-con}\}$$

For constrained Lasso type problems such as (2.2), we note that our algorithm can solve it directly instead of solving a series of subproblems of level set method [6].

We note that we focus on the conic programming with the following structure:

$$\min_{\mathbf{x}} \langle \mathbf{c}, \mathbf{x} \rangle, \quad \text{s.t. } \mathbf{x} \in \mathbb{S}_+^n, \mathbf{b}\mathbf{l} \leq \mathcal{A}\mathbf{x} \leq \mathbf{b}\mathbf{u}, \mathbf{l} \leq \mathbf{x} \leq \mathbf{u}. \quad (2.3) \quad \{\text{conic-general}\}$$

Model (2.3) also includes SDP with affian constraints such as the relaxation SDP problem arising from BIQ and clustering problems (RCP).

When  $p(\mathbf{x}) = \delta_{\mathcal{K}}(\mathbf{x})$  denotes the nonnegative cone, semidefinite cone, or second-order cone, the corresponding problem is linear programming (LP), semidefinite programming (SDP), or second-order cone programming (SOCP) respectively, i.e. the following classical conic programming:

$$\min_{\mathbf{x}} \langle \mathbf{c}, \mathbf{x} \rangle \quad \text{s.t. } \mathcal{A}\mathbf{x} = \mathbf{b}, \mathbf{x} \in \mathbb{S}_+^n. \quad (2.4) \quad \{\text{pro:conic}\}$$

If  $h, p$  are two singleton sets and  $\mathcal{P}_2 = \{\mathbf{b}\}$ , then it corresponds to the classical quadratic programming (QP):

$$\min_{\mathbf{x}} \frac{1}{2} \langle \mathbf{x}, Q(\mathbf{x}) \rangle + \langle \mathbf{c}, \mathbf{x} \rangle \quad \text{s.t. } \mathcal{A}\mathbf{x} = \mathbf{b}, \mathbf{l} \leq \mathbf{x} \leq \mathbf{u}. \quad (2.5) \quad \{\text{QP}\}$$

It is noted in [1] that it is more welcomed to transform QP into second order conic programming. This transformation is numerically more robust than the one for quadratic problems. However, in this paper, we consider solving the dual of (2.5). This approach can make use of the sparsity of the problem. Furthermore, it can be extended to solve the  $\ell_1$ -QP problem directly.

$$\min_{\mathbf{x}} \frac{1}{2} \langle \mathbf{x}, Q\mathbf{x} \rangle + \|\mathbf{x}\|_1 \quad \text{s.t. } \mathcal{A}\mathbf{x} = \mathbf{b}, \mathbf{l} \leq \mathbf{x} \leq \mathbf{u}. \quad (2.6) \quad \{\text{11-QP}\}$$

When  $Q = 0$  and  $\mathbf{c} = 0$ , the corresponding two-block composite optimization is:

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathcal{A}\mathbf{x} - \mathbf{b}\|^2 + p(\mathbf{x}), \quad (2.7) \quad \{\text{comp-lasso}\}$$

Specially, when  $f(\mathbf{x}) = \frac{1}{2} \|\mathcal{A}\mathbf{x} - \mathbf{b}\|^2$  and  $p(\mathbf{x})$  is a nonnegative positively homogeneous convex function such that  $p(0) = 0$ , i.e. a gauge function. Model (2.7) covers typical problems that arise in statistical learning. We list a few of them as follows:

- When  $p(\mathbf{x}) = \lambda \|\mathbf{x}\|_1$ , the problem corresponds to the classical Lasso problem.
- When  $p(\mathbf{x}) = \lambda_1 \|B\mathbf{x}\|_1 + \lambda_2 \|\mathbf{x}\|_1$ , where  $B\mathbf{x} = [\mathbf{x}_2 - \mathbf{x}_1, \dots, \mathbf{x}_n - \mathbf{x}_{n-1}]^\top \in \mathbb{R}^{n-1}$ , then the problem corresponds to the classical Fused Lasso problem.
- When  $p(\mathbf{x}) = \lambda_1 \|\mathbf{x}\|_1 + \lambda_2 \sum_{l=1}^G w_l \|\mathbf{x}_{G_l}\|_2$ , then the problem corresponds to the classical Group Lasso problem.

In addition to Lasso type problems, the robust PCA problem for video segmentation problem can be represented by:

$$\min_{\mathbf{X}} \|\mathbf{X}\|_* + \|\mathbf{D} - \mathbf{X}\|_1, \quad (2.8) \quad \{\text{RPCA}\}$$

where  $\mathbf{X}, \mathbf{D} \in \mathbb{R}^{m \times n}$ , Model (1.1) also covers problems arising in machine learning. The convex clustering model is presented as:

$$\min_{\mathbf{X} \in \mathbb{R}^{d \times n}} \frac{1}{2} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{a}_i\|^2 + \gamma \sum_{(i,j) \in \mathcal{E}} w_{ij} \|\mathbf{x}_i - \mathbf{x}_j\|_q, \quad (2.9) \quad \{\text{cluster}\}$$

where  $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ , denotes the classification feature,  $\gamma > 0$  is a tuning parameter,  $\mathcal{E} = \cup_{i=1}^n \{(i, j) | j \text{ is } i\text{'s } k\text{-nearest neighbors}, i < j \leq n\}$  is the edge set. Typically,  $p$  is chosen to be 1, 2 or  $\infty$ . After solving (2.9) and obtaining the optimal solution  $X^* = [\mathbf{x}_1^*, \dots, \mathbf{x}_n^*]$ , we assign the data vector  $\mathbf{a}_i$  and  $\mathbf{a}_j$  to the same cluster if and only if  $\mathbf{x}_i^* = \mathbf{x}_j^*$ . In other words,  $\mathbf{x}_i^*$  is the centroid for observation  $\mathbf{a}_i$ . Sparse PCA has the following form:

$$\min_{\mathbf{x}} -\langle \mathbf{L}, \mathbf{x} \rangle + \lambda \|\mathbf{x}\|_1, \text{ s.t. } \text{Tr}(\mathbf{x}) = 1, \mathbf{x} \succeq 0. \quad (2.10)$$

The problems examples are summarized in Table 2.

$f$	atom	$f^*$	$\text{prox}_{\lambda p}(\mathbf{x})[\mathbf{u}]$	$\partial \text{prox}_{\lambda p}(\mathbf{x})$ or $\nabla p^*(\mathbf{x})$	Assumptions
$\lambda \ \mathbf{x}\ ^2$	square	$\frac{1}{4\lambda} \ \mathbf{y}\ ^2$	-	$2\lambda \mathbf{x}$	-
$\lambda \sum_{i=1}^n e^{x_i}$	exp	$\sum_{i=1}^n \mathbf{y}_i \log \mathbf{y}_i - \mathbf{y}_i$	-	$\lambda \log(\mathbf{x}/\lambda)$	$(\text{dom}(f^*) = \mathbb{R}_+^n)$
$-\lambda \sum_{i=1}^n \log x_i$	log	$-n - \sum_{i=1}^n \log(-\mathbf{y}_i)$	-	$[\frac{-\lambda}{\log(\mathbf{x})_1}, \dots, \frac{-\lambda}{\log(\mathbf{x})_n}]$	$(\text{dom}(f^*) = \mathbb{R}_+^n)$
$\lambda \log(\sum_{i=1}^n e^{x_i})$	exp	$\sum_{i=1}^n \mathbf{y}_i \log(\mathbf{y}_i)$ , $\text{dom} = \Delta_n$ .	-	$[e^{x_1}, \dots, e^{x_n}] / \sum_{i=1}^n e^{x_i}$	$(\text{dom}(f^*) = \mathbb{R}_+^n)$
$\lambda \ \mathbf{x}\ _1$	l1	$\delta_{B_{\ \cdot\ _\infty}[0, \lambda]}(\mathbf{y})$	$( \mathbf{x}  - \lambda \mathbf{e})_+ \odot \text{sgn}(\mathbf{x})$	$\text{Diag}(\mathbf{u}), \mathbf{u}_i = \begin{cases} 0, & \text{if }  (\mathbf{x})_i  < \lambda, \\ 1, & \text{otherwise.} \end{cases}$	-
$\lambda \ \mathbf{x}\ _2$	l2	$\delta_{B_{\ \cdot\ _2}[0, \lambda]}(\mathbf{y})$	$\begin{cases} \mathbf{x} - \lambda \mathbf{x} / \ \mathbf{x}\ , & \text{if } \ \mathbf{x}\  > \lambda, \\ 0, & \text{otherwise.} \end{cases}$	$\begin{cases} I - \lambda(I - \mathbf{x}\mathbf{x}^T / \ \mathbf{x}\ ^2) / \ \mathbf{x}\ , & \text{if } \ \mathbf{x}\  > \lambda, \\ 0, & \text{otherwise.} \end{cases}$	-
$\lambda \ \mathbf{x}\ _\infty$	linfty	$\delta_{B_{\ \cdot\ _1}[0, \lambda]}(\mathbf{y})$	$\mathbf{x} - \lambda P_{B_{\ \cdot\ _1}[0, 1]}(\mathbf{x}/\lambda)$	$\text{Diag}(\mathbf{u}), \mathbf{u}_i = \begin{cases} 0, & \text{if }  \mathbf{x}/\lambda  > \mu_*, \text{ where } \mu_* \text{ satisfy } \mathbf{1}^T[\mathbf{x} - \mu_* \mathbf{1}]_+ = 1 \\ 1, & \text{otherwise.} \end{cases}$	-
$\delta_{1 \leq \mathbf{x} \leq \mathbf{u}}(\mathbf{x})$	box	$\langle \mathbf{u}, \max\{\mathbf{x}, 0\} \rangle + \langle \mathbf{1}, \min\{\mathbf{x}, 0\} \rangle$	$P_{1 \leq \mathbf{x} \leq \mathbf{u}}(\mathbf{x})$	$\text{Diag}(\mathbf{u}), \mathbf{u}_i = \begin{cases} 1, & \text{if } x_i/\lambda \in C, \\ 0, & \text{otherwise.} \end{cases}$	
$\delta_{\mathbf{x} \in \mathcal{K}}(\mathbf{x})$	cone	$\delta_{\mathbf{y} \in \mathcal{K}}(-\mathbf{y})$	$P_{\mathcal{K}}(\mathbf{x})$	Depend on $\mathcal{K}$	
$\lambda \max\{\mathbf{x}_i\}$	max	$\delta_{\Delta_n}(\mathbf{y})$	$\mathbf{x} - \lambda P_{\Delta_n}(\mathbf{x}/\lambda)$	$\text{Diag}(\mathbf{u}), \mathbf{u}_i = \begin{cases} 0, & \text{if } x_i/\lambda \in C, \\ 1, & \text{otherwise.} \end{cases}$	-
$\lambda \sum_{i=1}^k \mathbf{x}_{[i]}$	topk	-	$\mathbf{x} - \lambda P_{C_{e,k}}(\mathbf{x}/\lambda),$ $C = H_{e,k} \cap \text{Box}[\mathbf{0}, \mathbf{e}]$	$\text{Diag}(\mathbf{u}), \mathbf{u}_i = \begin{cases} 0, & \text{if } x_i/\lambda \in C, \\ 1, & \text{otherwise.} \end{cases}$	-
$\lambda \sum_{k=1}^n  \mathbf{x}_{[i]} $	topkabs	-	$\mathbf{x} - \lambda P_C(\mathbf{x}/\lambda)$ $C = B_{\ \cdot\ _1, [0, k]} \cap \text{Box}[-e, e]$	$\text{Diag}(\mathbf{u}), \mathbf{u}_i = \begin{cases} 0, & \text{if } x_i/\lambda \in C, \\ 1, & \text{otherwise.} \end{cases}$	-
$\lambda H_\mu(\mathbf{x})$	huber	-	-	$\lambda \mathbf{x}$	-
$\lambda_1 \ \mathbf{x}\ _1 + \lambda_2 \ B\mathbf{x}\ _1$	fused	$\delta_{\ \mathbf{y}\ _\infty < \lambda_1}(\mathbf{y}) + \delta_{\ B^T \mathbf{y}\ _\infty < \lambda_2}(\mathbf{y})$	(4.6)	(4.7)	-

