Convergence and Applications of ADMM on the Multi-convex Problems

Junxiang Wang and Liang Zhao jwan936@emory.edu lzhao41@emory.edu

Emory University, 201 Dowman Drive, Atlanta, Georgia, USA

Abstract. In recent years, although the Alternating Direction Method of Multipliers (ADMM) has been empirically applied widely to many multi-convex applications, delivering an impressive performance in areas such as nonnegative matrix factorization and sparse dictionary learning, there remains a dearth of generic work on proposed ADMM with a convergence guarantee under mild conditions. In this paper, we propose a generic ADMM framework with multiple coupled variables in both objective and constraints. Convergence to a Nash point is proven with a sublinear convergence rate o(1/k). Two important applications are discussed as special cases under our proposed ADMM framework. Extensive experiments on ten real-world datasets demonstrate the proposed framework's effectiveness, scalability, and convergence properties. We have released our code at https://github.com/xianggebenben/miADMM.

1 Introduction

Due to the advantages and popularity of non-differentiable regularized and distributive computing for complex optimization problems, the Alternating Direction Method of Multipliers (ADMM) has received a great deal of attention in recent years [6]. The standard ADMM was originally proposed to solve the following separable convex optimization problem:

 $\min_{x,z} f(x) + g(z) \quad s.t. \ Ax + Bz = c.$

where f(x) and g(z) are closed convex functions, A and B are matrices and c is a vector. There are extensive reports in the literature exploring the theoretical properties for convex optimization problems related to ADMM and its variants, including multi-block ADMM [12], Bregman ADMM [30], fast ADMM [14, 18], and stochastic ADMM [24]. ADMM has now been extended to cover a wide range of nonconvex problems, and it has achieved outstanding performance in many practical applications [40].

Unlike convex problems, the convergence theory on the nonconvex ADMM is much more difficult, and considerable progress has been made on this problem, please refer to Section 2 for a detailed summary. Recently, however, there has been an increasing number of real-world applications where the objective functions are multi-convex (i.e. nonconvex for all the variables but convex for each when all the others are fixed). For example, nonnegative matrix factorization, which aims to decompose a matrix into a product of two matrices, has been applied widely in computer vision, machine learning, and various other fields [19]; A bilinear matrix inequality problem has been designed for the analysis of linear and nonlinear uncertain systems [16].

All of the above applications can be considered as special cases of multi-convex optimization problems. However, such problems have not yet been rigorously and systematically investigated by ADMM. Moreover, the convergence properties of the ADMM required to solve such problems remain unknown. In this work, we propose mild conditions to ensure the convergence of ADMM to a Nash point on the multi-convex problems with a sublinear convergence rate o(1/k). We also discuss how our ADMM is applied to two important applications. Extensive experiments show the effectiveness of our proposed ADMM. Our contributions in this paper include:

- We propose an ADMM framework to solve the multi-convex problem, and we investigate the convergence properties of the proposed ADMM. Specifically, we prove that the objective value and the residual are convergent. Moreover, any limit point generated by the proposed ADMM is a Nash point of the original problem. The convergence rate of the proposed ADMM is o(1/k).
- We demonstrate two important and promising applications that are special cases of our proposed ADMM framework and benefit from its theoretical properties. Specifically, we show how these applications can be transformed equivalently to fit into the proposed ADMM framework.
- We conduct extensive experiments to validate our proposed ADMM. Experiments on ten real-world datasets demonstrate its effectiveness, scalability, and convergence properties.

The rest of this paper is summarized as follows: Section 2 summarizes previous work related to this paper. Section 3 introduces the ADMM algorithm and its convergence properties. In Section 4, the proposed ADMM algorithm is applied to several important applications. Extensive experiments are described in Section 5. The paper concludes with a summary of the work in Section 6.