$f$	atom	$f^*$	$\text{prox}_{\lambda p}(\mathbf{x})[\mathbf{u}]$	$\partial \text{prox}_{\lambda p}(\mathbf{x})$ or $\nabla p^*(\mathbf{x})$	Assumptions
$\delta_{\ \mathbf{x}\ _1 \leq \lambda}(\mathbf{x})$	l1con	$\lambda \ \mathbf{y}\ _\infty$	$\lambda P_{B_{\ \cdot\ _1}[0,1]}(\mathbf{x}/\lambda)$	$\text{diag}(\mathbf{u}), \mathbf{u}_i = \begin{cases} 0, & \text{if }  \mathbf{x}/\lambda  > \mu_*, \text{ where } \mu_* \\ & \text{satisfy } \mathbf{1}^T[\mathbf{x} - \mu_* \mathbf{1}]_+ = 1 \\ 1, & \text{otherwise.} \end{cases}$	-
$\delta_{\ \mathbf{x}\ _2 \leq \lambda}(\mathbf{x})$	l2com	$\lambda \ \mathbf{y}\ _2$	$\begin{cases} \lambda \mathbf{x} / \ \mathbf{x}\ , & \text{if } \ \mathbf{x}\  > \lambda, \\ \mathbf{x}, & \text{otherwise.} \end{cases}$	$\lambda \log(\mathbf{x}/\lambda)$	$(\text{dom}(f^*) = \mathbb{R}_+^n)$
$\delta_{\ \mathbf{x}\ _\infty \leq \lambda}(\mathbf{x})$	linftycon	$\lambda \ \mathbf{y}\ _1$	$P_{\ \mathbf{x}\ _\infty \leq \lambda}(\mathbf{x})$	$[e^{\mathbf{x}^1}, \dots, e^{\mathbf{x}^n}] / \sum_{i=1}^n e^{\mathbf{x}^i}$	$(\text{dom}(f^*) = \mathbb{R}_+^n)$
$\lambda \ \mathbf{X}\ _1$	l1	$\delta_{B_{\ \cdot\ _\infty}[0,\lambda]}(\mathbf{Y})$	$( \mathbf{X}  - \lambda \mathbf{E})_+ \odot \text{sgn}(\mathbf{X})$	$\mathbf{U}_{i,j} = \begin{cases} 0, & \text{if }  (X)_{i,j}  < \lambda, \\ 1, & \text{otherwise.} \end{cases}$	-
$\lambda \ \mathbf{X}\ _F^2$	frobenius	$\frac{1}{4\lambda} \ \mathbf{Y}\ _F^2$	-	$2\lambda \mathbf{X}$	-
$\lambda \ \mathbf{X}\ _F$	frobenius	$\delta_{B_{\ \cdot\ _F}[0,\lambda]}(\mathbf{Y})$	$\left(1 - \frac{\lambda}{\max\{\ \mathbf{X}\ _F, \lambda\}}\right) \mathbf{X}$	$\begin{cases} I - \lambda(I - \frac{\mathbf{X}\mathbf{X}^T}{\ \mathbf{X}\ _F^2}) / \ \mathbf{X}\ _F, & \text{if }  (X)_i  < \lambda, \\ 0, & \text{otherwise.} \end{cases}$	-
$\lambda \ \mathbf{X}\ _*$	nuclear	$\delta_{B_{\ \cdot\ _{S_\infty}}[0,\lambda]}(\mathbf{Y})$	(4.2)	(4.3)	-
$\lambda \ \mathbf{X}\ _{1,2}$	l1l2	$\delta_{\ \mathbf{Y}\ _{\infty,2} < \lambda}(\mathbf{Y})$	$[\max(0, 1 - \frac{\lambda}{\ \mathbf{X}_1\ _2})X_1, \dots, \max(0, 1 - \frac{\lambda}{\ \mathbf{X}_n\ _2})X_n]$	$\begin{cases} I - \lambda(I - \mathbf{X}_i \mathbf{X}_i^T / \ \mathbf{X}_i\ ^2) / \ \mathbf{X}_i\ , & \text{if } \ \mathbf{X}_i\  > \lambda, \\ 0, & \text{otherwise.} \end{cases}$	-
$\lambda \ \mathbf{X}\ _{1,\infty}$	l1linfty	$\delta_{\ \mathbf{Y}\ _{\infty,1} < \lambda}(\mathbf{Y})$	$\mathbf{X}_i - \lambda P_{B_{\ \cdot\ _1}[0,1]}(\mathbf{X}_i/\lambda)$	$\mathbf{u}_{i,j} = \begin{cases} 0, & \text{if }  \mathbf{X}_j/\lambda  > \mu_{*,j}, \text{ where } \mu_* \\ & \text{satisfy } \mathbf{1}^T[\mathbf{x} - \mu_{*,j} \mathbf{1}]_+ = 1 \\ 1, & \text{otherwise.} \end{cases}$	-
$-\lambda \log \det(\mathbf{X})$	logdet	$-n - \log \det(-\mathbf{Y})$	$\mathbf{U} \text{diag} \left( \frac{\lambda_j(\mathbf{X}) + \sqrt{\lambda_j(\mathbf{X})^2 + 4\lambda}}{2} \right) \mathbf{U}^T$	-	$\mathbb{S}_{++}^n$
$\lambda \sigma_1(\mathbf{X})$	mmax	$\delta_{B_{\ \cdot\ _*}[0,\lambda]}(\mathbf{Y})$	$\mathbf{U} \text{diag} (\lambda(\mathbf{X}) - \lambda P_{\Delta_n}(\lambda(\mathbf{X})/\lambda)) \mathbf{V}^T$	-	$\mathcal{Y}_n$

Table 1. Combined table of functions, their duals, proximal operators, and subdifferentials

Problem Type	Objective Function	Constraints	Function block	Remarks
General Problem	$\underbrace{\langle \mathbf{c}, \mathbf{x} \rangle}_{(1)} + \underbrace{\frac{1}{2} \langle \mathbf{x}, \mathbf{Q}(\mathbf{x}) \rangle}_{(2)} + \underbrace{f(\mathcal{B}(\mathbf{x})) + p(\mathbf{x})}_{(3)}$	$\underbrace{\mathbf{x} \in \mathcal{P}_1}_{(4)}, \underbrace{\mathcal{A}(\mathbf{x}) \in \mathcal{P}_2}_{(5)}$	(1) (2) (3) (4) (5)	Handles various optimization problems
Conic programming	$\langle \mathbf{c}, \mathbf{x} \rangle$	$\mathcal{A}(\mathbf{x}) = \mathbf{b}, \mathbf{x} \geq 0.$	(1)(5)	Linear programming
	$\langle \mathbf{C}, \mathbf{X} \rangle$	$\mathcal{A}(\mathbf{X}) = \mathbf{b}, \mathbf{X} \in \mathbb{S}_+^n.$	(1)(5)	semidefinite programming
	$\langle \mathbf{c}, \mathbf{x} \rangle$	$\mathcal{A}(\mathbf{x}) = \mathbf{b}, \mathbf{x} \in \mathcal{Q}^n.$	(1)(5)	quadratic cone programming
SDP with box constraints	$\langle \mathbf{C}, \mathbf{X} \rangle$	$\mathcal{A}(\mathbf{X}) = \mathbf{b}, \mathbf{x} \in \mathcal{P}_1, \mathbf{X} \in \mathbb{S}_+^n.$	(1)(4)(5)	SDP with box constraints
Lasso type Problems	$\frac{1}{2} \ \mathcal{B}(\mathbf{x}) - \mathbf{b}\ ^2 + \lambda \ \mathbf{x}\ _1$	-	(3)	Lasso problem
	$\frac{1}{2} \ \mathcal{B}(\mathbf{x}) - \mathbf{b}\ ^2 + \lambda_1 \ \mathbf{x}\ _1 + \lambda_2 \ D\mathbf{x}\ _1$	-	(3)	Fused lasso problem
	$\frac{1}{2} \ \mathcal{B}(\mathbf{x}) - \mathbf{b}\ ^2 + \lambda \ \mathbf{x}\ _2$	-	(3)	Group Lasso
	$\frac{1}{2} \ \mathcal{B}(\mathbf{x}) - \mathbf{b}\ ^2 + \lambda \sum_{i=1}^k \mathbf{x}_{[i]}$	-	(3)	Top-k regression
Matrix Completion	$\ \mathcal{B}(\mathbf{X}) - \mathbf{B}\ _{\mathbb{F}}^2 + \lambda \ \mathbf{X}\ _*$	-	(3)	Low-rank matrix recovery
QP	$\langle \mathbf{x}, \mathbf{Q}(\mathbf{x}) \rangle + \langle \mathbf{x}, \mathbf{c} \rangle$	$1 \leq \mathbf{x} \leq \mathbf{u}, \mathcal{A}(\mathbf{x}) = \mathbf{b}.$	(1)(2)(4)(5)	Quadratic Programming
QP with regularizer	$\langle \mathbf{x}, \mathbf{Q}(\mathbf{x}) \rangle + \lambda \ \mathbf{x}\ _1$	$1 \leq \mathbf{x} \leq \mathbf{u}, \mathcal{A}(\mathbf{x}) = \mathbf{b}.$	(1)(2)(3)(5)	QP with $\ell_1$ norm
convex constraint problems	$\ \mathbf{x}\ _1$	$\ \mathcal{B}\mathbf{x} - \mathbf{b}\ _{\infty} < \lambda$	(3)	$\ell_{\infty}$ constraint problem
	$\ \mathcal{B}(\mathbf{x}) - \mathbf{b}\ _1$	$\ \mathbf{x}\ _1 < \lambda$	(3)	$\ell_1$ constraint problem
Statistical learning	$-\langle \mathbf{L}, \mathbf{x} \rangle + \lambda \ \mathbf{x}\ _1,$	$\text{Tr}(\mathbf{x}) = 1, \mathbf{x} \succeq 0$	(3)	Sparse PCA
	$\ \mathbf{x}\ _1$	$\mathcal{A}(\mathbf{x}) = \mathbf{b},$	(5)	Base pursuit
	$\ \mathbf{x}_1\ _* + \mu \ \mathbf{x}_2\ _1$	$\mathbf{x}_1 + \mathbf{x}_2 = \mathbf{D}$	(3)(5)	Roubst PCA
	$-\log(\det(\mathbf{X})) + \text{Tr}(\mathbf{XS}) + \lambda \ \mathbf{X}\ _1$	-	(1)(3)	Sparse covariance matrix estimation
	$\frac{1}{2} \sum_{i=1}^n \ \mathbf{x}_i - \mathbf{a}_i\ ^2 + \gamma \sum_{(i,j) \in \mathcal{E}} w_{ij} \ \mathbf{x}_i - \mathbf{x}_j\ _q$	-	(3)	convex clustering problem

Table 2. The examples of problems Model (1.1) able to solve.

### 3 A primal dual semismooth Newton method for (1.1)

In this section, using the AL duality, we first transform the original problem (1.1) into a saddle point problem. Subsequently, a monotone nonlinear system induced derived from the saddle point problem is presented. Furthermore, such nonlinear system is semismooth and equivalent to the Karush–Kuhn–Tucker (KKT) optimality condition of problem (1.1).

#### 3.1 An equivalent saddle point problem

The dual of model (1.1) can be represented by

$$\begin{aligned} \min_{\mathbf{y}, \mathbf{z}, \mathbf{s}, \mathbf{r}, \mathbf{v}} \quad & \delta_{\mathcal{P}_2}^*(-\mathbf{y}) + f^*(-\mathbf{z}) + p^*(-\mathbf{s}) + \frac{1}{2} \langle \mathbf{Q}\mathbf{v}, \mathbf{v} \rangle + \delta_{\mathcal{P}_1}^*(-\mathbf{r}), \\ \text{s.t.} \quad & \mathcal{A}^* \mathbf{y} + \mathcal{B}^* \mathbf{z} + \mathbf{s} - \mathbf{Q}\mathbf{v} + \mathbf{r} = \mathbf{c}. \end{aligned} \quad (3.1)$$

Introducing the slack variable, the corresponding optimization problem is

$$\begin{aligned} \min_{\mathbf{y}, \mathbf{z}, \mathbf{s}, \mathbf{r}, \mathbf{v}} \quad & \delta_{\mathcal{P}_2}^*(-\mathbf{p}) + f^*(-\mathbf{q}) + p^*(-\mathbf{s}) + \frac{1}{2} \langle \mathbf{Q}\mathbf{v}, \mathbf{v} \rangle + \delta_{\mathcal{P}_1}^*(-\mathbf{t}), \\ \text{s.t.} \quad & \mathcal{A}^* \mathbf{y} + \mathcal{B}^* \mathbf{z} + \mathbf{s} - \mathbf{Q}\mathbf{v} + \mathbf{r} = \mathbf{c}, \quad \mathbf{y} = \mathbf{p}, \quad \mathbf{z} = \mathbf{q}, \quad \mathbf{r} = \mathbf{t}. \end{aligned} \quad (3.2)$$

The corresponding augment Lagrangian function of (3.1) is

$$\begin{aligned} \mathcal{L}_\sigma(\mathbf{y}, \mathbf{s}, \mathbf{z}, \mathbf{r}, \mathbf{v}, \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4) = & \delta_{\mathcal{P}_2}^*(-\mathbf{p}) + f^*(-\mathbf{q}) + p^*(-\mathbf{s}) + \frac{1}{2} \langle \mathbf{Q}(\mathbf{v}), \mathbf{v} \rangle \\ & + \delta_{\mathcal{P}_1}^*(-\mathbf{t}) + \frac{\sigma}{2} \left( \|\mathbf{p} - \mathbf{y} + \mathbf{x}_1/\sigma\|_{\mathbb{F}}^2 + \|\mathbf{q} - \mathbf{z} + \frac{1}{\sigma} \mathbf{x}_2\|_{\mathbb{F}}^2 + \|\mathbf{t} - \mathbf{r} + \frac{1}{\sigma} \mathbf{x}_3\|_{\mathbb{F}}^2 \right) \\ & + \frac{\sigma}{2} (\|\mathcal{A}^* \mathbf{y} + \mathcal{B}^* \mathbf{z} + \mathbf{s} - \mathbf{Q}\mathbf{v} + \mathbf{r} - \mathbf{c} + \frac{1}{\sigma} \mathbf{x}_4\|_{\mathbb{F}}^2) - \frac{1}{2\sigma} \sum_{i=1}^4 \|\mathbf{x}_i\|^2. \end{aligned}$$

Eliminating the variables  $\mathbf{p}, \mathbf{q}, \mathbf{s}, \mathbf{t}$ :

$$\begin{aligned} \mathbf{p} &= \text{prox}_{\delta_{\mathcal{P}_2}^*/\sigma}(\mathbf{x}_1/\sigma - \mathbf{y}), \quad \mathbf{q} = \text{prox}_{f^*/\sigma}(\mathbf{x}_2/\sigma - \mathbf{z}), \\ \mathbf{s} &= \text{prox}_{p^*/\sigma}(\mathcal{A}^* \mathbf{y} + \mathcal{B}^* \mathbf{z} - \mathbf{Q}\mathbf{v} + \mathbf{r} - \mathbf{c} + \frac{1}{\sigma} \mathbf{x}_4), \quad \mathbf{t} = \text{prox}_{\delta_{\mathcal{P}_1}^*/\sigma}(\mathbf{x}_3/\sigma - \mathbf{r}). \end{aligned} \quad (3.3)$$

Let  $\mathbf{w} = (\mathbf{y}, \mathbf{z}, \mathbf{r}, \mathbf{v}, \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4)$  and then we obtain that the modified augment Lagrangian function:

$$\begin{aligned} \Phi_\sigma(\mathbf{w}) = & \delta_{\mathcal{P}_2}^*(\text{prox}_{\delta_{\mathcal{P}_2}^*/\sigma}(\mathbf{x}_1/\sigma - \mathbf{y})) + f^*(\text{prox}_{f^*/\sigma}(\mathbf{x}_2/\sigma - \mathbf{z})) \\ & + \frac{1}{2} \langle \mathbf{Q}\mathbf{v}, \mathbf{v} \rangle + \delta_{\mathcal{P}_1}^*(\text{prox}_{\delta_{\mathcal{P}_1}^*/\sigma}(\mathbf{x}_3/\sigma - \mathbf{r})) \\ & + \frac{1}{2\sigma} \|\Pi_{\mathcal{P}_2}(\mathbf{x}_1 - \sigma\mathbf{y})\|^2 + p^*(\text{prox}_{p^*/\sigma}(\mathbf{x}_4/\sigma - \mathcal{A}^* \mathbf{y} - \mathcal{B}^* \mathbf{z} - \mathbf{Q}\mathbf{v} - \mathbf{r})) \\ & + \frac{1}{2\sigma} \|\text{prox}_{\sigma p}(\mathbf{x}_4 + \sigma(\mathcal{A}^* \mathbf{y} + \mathcal{B}^* \mathbf{z} - \mathbf{Q}\mathbf{v} + \mathbf{r} - \mathbf{c}))\|^2 + \frac{1}{2\sigma} \|\Pi_{\mathcal{P}_1}(\mathbf{x}_3 - \sigma\mathbf{r})\|^2 \\ & + \frac{1}{2\sigma} \|\text{prox}_{\sigma f}(\mathbf{x}_2 - \sigma\mathbf{z})\|^2 - \frac{1}{2\sigma} \sum_{i=1}^4 \|\mathbf{x}_i\|^2. \end{aligned} \quad (3.4)$$

Then the corresponding saddle point problem is presented as

$$\min_{\mathbf{y}, \mathbf{z}, \mathbf{r}, \mathbf{v}} \max_{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4} \Phi(\mathbf{y}, \mathbf{z}, \mathbf{r}, \mathbf{v}; \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4). \quad (3.5)$$



ASSUMPTION 1. The dual problem (1.1) has an optimal solution  $\mathbf{y}_*, \mathbf{z}_*, \mathbf{s}_*$ . Furthermore, Slater's condition holds for the dual problem (3.1), i.e., there exists  $-\mathbf{y} \in \text{ri}(\text{dom}(\delta_Q^*))$ ,  $-\mathbf{s} \in \text{ri}(\text{dom}(\delta_{\mathcal{K}}^*))$ , and  $-\mathbf{z} \in \text{ri}(\text{dom}(h^*))$  such that  $\mathcal{A}^*(\mathbf{y}) + \mathbf{s} + \mathbf{z} = \mathbf{c}$ .

LEMMA 3.1 (STRONG DUALITY). Suppose that Assumption 1 holds. Given  $\sigma > 0$ , the strong duality holds for (??), i.e.

$$\min_{\mathbf{y}, \mathbf{z}} \max_{\mathbf{x}, \mathbf{u}, \mathbf{q}} \Phi(\mathbf{y}, \mathbf{z}, \mathbf{x}, \mathbf{u}, \mathbf{q}) = \max_{\mathbf{x}, \mathbf{u}, \mathbf{q}} \min_{\mathbf{y}, \mathbf{z}} \Phi(\mathbf{y}, \mathbf{z}, \mathbf{x}, \mathbf{u}, \mathbf{q}), \quad (3.6) \quad \{\text{lemma:strong}\}$$

where both sides of (3.6) are equivalent to problem (1.1).

The gradient of the saddle point problem can be represented by

$$\begin{aligned} \nabla_{\mathbf{y}} \Phi_{\sigma}(\mathbf{w}) &= \mathcal{A} \text{prox}_{\sigma p}(\mathbf{x}_4 + \sigma(\mathcal{A}^* \mathbf{y} + \mathcal{B}^* \mathbf{z} - Q\mathbf{v} + \mathbf{r} - \mathbf{c})) - \Pi_{\mathcal{P}_2}(\mathbf{x}_1 - \sigma \mathbf{y}), \\ \nabla_{\mathbf{z}} \Phi_{\sigma}(\mathbf{w}) &= \mathcal{B} \text{prox}_{\sigma p}(\mathbf{x}_4 + \sigma(\mathcal{A}^* \mathbf{y} + \mathcal{B}^* \mathbf{z} - Q\mathbf{v} + \mathbf{r} - \mathbf{c})) - \text{prox}_{\sigma f}(\mathbf{x}_2 - \sigma \mathbf{z}), \\ \nabla_{\mathbf{r}} \Phi_{\sigma}(\mathbf{w}) &= \text{prox}_{\sigma p}(\mathbf{x}_4 + \sigma(\mathcal{A}^* \mathbf{y} + \mathcal{B}^* \mathbf{z} - Q\mathbf{v} + \mathbf{r} - \mathbf{c})) - \Pi_{\mathcal{P}_1}(\mathbf{x}_3 - \sigma \mathbf{r}), \\ \nabla_{\mathbf{v}} \Phi_{\sigma}(\mathbf{w}) &= -Q \text{prox}_{\sigma p}(\mathbf{x}_4 + \sigma(\mathcal{A}^* \mathbf{y} + \mathcal{B}^* \mathbf{z} - Q\mathbf{v} + \mathbf{r} - \mathbf{c})) + Q\mathbf{v}, \\ \nabla_{\mathbf{x}_1} \Phi_{\sigma}(\mathbf{w}) &= \frac{1}{\sigma} \Pi_{\mathcal{P}_2}(\mathbf{x}_1 - \sigma \mathbf{y}) - \frac{1}{\sigma} \mathbf{x}_1, \\ \nabla_{\mathbf{x}_2} \Phi_{\sigma}(\mathbf{w}) &= \frac{1}{\sigma} \text{prox}_{\sigma f}(\mathbf{x}_2 - \sigma \mathbf{z}) - \frac{1}{\sigma} \mathbf{x}_2, \\ \nabla_{\mathbf{x}_3} \Phi_{\sigma}(\mathbf{w}) &= \frac{1}{\sigma} \Pi_{\mathcal{P}_1}(\mathbf{x}_3 - \sigma \mathbf{r}) - \frac{1}{\sigma} \mathbf{x}_3, \\ \nabla_{\mathbf{x}_4} \Phi_{\sigma}(\mathbf{w}) &= \frac{1}{\sigma} \text{prox}_{\sigma p}(\mathbf{x}_4 + \sigma(\mathcal{A}^* \mathbf{y} + \mathcal{B}^* \mathbf{z} - Q\mathbf{v} + \mathbf{r} - \mathbf{c})) - \frac{1}{\sigma} \mathbf{x}_4. \end{aligned} \quad (3.7) \quad \{\text{eqn:gradient}\}$$