2 Related Work

Multi-convex optimization problems: Some works studied multi-convex problems. The earliest work required that the objective function was differentiable continuous and strictly convex [38]. Various conditions on separability and regularity on the objective functions have been discussed in [28, 29]. In the most recent work, Xu and Yin presented three types of multi-convex algorithms and analyzed convergence based on either Lipschitz differentiability or strong convexity assumption [39]. For a comprehensive survey, please refer to [26].

Convergence analysis of ADMM: Existing literature on the convergence analysis of ADMM can be categorized into two classes: the convex ADMM and the nonconvex ADMM. The convex ADMM is investigated relatively well compared with the nonconvex ADMM. Existing works either study suitable stepsizes of the convex ADMM or extend ADMM to the stochastic version. For example, Bai et al. proposed a generalized symmetric ADMM to solve the multi-block separable objective by updating the Lagrange multiplier twice with suitable stepsizes [3]; Gu et al. extended contractive

Peaceman-Rachford splitting method to ADMM with larger stepsizes [15]; Ouyang et al. proposed a stochastic ADMM with a convergence rate $O(\frac{1}{\sqrt{t}})$. Despite the outstanding performance of the nonconvex ADMM, its convergence theory is not well established due to the complexity of both coupled objectives and various (inequality and equality) constraints. Most existing works discussed the convergence of the nonconvex ADMM on separable objectives: they provided convergence guarantee to the stationary solutions with different assumptions [5, 9, 10, 20]. Some works explored more difficult cases where the objectives are coupled: for example, Wang et al. presented mild convergence conditions of the nonconvex ADMM where the objective can be nonsmooth [37]; Gao et al. explored the convergence proofs of ADMM in the nonconvex deep learning problems [31, 33, 34]; while experiments by Wang and Zhao showed that the ADMM was not necessarily convergent in the nonlinear-constrained problems [35].

3 ADMM on the Multi-convex Problems

In this section, we present an ADMM framework to solve Problem 1.

3.1 Preliminaries

First, the definition of Lipschitz differentiability is shown as follows [8]:

Definition 1 (Lipschitz Differentiability). Any arbitrary differentiable function G_1 : $\mathbb{R}^m \to \mathbb{R}$ is Lipschitz differentiable if for any $x', x'' \in \mathbb{R}^m$,

$$\|\nabla G_1(x') - \nabla G_1(x'')\| \le D \|x' - x''\|.$$

where $D \ge 0$ is constant and $\nabla G_1(x)$ denotes the gradient of $G_1(x)$.

The following defines strong convexity, which is indispensable for the proof of convergence to a Nash point.

Definition 2 (Strong Convexity). A convex function $G_4(x)$ is strongly convex if there exists H > 0 such that for $\forall x', x'' \in dom(G_4)$, the following holds

$$G_4(x^{''}) \ge G_4(x^{'}) + (v^{'})^T(x^{''} - x^{'}) + (H/2) \|x^{''} - x^{'}\|_2^2.$$

where $\forall v' \in \partial G_4(x')$ is a subdifferential of G_4 at x'.

Finally, the Nash point is defined as follows [39]:

Definition 3 (Nash Point). Given $G_5(a_1, a_2, \dots, a_m)$, a Nash point $(a_1^*, a_2^*, \dots, a_m^*)$ satisfies the following property:

$$G_5(a_1^*, \cdots, a_{i-1}^*, a_i^*, a_{i+1}^*, \cdots, a_m^*) \le G_5(a_1^*, \cdots, a_{i-1}^*, a_i, a_{i+1}^*, \cdots, a_m^*),$$

$$\forall (a_1^*, \cdots, a_{i-1}^*, a_i, a_{i+1}^*, \cdots, a_m^*) \in dom(G_5), \ (i = 1, \cdots, m).$$

Naturally, when we optimize one variable while fixing others, the Nash point ensures the optimality of this variable [39]. Without loss of generality, we assume that Problem 1 has at least a Nash point, and in the next section, we will prove that any limit point generated by ADMM converges to a Nash point.

3.2 The ADMM algorithm

The following problem is our focus in this paper:

Problem 1.

 $\min_{x_1,\dots,x_n,z} F(x_1,\dots,x_n,z) = f(x_1,\dots,x_n) + h(z) \quad s.t. \sum_{i=1}^n A_i x_i - z = 0.$ where $x_i \in \mathbb{R}^{p_i} (i = 1,\dots,n), z \in \mathbb{R}^q, f(x_1,\dots,x_n) : \mathbb{R}^p \to \mathbb{R} \cup \{\infty\} (p = \sum_{i=1}^n p_i)$ are proper, continuous, multi-convex and possibly nonsmooth functions, h(z) is a proper, differentiable and convex function. $A_i \in \mathbb{R}^{q \times p_i} (i = 1,\dots,n)$ are matrices. Obviously, the domain of F is $dom(F) = \{(x_1,\dots,x_n,z) \mid \sum_{i=1}^n A_i x_i - z = 0\}.$ Without the loss of generality, the objective of Problem 1 is assumed to be bounded from below.

To ensure the convergence of the proposed ADMM, some mild assumptions are imposed, which are shown as follows:

Assumption 1 (Lipschitz Differentiability) h(z) is Lipschitz differentiable with constant $H \ge 0$.

Most loss functions such as the cross-entropy loss and the square loss are Lipschitz differentiable [34]. In order to propose the ADMM algorithm, the augmented Lagrangian function can be formulated mathematically as follows:

$$L_{\rho}(x_1, \cdots, x_n, z, y) = F(x_1, \cdots, x_n, z) + y^T (\sum_{i=1}^n A_i x_i - z) + (\rho/2) \|\sum_{i=1}^n A_i x_i - z\|_2^2$$
(1)

where y is a dual variable and $\rho > 0$ is a penalty parameter. The proposed ADMM aims to optimize the following n + 1 subproblems alternately.

$$x_{i}^{k+1} \leftarrow \arg\min_{x_{i}} f(\cdots, x_{i-1}^{k+1}, x_{i}, x_{i+1}^{k}, \cdots) + (y^{k})^{T} A_{i} x_{i} + (\rho/2) \| \sum_{j=1}^{i-1} A_{j} x_{j}^{k+1} + A_{i} x_{i} + \sum_{j=i+1}^{n} A_{j} x_{j}^{k} - z^{k} \|_{2}^{2}.$$
(2)

$$z^{k+1} \leftarrow \arg\min_{z} L_{\rho}(\cdots, x_n^{k+1}, z, y^k)$$
(3)

$$= \arg\min_{z} h(z) - (y^{k})^{T} z + (\rho/2) \| \sum_{i=1}^{n} A_{i} x_{i}^{k+1} - z \|_{2}^{2}.$$

Algorithm 1 is presented for Problem 1. Concretely, Lines 3-5 and 6 update $x_i^{k+1}(i = 1, \dots, n)$ and z^{k+1} , respectively. Line 7 updates the primal residual r^{k+1} , which is defined in accordance with the standard ADMM [6]: it measures how the linear constraint $\sum_{i=1}^{n} A_i x_i - z = 0$ is violated. Line 8 updates the dual variable y^{k+1} , which follows the routine of the standard ADMM. Line 10 uses the norm of the primal residual r as a condition to terminate the algorithm, where $\delta > 0$ is a threshold. Each subproblem is convex and

Algorithm 1 The Proposed ADMM to Solve Problem 1

```
Require: A_i (i = 1, \cdots, n), \delta > 0.
Ensure: x_i (i = 1, \dots, n), z.
 1: Initialize \rho, k = 0.
 2:
     repeat
          peat
for i=1 to n do
x^{k+1} in Equation (2).
 3:
 4:
          end for
Update z^{k+1} in Equation (3).
 5:
6:
7:
          r^{k+1} \leftarrow \sum_{i=1}^{n} A_i x_i^{k+1} - z^{k+1}. # update primal
7. r \leftarrow \sum_{i=1}^{r} \frac{1}{2}
residual
8: y^{k+1} \leftarrow y^k + \rho r
9: k \leftarrow k+1.
10: until ||r^{k+1}|| \le \delta.
                    \leftarrow y^k + \rho r^{k+1}
11: Output x_i (i = 1, \cdots, n), z_i
```

threshold. Each subproblem is convex and implicitly assumed to be solvable.