We note that if  $f^*$  is differentiable, then  $\mathbf{x}_2$  is not exist and the corresponding gradient is  $\nabla_{\mathbf{z}} \Phi_{\sigma}(\mathbf{w}) = \mathcal{B} \text{prox}_{\sigma p}(\mathbf{x}_4 + \sigma(\mathcal{A}^* \mathbf{y} + \mathcal{B}^* \mathbf{z} - Q\mathbf{v} + \mathbf{r} - \mathbf{c})) - \nabla f^*(-\mathbf{z})$ . We make the following existence condition of  $\mathbf{x}_1, \mathbf{x}_2$ , and  $\mathbf{x}_3$  that will be used in the subsequent analysis.

- $\mathbf{y}$  exist if and only if  $\mathcal{P}_2$  is nontrivial,  $\mathbf{x}_1$  exist if and only if  $\mathcal{P}_2$  is not a singleton set.
- $\mathbf{r}$  and  $\mathbf{x}_3$  exist if and only if  $\mathcal{P}_1$  is nontrivial.
- $\mathbf{x}_2$  exist if and only if  $f$  is nonsmooth i.e. needs proximal operator.
- $\mathbf{v}$  exist if and only if  $Q$  is nontrivial.

The nonlinear operator  $F(\mathbf{w})$  is defined as

$$F(\mathbf{w}) = \left( \nabla_{\mathbf{y}} \Phi(\mathbf{w}); \nabla_{\mathbf{z}} \Phi(\mathbf{w}); \nabla_{\mathbf{r}} \Phi(\mathbf{w}); \nabla_{\mathbf{v}} \Phi(\mathbf{w}); -\nabla_{\mathbf{x}_1} \Phi(\mathbf{w}); -\nabla_{\mathbf{x}_2} \Phi(\mathbf{w}); -\nabla_{\mathbf{x}_3} \Phi(\mathbf{w}); -\nabla_{\mathbf{x}_4} \Phi(\mathbf{w}) \right). \quad (3.8) \quad \{\text{eq:def:F}\}$$

Definition 3.2. [9] Let  $F$  be a locally Lipschitz continuous mapping.  $F$  is called semismooth at  $\mathbf{x}$  if  $F$  is directional differentiable at  $\mathbf{x}$  and for any  $\mathbf{d}$  and  $J \in \partial F(\mathbf{x} + \mathbf{d})$ , it holds that

$$\|F(\mathbf{x} + \mathbf{d}) - F(\mathbf{x}) - J\mathbf{d}\| = o(\|\mathbf{d}\|), \quad \mathbf{d} \rightarrow 0.$$

$F$  is said to be strongly semismooth at  $\mathbf{x}$  if  $F$  is semismooth at  $\mathbf{x}$  and

$$\|F(\mathbf{x} + \mathbf{d}) - F(\mathbf{x}) - J\mathbf{d}\| = O(\|\mathbf{d}\|^2), \quad \mathbf{d} \rightarrow 0.$$

We say  $F$  is semismooth (respectively, strongly semismooth) if  $F$  is semismooth (respectively, strongly semismooth) for any  $\mathbf{x}$ .

We first compute the generalized Jacobian of  $F(\mathbf{w})$ . Note that for a convex function  $h$ , its proximal operator  $\text{prox}_{th}$  is Lipschitz continuous. Then, by Definition 3.2, the following sets:

$$\begin{aligned} D_{\Pi_1} &:= \partial \text{prox}(x_3 - \sigma \mathbf{r}), \quad D_{\Pi_2} = \partial \text{prox}(x_1 - \sigma \mathbf{y}), \quad D_f := \partial \text{prox}_{\sigma f}(x_2 - \sigma \mathbf{z}), \\ D_p &:= \partial \text{prox}_{\sigma p}(x_4 + \sigma(\mathcal{A}^* \mathbf{y} + \mathcal{B}^* \mathbf{z} - Q\mathbf{v} + \mathbf{r} - \mathbf{c})) \end{aligned} \quad (3.9)$$

Consequently, the corresponding generalized Jacobian can be represented by

$$\hat{\partial} F(\mathbf{w}) := \left\{ \begin{pmatrix} \mathcal{H}_{11} & \mathcal{H}_{12} \\ -\mathcal{H}_{12}^\top & \mathcal{H}_{22} \end{pmatrix} \right\}. \quad (3.10)$$

where

$$\begin{aligned} \mathcal{H}_{11} &= \sigma(\mathcal{A}, \mathcal{B}, I, -Q)^\top D_p(\mathcal{A}, \mathcal{B}, I, -Q) + \sigma \text{blkdiag}(D_{\Pi_1}, D_f, D_{\Pi_2}, Q), \\ \mathcal{H}_{12} &= [(-\text{blkdiag}([D_{\Pi_1}, D_f, D_{\Pi_2}]); \mathbf{0}), (\mathcal{A}, \mathcal{B}, I, -Q)^\top D_{x_4}], \\ \mathcal{H}_{22} &= \text{blkdiag}\left\{ \frac{1}{\sigma}(I - D_{\Pi_1}), \frac{1}{\sigma}(I - D_h), \frac{1}{\sigma}(I - D_{\Pi_2}), \frac{1}{\sigma}(I - D_p) \right\}. \end{aligned} \quad (3.11)$$

The linear system can be represented by:

$$\begin{pmatrix} \mathcal{H}_{11} + \tau I & \mathcal{H}_{12} \\ -\mathcal{H}_{12}^\top & \mathcal{H}_{22} + \tau I \end{pmatrix} \begin{pmatrix} \mathbf{d}_1 \\ \mathbf{d}_2 \end{pmatrix} = \begin{pmatrix} -F_1 \\ -F_2 \end{pmatrix}, \quad (3.12)$$

For a given  $\mathbf{d}_1$ , the direction  $\mathbf{d}_2$  can be calculated by

$$\mathbf{d}_2 = (\mathcal{H}_{22} + \tau I)^{-1}(\mathcal{H}_{12}^\top \mathbf{d}_1 - F_2). \quad (3.13)$$

Hence, (3.12) can be reduced to a linear system with respect to  $\mathbf{d}_1$ :

$$\tilde{\mathcal{H}}_{11} \mathbf{d}_1 = -\tilde{F}_1, \quad (3.14)$$

where  $\tilde{F}_1 := \mathcal{H}_{12}(\mathcal{H}_{22} + \tau I)^{-1} F_2 - F_1$  and  $\tilde{\mathcal{H}}_{11} := (\mathcal{H}_{11} + \mathcal{H}_{12}(\mathcal{H}_{22} + \tau I)^{-1} \mathcal{H}_{12}^\top + \tau I)$ . The definition of  $\mathcal{H}_{12}$  in (3.10) yields

$$\tilde{\mathcal{H}}_{11} = \sigma(\mathcal{A}, \mathcal{B}, I, Q)^\top \bar{D}_p(\mathcal{A}, \mathcal{B}, I, Q) + \sigma \text{blkdiag}(\bar{D}_{\Pi_1}, \bar{D}_f, \bar{D}_{\Pi_2}, Q) \quad (3.15)$$

where  $D_p = \sigma \hat{D}_p + \bar{D}_p$ ,  $\tilde{D}_p = \hat{D}_p(\frac{1}{\sigma}I - \frac{1}{\sigma}\hat{D}_p + \tau I)^{-1}\hat{D}_p$ ,  $\bar{D}_{\Pi_1}, \bar{D}_f$  and  $\bar{D}_{\Pi_2}$  are defined analogously.

## 4 Proximal operator

According to (3.13), we need the explicit calculation process of  $(\frac{1}{\sigma}(I - D_p) + \tau I)^{-1}$  and. Furthermore, if  $\mathbf{x}$  and  $\mathcal{B}(\mathbf{x})$  are replaced by  $\mathbf{x} - \mathbf{b}_1$  or  $\mathcal{B}(\mathbf{x}) - \mathbf{b}_2$  respectively, the corresponding needs to be correct by a shift term. We demonstrate the computation details of some proximal operators in this section.

### 4.1 Handling shift term

For problems that have terms such as  $p(\mathbf{x} - \mathbf{b}_1)$  or  $f(\mathcal{B}(\mathbf{x}) - \mathbf{b}_2)$ , the corresponding dual problem is

$$\begin{aligned} \min_{\mathbf{y}, \mathbf{z}, \mathbf{s}, \mathbf{r}, \mathbf{v}} \quad & \delta_{\mathcal{P}_2}^*(-\mathbf{p}) + f^*(-\mathbf{q}) - \langle \mathbf{b}_2, \mathbf{q} \rangle - \langle \mathbf{b}_1, \mathbf{p} \rangle + p^*(-\mathbf{s}) + \frac{1}{2} \langle Q\mathbf{v}, \mathbf{v} \rangle + \delta_{\mathcal{P}_1}^*(-\mathbf{t}), \\ \text{s.t.} \quad & \mathcal{A}^* \mathbf{y} + \mathcal{B}^* \mathbf{z} + \mathbf{s} - Q\mathbf{v} + \mathbf{r} = \mathbf{c}, \quad \mathbf{y} = \mathbf{p}, \quad \mathbf{z} = \mathbf{q}, \quad \mathbf{r} = \mathbf{t}. \end{aligned} \quad (4.1)$$

If  $f^*$  is differentiable, then the corresponding gradient is  $-\nabla f^*(-\mathbf{q}) - \mathbf{b}_2$ . If  $f$  is nonsmooth, it follows from the property of the proximal operator that  $\mathbf{q} = \text{prox}_{f^*/\sigma}(\mathbf{x}_2/\sigma - \mathbf{z} - \mathbf{b}_2/\sigma)$ . Consequently, we only need to replace the

prox<sub>σf</sub> in (3.7) of prox<sub>σf</sub> - b<sub>2</sub>. For p<sup>\*</sup>(-s), the corresponding prox<sub>σp</sub>(x<sub>4</sub> + σ(ℳ<sup>\*</sup>y + ℬ<sup>\*</sup>z - Qv + r - c)) is replaced by prox<sub>σp</sub>(x<sub>4</sub> + σ(ℳ<sup>\*</sup>y + ℬ<sup>\*</sup>z - Qv + r - c)) - b<sub>1</sub>. Consequently, we do not need to introduce slack variables when adding shift term in p(x) or f(ℬ(x)).

## 4.2 ℓ<sub>2</sub> norm

For ℓ<sub>2</sub> norm, i.e. p(x) = λ||x||<sub>2</sub>, where x ∈ ℝ<sup>n</sup>. The corresponding proximal operator is prox<sub>λ||·||<sub>2</sub></sub> =  $\begin{cases} x - \lambda x / \|x\|, & \text{if } \|x\| > \lambda, \\ 0, & \text{otherwise.} \end{cases}$

According to the SMW formular: (A - uu<sup>T</sup>)<sup>-1</sup> = A<sup>-1</sup> +  $\frac{A^{-1}uu^TA^{-1}}{1-u^TA^{-1}u}$ . Consequently, for D ∈ ∂prox<sub>λ||·||<sub>2</sub></sub> · ||·||<sub>2</sub> we have

$$\begin{aligned} \left( \tau I + \frac{1}{\sigma} (I - D) \right)^{-1} &= \left( \tau I + \lambda (I - xx^T / \|x\|^2) / \|x\| \right)^{-1} \\ &= \left( \left( \tau + \frac{\lambda}{\|x\|} \right) I - \frac{\lambda}{\|x\|} xx^T / \|x\|^2 \right)^{-1} \\ &= \frac{1}{\left( \tau + \frac{\lambda}{\|x\|} \right)} \left( I + \frac{\lambda}{\tau \|x\|} x^T x / \|x\|^2 \right). \end{aligned}$$

## 4.3 Second order cone

Let Q ⊂ ℝ<sup>n</sup> denote the second-order cone (SOC), where a vector x ∈ Q is structured as x = [x<sub>0</sub>; x̄] with x<sub>0</sub> ∈ ℝ and x̄ ∈ ℝ<sup>n-1</sup>. The SOC is characterized by the condition Q = {x ∈ ℝ<sup>n</sup> : ||x̄|| ≤ x<sub>0</sub>}. For x ∈ int(Q), the determinant is defined as det(x) = x<sub>0</sub><sup>2</sup> - ||x̄||<sup>2</sup>, and the inverse of x, when det(x) ≠ 0, is given by x<sup>-1</sup> =  $\frac{1}{\det(x)} [x_0 \quad -\bar{x}^T]^T$ . For p(x) = δ<sub>Q</sub>(x), we utilize a barrier function p(x) = μg(x) to replace it, where μ > 0. Define the logarithmic barrier function for any x ∈ K by g : ℝ<sup>n</sup> ↦ ℝ with

$$g(x) = \begin{cases} -\frac{1}{2} \log(\det(x)), & \text{if } \|x\| < x_0, \\ +\infty, & \text{otherwise.} \end{cases}$$

We note that lim<sub>μ→0</sub> μg(x) = δ<sub>Q</sub>(x). For the smoothing function μg(x), we have the following lemma.

LEMMA 4.1. (i) The proximal mapping of μ · g(x) is given by prox<sub>μg</sub> : ℝ<sup>n</sup> ↦ int(K) with

$$\text{prox}_{\mu g}(z) = \begin{bmatrix} \frac{1}{2} \left( z_0 + \sqrt{\frac{1}{2}(\|z\|^2 + 4\mu + \Delta)} \right) \\ \frac{z}{2} \left( 1 + \frac{\sqrt{2}z_0}{\sqrt{\|z\|^2 + 4\mu + \Delta}} \right) \end{bmatrix}, \quad z \in \mathbb{R}^n$$

where Δ = √det(z)<sup>2</sup> + 8μ||z||<sup>2</sup> + 16μ<sup>2</sup>. Furthermore, the inverse function of the proximal mapping is given by prox<sub>μg</sub><sup>-1</sup> : int(K) → ℝ with

$$\text{prox}_{\mu g}^{-1}(x) = x - \mu x^{-1}, \quad x \in \text{int}(K).$$

(ii) The projection function is the limit of the proximal mapping as taking μ to 0. That is,

$$\lim_{\mu \rightarrow 0} \text{prox}_{\mu g}(z) = \Pi_K(z), \quad z \in \mathbb{R}^n.$$

(iii) For z ∈ ℝ<sup>n</sup>, let x = prox<sub>μg</sub>(z). The inverse matrix of the derivative of the proximal mapping at the point z is given by

$$(\partial_z \text{prox}_{\mu g}(z))^{-1} = I - \mu \partial_x(x^{-1}),$$

$$\text{where } \partial_x(x^{-1}) = \frac{1}{\det(x)} \begin{bmatrix} 1 & \\ & -I_{n-1} \end{bmatrix} - 2(x^{-1})(x^{-1})^\top.$$

(iv) The derivative of the proximal mapping at the point  $z$  is given by

$$\begin{aligned} \partial_z \text{prox}_{\mu g}(z) &= \begin{bmatrix} \frac{\det(x)}{\det(x)-\mu} & \\ & \frac{\det(x)}{\det(x)+\mu} I_{n-1} \end{bmatrix} - \frac{2\mu \det(x) \begin{bmatrix} \frac{x_0}{\det(x)-\mu} \\ \frac{-x}{\det(x)+\mu} \end{bmatrix} \begin{bmatrix} \frac{x_0}{\det(x)+\mu} \\ \frac{-x}{\det(x)+\mu} \end{bmatrix}^\top}{\det(x) + 2\mu \left( \frac{x_0^2}{\det(x)-\mu} + \frac{\|\tilde{x}\|^2}{\det(x)+\mu} \right)} \\ &:= \Lambda + auu^\top. \end{aligned}$$

The corresponding linear system is  $\sigma D + D \left( \frac{1}{\sigma} (I - D) + \tau_2 I \right)^{-1} D$ , where  $\Lambda = \begin{bmatrix} a_0 & \\ & a_1 I_{n-1} \end{bmatrix}$ . According to the SMW formula

$$(\Lambda - uu^\top)^{-1} = \Lambda^{-1} + \frac{\Lambda^{-1} uu^\top \Lambda^{-1}}{1 - u^\top \Lambda^{-1} u},$$

we have that  $\frac{1}{\sigma} (I - D) + \tau_2 I = \frac{1}{\sigma} \left( \begin{bmatrix} b_0 & \\ & b_1 I \end{bmatrix} - auu^\top \right) := \frac{1}{\sigma} (\Lambda_1 - auu^\top)$ ,  $\Lambda_1 = (1 + \sigma \tau_2) I - \Lambda$ . Consequently, it follows that  $(\frac{1}{\sigma} (I - D) + \tau_2 I)^{-1} = \sigma (\Lambda_1^{-1} + c \Lambda_1^{-1} uu^\top \Lambda_1^{-1})$ ,  $c = \frac{a}{1 - au^\top \Lambda_1^{-1} u}$  is a constant. Hence we have

$$\begin{aligned} &\sigma D + D \left( \frac{1}{\sigma} (I - D) + \tau_2 I \right)^{-1} D \\ &= \sigma \left( \Lambda + auu^\top \right) + \sigma (\Lambda + auu^\top) (\Lambda_1^{-1} + c \Lambda_1^{-1} uu^\top \Lambda_1^{-1}) (\Lambda + auu^\top) \\ &= \sigma \left( \Lambda + auu^\top \right) + \sigma (\Lambda \Lambda_1^{-1} + c \Lambda \Lambda_1^{-1} uu^\top \Lambda_1^{-1} + auu^\top \Lambda_1^{-1} + ac \gamma uu^\top \Lambda_1^{-1}) (\Lambda + auu^\top), \\ &= \sigma \left( \Lambda \Lambda_1^{-1} \Lambda + \Lambda + c \Lambda \Lambda_1^{-1} uu^\top \Lambda_1^{-1} \Lambda + a(1 + c\gamma) \Lambda \Lambda_1^{-1} uu^\top + a(1 + c\gamma) uu^\top \Lambda_1^{-1} \Lambda + (a + a^2 \gamma + a^2 c \gamma^2) uu^\top \right) \\ &= \sigma \left( \tilde{\Lambda} + \begin{bmatrix} b_0 u_0^2 & b_1 u_0 u_1^\top \\ b_1 u_0 u_1 & b_2 u_1 u_1^\top \end{bmatrix} \right). \end{aligned}$$

where  $\tilde{\Lambda} = \Lambda \Lambda_1^{-1} \Lambda + \Lambda$ ,  $\gamma = u^\top \Lambda_1^{-1} u$ , and  $b_0, b_1, b_2$  are constants. Denote  $\Lambda_1^{-1} \Lambda = \begin{bmatrix} c_0 & \\ & c_1 I \end{bmatrix}$ , we have  $b_0 = cc_0^2 + 2a(1 + c\gamma)c_0 + a + a^2 \gamma + a^2 c \gamma^2$ ,  $b_1 = cc_0 c_1 + a(1 + c\gamma)(c_0 + c_1) + a + a^2 \gamma + a^2 c \gamma^2$ ,  $b_2 = cc_1^2 + 2a(1 + c\gamma)c_1 + a + a^2 \gamma + a^2 c \gamma^2$ . Consequently, if we set  $\tilde{u} = [b_1/b_2 u_0; u_1]$ , we have  $b_2 \tilde{u} \tilde{u}^\top = \begin{bmatrix} b_1^2/b_2 u_0^2 & b_1 u_0 u_1^\top \\ b_1 u_0 u_1 & b_2 u_1 u_1^\top \end{bmatrix}$ . It follows that

$$\sigma D + D \left( \frac{1}{\sigma} (I - D) + \tau_2 I \right)^{-1} D = \sigma \left( \tilde{\Lambda} + \begin{bmatrix} (b_0 - b_1^2/b_2) u_0^2 & 0 \\ 0 & 0 \end{bmatrix} \right) + \sigma b_2 \tilde{u} \tilde{u}^\top.$$

Consequently, can be represented as a diagonal matrix plus a rank 1 matrix, which is significant in constructing the Schur matrix when using .

#### 4.4 Nuclear norm

For nuclear norm  $\|X\|_*$ , let the singular value decomposition of  $X$  denoted by  $X = U\Sigma V^T$ , where the its proximal operator can be presented by:

$$\text{prox}_{\lambda\|\cdot\|_*}(X) = U \text{diag}(T_\alpha(\lambda(X))) V^T. \quad (4.2)$$

Denote

$$\hat{D}_2(G) = U \left[ \frac{\Omega_{\sigma,\sigma}^\mu + \Omega_{\sigma,-\sigma}^\mu}{2} \odot G_1 + \frac{\Omega_{\sigma,\sigma}^\mu - \Omega_{\sigma,-\sigma}^\mu}{2} \odot G_1^T, (\Omega_{\sigma,0}^\mu \odot (G_2)) \right] V^T, \quad (4.3)$$

where  $\sigma = [\sigma^{(1)}, \dots, \sigma^{(m)}]$  is the tensor singular value of  $X$ ,  $G_1 = U^T G V_1 \in \mathbb{R}^{n_1 \times n_1}$ ,  $G_2 = U^T G V_2 \in \mathbb{R}^{n_1 \times (n_2 - n_1)}$  and  $(\Omega_{\sigma,\sigma}^\mu)$  is defined by:

$$(\Omega_{\sigma,\sigma}^\mu)_{ij} := \begin{cases} \partial_B \text{prox}_{\mu\|\cdot\|_1}(\sigma_i), & \text{if } \sigma_i = \sigma_j, \\ \frac{\text{prox}_{\mu\|\cdot\|_1}(\sigma_i) - \text{prox}_{\mu\|\cdot\|_1}(\sigma_j)}{\sigma_i - \sigma_j}, & \text{otherwise.} \end{cases} \quad (4.4)$$

Then for any  $G \in \mathbb{R}^{m \times n}$ , we have  $\hat{D}[G] = D[G]$ .

LEMMA 4.2. For operator  $\mathcal{T} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n} : \mathcal{T}(\mathcal{G}) = \Omega_1 \odot \mathcal{G} + \Omega_2 \odot \mathcal{G}^T$ . The inverse of  $\mathcal{T}$  is

$$\mathcal{T}^{-1}(\mathcal{G}) = (\Omega_s + \Omega_a) \odot \mathcal{G} + (\Omega_s - \Omega_a) \odot \mathcal{G}^T,$$

where  $\Omega_s = 1./[2(\Omega_1 + \Omega_2)]$ ,  $\Omega_a = 1./[2(\Omega_1 - \Omega_2)]$  and  $./$  means element wise division.

From the above lemma,  $(D_2^{\tau_4})^{-1}$  can be represented as

$$\begin{aligned} (D_2^{\tau_4})^{-1}(\mathcal{G}) &= \mathcal{U} \left[ (\Omega_s^{\tau_4} + \Omega_a^{\tau_4}) \odot \mathcal{G}_1 + (\Omega_s^{\tau_4} - \Omega_a^{\tau_4}) \odot \mathcal{G}_1^T, (1./\Omega_3^{\tau_4} \odot \mathcal{G}_2) \right] \mathcal{V}^T \\ &\quad + \sigma/(1 + \sigma\tau_4)\mathcal{G}. \end{aligned} \quad (4.5)$$

where  $\Omega_s^{\tau_4} = 1./[2(\Omega_1^{\tau_4} + \Omega_2^{\tau_4})] - \sigma/(1 + \sigma\tau_4)E$ ,  $\Omega_a^{\tau_4} = 1./[2(\Omega_1^{\tau_4} - \Omega_2^{\tau_4})]$ ,  $E$  is the matrix of ones with correct size. The detail of the computation process is summarized in Algorithm 1. Then the low-rank structure can be fully exploited and the total computational cost for each inner iteration is reduced to  $O(n_1 r^2)$

#### 4.5 Fused reuglarizer

In this case,  $p(x) = \lambda_1 \|x\|_1 + \lambda_2 \|Bx\|_1$ , where  $B(x) = [x_2 - x_1, \dots, x_n - x_{n-1}]$ . According to lemma [5], the corresponding proximal operator is

$$\text{prox}_p(v) = \text{prox}_{\lambda_1\|\cdot\|_1}(x_{\lambda_2}(v)) = \text{prox}_{\lambda_1\|\cdot\|_1}(v - B^T z_{\lambda_2}(Bv)), \quad (4.6) \quad \{\text{prox-fused}\}$$

where  $z_{\lambda_2}(u) := \text{argmin} \left\{ \frac{1}{2} \|B^T z\|^2 - \langle z, u \rangle \mid \|z\|_\infty \leq \lambda_2 \right\}$ ,  $\forall u \in \mathbb{R}^{n-1}$ . Define the multifunction  $\mathcal{M} : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times n}$  by:

$$\mathcal{M}(v) := \{M \in \mathbb{R}^{n \times n} \mid M = \Theta P, \Theta \in \partial_B \text{prox}_{\lambda_1\|\cdot\|_1}(x_{\lambda_2}), P \in \mathcal{P}_x(v)\}, \quad (4.7) \quad \{\text{genJacobian:fused}\}$$

and  $\mathcal{P}_x(v) := \{\hat{P} \in \mathbb{R}^{(n-1) \times (n-1)} \mid \hat{P} = I - B^T(\Sigma_K B B^T \Sigma_K)^{\dagger} B, K \in \mathcal{K}_z(v)\}$ , where  $\Sigma_K = \text{Diag}(\sigma_K) \in \mathbb{R}^{(n-1) \times (n-1)}$  with

$$(\sigma_K)_i = \begin{cases} 0, & \text{if } i \in K \\ 1, & \text{otherwise, } i = 1, \dots, n-1. \end{cases}$$

---

**Algorithm 1** The procedure of computing  $(D_2^{\tau_4})^{-1}(\mathcal{G})$ .

---

**Require:**  $\mathcal{G}, (\Omega_1^{\tau_4})_{\alpha\alpha}, (\Omega_1^{\tau_4})_{\alpha\bar{\alpha}}, (\Omega_2^{\tau_4})_{\alpha\alpha}, (\Omega_2^{\tau_4})_{\alpha\bar{\alpha}}, U = [U_\alpha, U_{\bar{\alpha}}], V = [V_\alpha, V_{\bar{\alpha}}, V_\beta], (\Omega_1^{\tau_4})_{\alpha\beta}$ , penalty parameter  $\sigma$  and  $\tau_4$ .

**Ensure:**  $(D_2^{\tau_4})^{-1}(\mathcal{G})$

1: Compute  $\bar{\mathcal{G}} = \text{fft}(\mathcal{G}, [], 3)$ ,

2: Compute  $(\mathcal{G}_1)_{\alpha\alpha}, (\mathcal{G}_1)_{\alpha\bar{\alpha}}, (\mathcal{G}_1)_{\bar{\alpha}\alpha}, (\mathcal{G}_1)_{\alpha\beta}$  where

$$(\mathcal{G}_1)_{\alpha\alpha} = U_\alpha^T \bar{\mathcal{G}} (V_\alpha), \quad (\mathcal{G}_1)_{\alpha\bar{\alpha}} = U_\alpha^T \bar{\mathcal{G}} (V_{\bar{\alpha}}),$$

$$(\mathcal{G}_1)_{\bar{\alpha}\alpha} = U_{\bar{\alpha}}^T \bar{\mathcal{G}} (V_\alpha), \quad (\mathcal{G}_1)_{\alpha\beta} = U_\alpha^T \bar{\mathcal{G}} (V_\beta).$$

3: Compute

$$\mathcal{G}_2 = \begin{bmatrix} (\mathcal{G}_2)_{\alpha\alpha} & (\mathcal{G}_2)_{\alpha\bar{\alpha}} & (\mathcal{G}_2)_{\alpha\beta} \\ (\mathcal{G}_2)_{\bar{\alpha}\alpha} & 0 & 0 \end{bmatrix},$$

where

$$(\mathcal{G}_2)_{\alpha\alpha} = (\Omega_1^{\tau_4})_{\alpha\alpha} \odot (\mathcal{G}_1)_{\alpha\alpha} + (\Omega_2^{\tau_4})_{\alpha\alpha} \odot ((\mathcal{G}_1)_{\alpha\alpha})^T,$$

$$(\mathcal{G}_2)_{\alpha\bar{\alpha}} = (\Omega_1^{\tau_4})_{\alpha\bar{\alpha}} \odot (\mathcal{G}_1)_{\alpha\bar{\alpha}} + (\Omega_2^{\tau_4})_{\alpha\bar{\alpha}} \odot ((\mathcal{G}_1)_{\alpha\bar{\alpha}})^T,$$

$$(\mathcal{G}_2)_{\bar{\alpha}\alpha} = ((\Omega_1^{\tau_4})_{\alpha\bar{\alpha}})^T \odot (\mathcal{G}_1)_{\bar{\alpha}\alpha} + ((\Omega_2^{\tau_4})_{\alpha\bar{\alpha}})^T \odot ((\mathcal{G}_1)_{\alpha\bar{\alpha}})^T,$$

$$(\mathcal{G}_2)_{\alpha\beta} = (\Omega_1^{\tau_4})_{\alpha\beta} \odot (\mathcal{G}_1)_{\alpha\beta}.$$

4: Compute  $\mathcal{G}_3 = \mathcal{G}_{12} + \mathcal{G}_{11} + \mathcal{G}_{21} + \mathcal{G}_{13}$  where

$$\mathcal{G}_{11} = U_\alpha (\mathcal{G}_2)_{\alpha\alpha} V_\alpha^T, \quad \mathcal{G}_{12} = U_\alpha (\mathcal{G}_2)_{\alpha\bar{\alpha}} V_{\bar{\alpha}}^T,$$

$$\mathcal{G}_{21} = U_{\bar{\alpha}} (\mathcal{G}_2)_{\bar{\alpha}\alpha} V_\alpha^T, \quad \mathcal{G}_{13} = U_\alpha (\mathcal{G}_2)_{\alpha\beta} V_\beta^T.$$

5: Compute  $(D_2^{\tau_4})^{-1}(\mathcal{G}) = \text{ifft}(\mathcal{G}, [], 3) + \sigma/(1 + \sigma\tau_4)\mathcal{G}$ .

---

In fact, Define  $\Gamma := I_n - B^T(\Sigma B B^T \Sigma)^{\dagger} B = \text{Diag}(\Gamma_1, \dots, \Gamma_N)$ , where for  $i = 1, \dots, N$ ,

$$\Gamma_i = \begin{cases} \frac{1}{n_i+1} \mathbf{E}_{n_i+1}, & \text{if } i \in J \\ I_{n_i}, & \text{if } i \in \{1, N\} \\ I_{n_i-1}, & \text{otherwise.} \end{cases}$$

Moreover,  $\Gamma = H + U U^T = H + U_J U_J^T$ , where  $H \in \mathbb{R}^{n \times n}$  is a N-block diagonal matrix given by  $H =$ . Then  $D = \Theta P$ , with  $P = I - B^T(\Sigma B B^T \Sigma)^{\dagger} B$ , where

$$\theta_i = \begin{cases} 0, & \text{if } |(x_{\lambda_2}(v))_i| \leq \lambda_1 \\ 1, & \text{otherwise, } i = 1, \dots, n, \end{cases}$$

and let  $I_z(v) := \{i \mid |(z_{\lambda_2}(Bv))_i| = \lambda_2, i = 1, \dots, n-1\}$ . Then  $\sigma = \text{Diag}(\sigma) \in \mathbb{R}^{(n-1) \times (n-1)}$  with

$$\sigma_i = \begin{cases} 0, & \text{if } i \in I_z(v) \\ 1, & \text{otherwise, } i = 1, \dots, n-1. \end{cases}$$

It follows that  $\Theta \in \partial_{B \text{prox}_{\lambda_1 \|\cdot\|_1}}(x_{\lambda_2}(v))$  and  $P \in \mathcal{P}_x(v)$ . We first take a Clarke's generalized Jacobian:  $J^k \in \hat{\partial} F(\mathbf{w}^k)$ . The problem is transformed into solving the following problem by Gaussian elimination:

$$((\tau_1 + 1)I + \mathcal{A} \tilde{D} \mathcal{A}^T) \mathbf{d}_y^k = \mathcal{A} D \hat{D}^{-1} F_x - F_y,$$

where  $\hat{D} = D\hat{D}^{-1}D + \sigma D$ ,  $\hat{D} = \frac{1}{\sigma}(I - D) + \tau_2 I$  and  $D \in \hat{\partial}\text{prox}_{\sigma p}(\mathbf{x} - \sigma \mathcal{A}^* \mathbf{y})$  is an element of the generalized Jacobian of  $\text{prox}_{\sigma p}$ . Then

$$d_{\mathbf{x}}^k = (\hat{D})^{-1}(-F_{\mathbf{x}} + D\mathcal{A}^T d_{\mathbf{y}}^k).$$

Therefore,  $M = \Theta(H + U_J U_J^T) = \Theta(H + U_J U_J^T)\Theta$ ,  $\Theta^2 = \Theta$ ,  $H^2 = H$ ,  $\Theta H = \Theta H \Theta$ . Define the following index sets:

$$\alpha := \{i | \theta_i = 1, i \in \{1, \dots, n\}\}, \quad \beta := \{i | h_i = 1, i \in \alpha\},$$

where  $\theta_i$  and  $h_i$  are the  $i$ -th diagonal entries of matrices  $\Theta$  and  $H$  respectively. Then we have

$$A\Theta H A^T = A\Theta H \Theta A^T = A_{\alpha} H A_{\alpha}^T = A_{\beta} A_{\beta}^T,$$

where  $A_{\alpha} \in \mathbb{R}^{m \times |\alpha|}$  and  $A_{\beta} \in \mathbb{R}^{m \times |\beta|}$  are two submatrix obtained from  $A$  by extracting those columns with indices in  $\alpha$  and  $\beta$ . Meanwhile, we have

$$A\Theta(U_J U_J^T)A^T = A\Theta(U_J U_J^T)\Theta A^T = A_{\alpha} \tilde{U} \tilde{U}^T A_{\alpha}^T,$$

where  $\tilde{U} \in \mathbb{R}^{|\alpha| \times r}$  is a submatrix obtained from  $\Theta U_J$  by extracting those rows with indices in  $\alpha$  and the zero columns in  $\Theta U_J$  are removed. Therefore, by exploiting the structure in  $D$ , we can express  $ADA^T$  in the following form:

$$ADA^T = A_{\beta} A_{\beta}^T + A_{\alpha} \tilde{U} \tilde{U}^T A_{\alpha}^T.$$

For  $D\hat{D}^{-1}D$  case, where  $\hat{D} = \frac{1}{\sigma}(I - D) + \tau_2 I$ ,  $D = \Theta H \Theta$ , we should note that  $D = \Theta H \Theta = \Theta H$  holds since  $\Theta = \text{Diag}(\Theta_1, \dots, \Theta_N)$ . Thus  $M = \text{Diag}(\Theta_1 \Gamma_1, \dots, \Theta_N \Gamma_N)$ . Define  $J := \{j | \Gamma_j \text{ is not an identity matrix}, 1 \leq j \leq N\}$ . Since  $\text{supp}(Bx_{\lambda_2}(v)) \subset K$ , we have that

$$\Theta_j = \mathbf{0}_{n_j+1} \text{ or } I_{n_j+1} \forall j \in J,$$

which implies  $\Theta_j \Gamma_j \in \mathbb{S}_+^{n_j+1} \forall j \in J$ . Therefore,  $D \in \mathbb{S}_+^n$ . As a consequence,

$$D = \text{Diag}(D_1, \dots, D_n),$$

where

$$D_i = \begin{cases} \frac{1}{n_i+1} \mathbf{E}_{n_i+1}, & \text{if } i \in J \text{ and } \Theta_i = I_{n_i} \\ I_{n_i}, & \text{if } i \notin J \text{ and } i \in \{i, N\}, \\ 0, & \text{if } \Theta_i = 0, \\ I_{n_i-1}, & \text{otherwise.} \end{cases}$$

According to the SMW formular:  $(A - vv^T)^{-1} = A^{-1} + \frac{A^{-1}vv^T A^{-1}}{1 - v^T A^{-1}v}$ , the inverse of  $\hat{D}^{-1}$  has explicit solution:

$$\begin{aligned} ((\frac{1}{\sigma} + \tau_2)I - \frac{1}{\sigma}D)^{-1} &= ((\frac{1}{\sigma} + \tau_2)I - \frac{1}{\sigma}\Theta(H + U_J U_J^T)\Theta)^{-1} \\ &= \begin{cases} \frac{\sigma}{1+\tau_2\sigma} I_{n_i+1} + \frac{1}{(1+\tau_2\sigma)(\tau_2)(n_i+1)} \mathbf{E}_{n_i+1}, & \text{if } 1 \\ \frac{1}{\tau_2} I_{n_i}, & \text{if } i \notin J \text{ and } i \in \{i, N\}, \\ \frac{\sigma}{1+\sigma\tau_2} I_{n_i}, & \text{if } \Theta_i = 0, \\ \frac{1}{\tau_2} I_{n_i-1}, & \text{otherwise.} \end{cases} \end{aligned}$$

As a consequence,  $\tilde{D}$  has the following formula:

$$\tilde{D} = \begin{cases} \left(\frac{1}{\tau_2} + \sigma\right) \frac{1}{n_i+1} \mathbf{E}_{n_i+1}, & \text{if } i \in J \text{ and } \Theta_i = I_{n_i} \\ \frac{1}{\tau_2} + \sigma I_{n_i}, & \text{if } i \notin J \text{ and } i \in \{i, N\}, \text{ and } \Theta_i = I_{n_i} \\ 0, & \text{if } \Theta_i = 0, \\ \frac{1}{\tau_2} + \sigma I_{n_i-1}, & \text{otherwise.} \end{cases}$$

and  $D\hat{D}^{-1}$  has the following formula:

$$D\hat{D}^{-1} = \begin{cases} \left(\frac{1}{(\tau_2)(n_i+1)}\right) \mathbf{E}_{n_i+1}, & \text{if } i \in J \text{ and } \Theta_i = I_{n_i} \\ \frac{1}{\tau_2} I_{n_i}, & \text{if } i \notin J, i \in \{i, N\}, \text{ and } \Theta_i = I_{n_i} \\ 0, & \text{if } \Theta_i = 0, \\ \frac{1}{\tau_2} I_{n_i-1}, & \text{otherwise.} \end{cases}$$

Define  $\tilde{D} = \Theta \tilde{H}$  and the following index sets:

$$\alpha := \{i | \theta_i = 1, i \in \{1, \dots, n\}\}, \beta := \{i | h_i = 1, i \in \alpha\},$$

Then we have  $A\Theta\tilde{H}A^T = A\Theta\tilde{H}\Theta A^T = A_\alpha \tilde{U} \tilde{U} A_\alpha^T + \tilde{A}_\beta \tilde{A}_\beta^T$ , where  $\tilde{U}$  and  $\tilde{A}$  are the scaling matrix of  $U$  and  $A$ . Then we have the following decomposition:

$$A\tilde{D}A^T = W_1 W_2^T,$$

where  $W_1 := [\tilde{A}_\beta, A_\alpha \tilde{U} \tilde{U}^T] \in \mathbb{R}^{m \times (|\alpha| + |\beta|)}$ ,  $W_2 = [\tilde{A}_\beta, A_\alpha]$ . Using the above decomposition and we obtain

$$((\tau_1 + 1)I + A\tilde{D}A^T)^{-1} = \frac{1}{\tau_1 + 1} I_m - \frac{1}{\tau_1 + 1} \tilde{W}_1 ((\tau_1 + 1)I_{|\alpha| + |\beta|} + \tilde{W}_2^T \tilde{W}_2)^{-1} \tilde{W}_2^T.$$

Thus, we only need to factorize an  $(|\alpha| + |\beta|) \times (|\alpha| + |\beta|)$  matrix and the total computational cost is merely  $\mathcal{O}(|\alpha| + |\beta|)^3 + \mathcal{O}(m(|\alpha| + |\beta|)^2)$ .

## 5 Numerical experiments

The criteria to measure the accuracy of  $(\mathbf{y}, \mathbf{z}, \mathbf{v}, \mathbf{r}, \mathbf{x})$  is based on KKT optimality: conditions

$$\eta = \max\{\eta_P, \eta_D, \eta_K, \eta_{\mathcal{P}}\},$$

where

$$\begin{aligned} \eta_P &:= \frac{\|\mathcal{A}\mathbf{x} - \Pi_{\mathcal{P}_2}(\mathcal{A}\mathbf{x} - \mathbf{y})\|}{1 + \|\mathbf{x}\|}, \eta_D := \frac{\|\mathcal{A}^*(\mathbf{y}) + \mathcal{B}^*(\mathbf{z}) + \mathbf{s} - \mathcal{Q}(\mathbf{v}) - \mathbf{c}\|}{1 + \|\mathbf{c}\|}, \\ \eta_K &:= \min \left\{ \frac{\|\mathbf{x} - \text{prox}_{\mathcal{P}}(\mathbf{x} - \mathbf{s})\|}{1 + \|\mathbf{s}\| + \|\mathbf{x}\|}, \frac{\|\mathcal{Q}(\mathbf{v}) - \mathcal{Q}(\mathbf{x})\|_{\text{F}}}{1 + \|\mathcal{Q}(\mathbf{v})\| + \|\mathcal{Q}(\mathbf{x})\|} \right\}, \\ \eta_{\mathcal{P}} &= \min \left\{ \frac{\|\text{prox}_{\mathcal{F}}(\mathcal{B}\mathbf{x} - \mathbf{z}) - \mathcal{B}(\mathbf{x})\|}{1 + \|\mathcal{B}(\mathbf{x})\| + \|\mathbf{z}\|} \text{ or } \frac{\|-\nabla f^*(-\mathbf{z}) - \mathcal{B}(\mathbf{x})\|}{1 + \|\mathcal{B}(\mathbf{x})\| + \|\mathbf{z}\|}, \frac{\|\Pi_{\mathcal{P}_1}(\mathbf{x} - \mathbf{r}) - \mathbf{x}\|}{1 + \|\mathbf{x}\| + \|\mathbf{r}\|} \right\}. \end{aligned}$$

We also compute the relative gap by

$$\eta_g = \frac{|\text{pobj} - \text{dobj}|}{1 + |\text{pobj}| + |\text{dobj}|}.$$

For given accuracy  $\eta$ , we terminate SSNCVX when  $\eta < \text{opts.tol}$ .



## 5.1 Dataset Source

We are going to test the problem in data from Random Gaussian, UCI data<sup>1</sup> and LIBSVM dataset<sup>2</sup>. These data sets are collected from 10-K Corpus [4] and the UCI data repository [7]. As suggested in [3], for the data sets **pyrim**, **triazines**, **abalone**, **bodyfat**, **housing**, **mpg**, and **space\_ga**, we expand their original features by using polynomial basis functions over those features. For example, the last digit in **pyrim5** indicates that an order 5 polynomial is used to generate the basis functions. This naming convention is also used in the rest of the expanded data sets. These test instances, shown in Table 5, can be quite difficult in terms of the problem dimensions and the largest eigenvalue of  $\mathcal{A}\mathcal{A}^*$ , which is denoted as  $\lambda_{\max}(\mathcal{A}\mathcal{A}^*)$ .

## 5.2 Lasso case

In the following table,  $m$  denotes the number of number of samples,  $n$  denotes the number of features, and “nnz” denotes the number of nonzeros in the solution  $x$  obtained by the optimal solution using the following estimation:

$$\text{nnz} := \min\{k \mid \sum_{i=1}^k |\hat{x}_i| \geq 0.999\|x\|_1\},$$

where  $\hat{x}$  is obtained by sorting  $x$  such that  $|\hat{x}_1| \geq |\hat{x}_2| \geq \dots \geq |\hat{x}_n|$ . The stop criterion is set by the following relative KKT residual:

$$\eta = \frac{\|\tilde{x} - \text{prox}_{\lambda\|\cdot\|_1}(\tilde{x} - \mathcal{A}^*(\mathcal{A}\tilde{x} - b))\|_1}{1 + \|\tilde{x}\| + \|\mathcal{A}\tilde{x} - b\|}.$$

We stop the tested algorithms when  $\eta < 1e - 6$ . The algorithms we want to compare is SSNAL: SSN based on ALM (SOTA, Best Paper Prize for Young Researchers in Continuous Optimization 2019).<sup>3</sup> FPC\_AS.<sup>4</sup> A accelerated proximal gradient based method, ADMM algorithm, APG algorithm. The test results for different lambda choice, i.e.  $\lambda = 10^{-3}\|\mathcal{A}^T b\|_\infty$  and  $\lambda = 10^{-4}\|\mathcal{A}^T b\|_\infty$  and different algorithms are given in Table 3 and 4. We can see that, both SSNCP and SSNAL has successfully solved all problems while other first-order methods can not. Furthermore,

id	nnz	SSNCP		SSNAL		FPC_AS		ADMM	
		it	time	it	time	$\eta$	time	$\eta$	time
uci_CT	13	25	<b>0.57</b>	31	0.64	<b>1.4-4</b>	58:25	9.1-7	10:31
log1p.E2006.train	5	35	<b>11.57</b>	48	26.62	<b>2.6-1</b>	7:00:01	9.9-7	33:47
E2006.test	1	17	<b>0.19</b>	18	0.28	<b>3.7-4</b>	19:25	6.7-7	18
log1p.E2006.test	8	33	<b>3.27</b>	45	5.12	<b>6.9-1</b>	4:56:49	9.9-7	3:10
pyrim5	72	57	<b>1.37</b>	97	2.16	<b>8.8-1</b>	1:12:55	<b>3.2-5</b>	21:14
triazines4	519	64	<b>9.91</b>	92	11.23	<b>9.6-1</b>	7:00:05	<b>3.2-4</b>	2:13:00
abalone7	24	28	<b>0.72</b>	57	1.06	<b>Error</b>	<b>Error</b>	9.9-7	9:52
bodyfat7	2	29	<b>0.79</b>	40	1.08	<b>3.6-1</b>	1:12:02	9.9-7	1:49
housing7	158	57	<b>1.65</b>	57	1.74	<b>8.4-1</b>	1:41:01	9.9-7	22:12
mpg7	47	40	<b>0.07</b>	50	0.11	<b>Error</b>	<b>Error</b>	9.9-7	07
spacega9	14	26	<b>0.26</b>	38	1.01	<b>1.5-2</b>	30:55	9.9-7	37
E2006.train	1	24	<b>0.82</b>	18	0.87	<b>1.4-3</b>	1:24:18	9.1-7	10:41

Table 3. The results of our algorithm on Lasso problem with real data( $\lambda = 10^{-3}\|\mathcal{A}^T b\|_\infty$ ), the

<sup>1</sup><https://archive.ics.uci.edu/>

<sup>2</sup><https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

<sup>3</sup><https://github.com/MatOpt/SuiteLasso>

<sup>4</sup>[http://www.caam.rice.edu/optimization/L1/FPC\\_AS/](http://www.caam.rice.edu/optimization/L1/FPC_AS/)

id	nnz	SSNCP		SSNAL		FPC_AS		ADMM	
		it	time	it	time	$\eta$	time	$\eta$	time
uci_CT	44	44	<b>1.26</b>	53	1.75	<b>2.6-1</b>	0.31	9.9-7	0.31
log1p.E2006.train	599	55	<b>48.92</b>	123	78.49	<b>1.3-1</b>	7:00:01	<b>1.7-1</b>	35:21
E2006.test	1	18	<b>0.20</b>	18	0.26	<b>2.7-4</b>	20:21	6.3-7	21
log1p.E2006.test	1081	114	<b>25.24</b>	176	34.33	<b>8.7-1</b>	4:48:29	9.9-7	2:28
pyrim5	78	77	<b>2.01</b>	113	2.53	<b>9.8-1</b>	54:20	<b>1.6-3</b>	21:34
triazines4	260	105	<b>18.48</b>	147	28.06	<b>9.8-1</b>	7:47:03	<b>1.1-3</b>	1:55:16
abalone7	59	52	<b>1.63</b>	63	2.00	<b>Error</b>	<b>Error</b>	9.9-7	9:36
bodyfat7	3	41	<b>1.14</b>	56	1.51	<b>1.9-1</b>	1:13:08	9.8-7	4:05
housing7	281	70	<b>2.51</b>	68	2.52	<b>4.2-1</b>	1:39:36	<b>6.6-6</b>	42:36
mpg7	128	49	<b>0.11</b>	56	0.18	<b>Error</b>	<b>Error</b>	9.9-7	11
spacega9	38	48	<b>0.53</b>	54	0.72	<b>4.0-2</b>	30:26	9.7-7	0.31
E2006.train	1	21	<b>0.75</b>	49	0.88	<b>9.9-4</b>	1:27:29	9.9-7	10:53

Table 4. The results of our algorithm on Lasso problem with real data( $\lambda = 10^{-4} \|\mathcal{A}^T \mathbf{b}\|_\infty$ ), the

{tab:1e4}

Table 5. Statistics of the UCI test instances .

Probname	$(m, n)$	$\lambda_{\max}(\mathcal{A}\mathcal{A}^*)$
E2006.train	(3308, 72812)	1.912+05
log1p.E2006.train	(16087, 4265669)	5.86e+07
E2006.test	(3308, 72812)	4.79e+04
log1p.E2006.test	(3308, 1771946)	1.46e+07
pyrim5	(74, 169911)	1.22e+06
triazines4	(186, 557845)	2.07e+07
abalone7	(4177, 6435)	5.21e+05
bodyfat7	(252, 116280)	5.29e+04
housing7	(506, 77520)	3.28e+05
mpg7	(392, 3432)	1.28e+04
spacega9	(3107, 5005)	4.01e+03

{tab:tensor-image}

### 5.3 Fused Lasso case

For fused lasso case, part of the numerical results are listed in table 6.

id	nnz(x)	nnz(Bx)	SSN		SSNAL		iADMM		ADMM	
			it	time	it	time	$\eta$	time	$\eta$	time
uci_CT	8	1	44	<b>0.25</b>	53	0.42	7.6-7	3:29	2.5-7	2:55
log1p.E2006.train	32	5	25	<b>12.6</b>	11	19.6	7.6-7	3:29	2.5-7	2:55
E2006.test	1	1	9	<b>0.17</b>	5	0.33	5.3-7	2:21	6.3-7	4:56
log1p.E2006.test	32	5	25	<b>5.75</b>	13	6.37	5.3-7	2:21	6.3-7	4:56
pyrim5	1135	74	62	<b>2.34</b>	22	3.06	6.0-7	7:59	9.8-7	2:22
triazines4	2670	207	41	<b>7.2</b>	56	9.45	6.0-7	7:59	9.8-7	2:22
bodyfat7	63	8	44	<b>0.81</b>	14	0.89	6.0-7	7:59	9.8-7	2:22
abalone7	1	1	33	<b>1.02</b>	14	1.17	6.0-7	7:59	9.8-7	2:22
housing7	213	25	70	<b>1.8</b>	68	2.4	9.2-7	17:24	9.9-7	4:20
mpg7	42	11	48	<b>0.11</b>	54	0.15	7.6-7	30:26	9.7-7	7:33
spacega9	24	11	48	<b>0.37</b>	54	0.44	7.6-7	30:26	9.7-7	7:33
E2006.train	1	1	9	<b>0.58</b>	11	1.12	5.3-7	36:42	9.9-7	10:53

Table 6. The results of our algorithm on Lasso problem with real data( $\lambda_1 = 10^{-3} \|A^* \mathbf{b}\|_\infty$  and  $\lambda_2 = 5\lambda_1$ .)

{tab:fused}

904 **5.4 QP case**

905 **5.5 SDP**

906 **5.6 LP case**

907 **5.7 SOCP case**

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

## References

- [1] MOSEK ApS. 2019. *The MOSEK optimization toolbox for MATLAB manual. Version 10.1.0*. <http://docs.mosek.com/10.1/toolbox/index.html>
- [2] Stephen R Becker, Emmanuel J Candès, and Michael C Grant. 2011. Templates for convex cone problems with applications to sparse signal recovery. *Mathematical programming computation* 3 (2011), 165–218.
- [3] Ling Huang, Jinzhu Jia, Bin Yu, Byung-Gon Chun, Petros Maniatis, and Mayur Naik. 2010. Predicting execution time of computer programs using sparse polynomial regression. *Advances in neural information processing systems* 23 (2010).
- [4] Shimon Kogan, Dmitry Levin, Bryan R Routledge, Jacob S Sagi, and Noah A Smith. 2009. Predicting risk from financial reports with regression. In *Proceedings of human language technologies: the 2009 annual conference of the North American Chapter of the Association for Computational Linguistics*. 272–280.
- [5] Xudong Li, Defeng Sun, and Kim-Chuan Toh. 2018. A highly efficient semismooth Newton augmented Lagrangian method for solving Lasso problems. *SIAM Journal on Optimization* 28, 1 (2018), 433–458.
- [6] Xudong Li, Defeng Sun, and Kim-Chuan Toh. 2018. On efficiently solving the subproblems of a level-set method for fused Lasso problems. *SIAM Journal on Optimization* 28, 2 (2018), 1842–1866.
- [7] Moshe Lichman et al. 2013. UCI machine learning repository.
- [8] Jun Liu, Shuiwang Ji, Jieping Ye, et al. 2009. SLEP: Sparse learning with efficient projections. *Arizona State University* 6, 491 (2009), 7.
- [9] Robert Mifflin. 1977. Semismooth and semiconvex functions in constrained optimization. *SIAM Journal on Control and Optimization* 15, 6 (1977), 959–972.