3.3 Convergence Analysis

This section focuses on the convergence of the proposed ADMM algorithm. Specifically, the first lemma states that the augmented Lagrangian L_{ρ} keeps decreasing, which is stated as follows.

Lemma 1 (Objective Descent). If $\rho > 2H$ so that $C_1 = \rho/2 - H/2 - H^2/\rho > 0$, then there exists $C_2 = \min(\rho/2, C_1)$ such that

$$L_{\rho}(x_{1}^{k}, \cdots, x_{n}^{k}, z^{k}, y^{k}) - L_{\rho}(x_{1}^{k+1}, \cdots, x_{n}^{k+1}, z^{k+1}, y^{k+1})$$

$$\geq C_{2}(\|z^{k+1} - z^{k}\|_{2}^{2} + \sum_{i=1}^{n} \|A_{i}(x_{i}^{k+1} - x_{i}^{k})\|_{2}^{2}).$$
(4)

Lemma 1 holds under Assumption 1, and its proof can be found in Section B in the supplementary materials¹. The next lemma states that the augmented Lagrangian is bounded from below, as shown below:

Lemma 2 (Objective Bound). If $\rho > 2H$, then $L_{\rho}(x_1^k, \cdots, x_n^k, z^k, y^k)$ is lower bounded.

The proof of Lemma 2 can be found in Section B in the supplementary materials ¹. Now we can prove that the proposed ADMM converges globally in the following theorem.

Theorem 1 (Residual and Objective Convergence). If $\rho > 2H$, then for the bounded sequence $(x_1^k, \dots, x_n^k, z^k, y^k)$, then it has the following properties: a). Residual convergence. This means that as $k \to \infty$, $r^k \to 0$, where r^k is defined in Algorithm 1.

b). Objective convergence. This means that as $k \to \infty$, $F(x_1^k, \dots, x_n^k, z^k)$ converges.

Theorem 1 guarantees the convergence of the proposed ADMM, whose proof is in Section C in the supplementary materials ¹. However, $x_i^k (i = 1, \dots, n)$ and z^k are not necessarily shown to be convergent. The next theorem states that any limit point is a feasible Nash Point of Problem 1.

Theorem 2 (Convergence to a Nash Point). Let $\rho > 2H$, if either of two assumptions hold: (a). $A_i(i = 1, \dots, n)$ have full rank. (b). F is strongly convex with regard to x_i . Then for bounded variables $(x_1^k, \dots, x_n^k, z^k)$, it has at least a limit point $(x_1^*, \dots, x_n^*, z^*)$, and any limit point $(x_1^*, \dots, x_n^*, z^*)$ is a feasible Nash point of F defined in Problem 1. That is

$$\begin{split} &\sum A_{i}x_{i}^{*}-z^{*}=0. \ (feasibility) \\ &F(x_{1}^{*},\cdots,x_{n}^{*},z^{*}) \leq F(x_{1}^{*},\cdots,x_{i-1}^{*},x_{i},x_{i+1}^{*},\cdots,x_{n}^{*},z^{*}), \\ &\forall (x_{1}^{*},\cdots,x_{i-1}^{*},x_{i},x_{i+1}^{*},\cdots,x_{n}^{*},z^{*}) \in dom(F), (i=1,\cdots,n). \\ &F(x_{1}^{*},\cdots,x_{n}^{*},z^{*}) \leq F(x_{1}^{*},\cdots,x_{n}^{*},z) \forall (x_{1}^{*},\cdots,x_{n}^{*},z) \in dom(F) \ (Nash \ point) \end{split}$$

¹The supplementary materials are available at https://github.com/ xianggebenben/miADMM/blob/main/multi_convex_ADMM-13-18.pdf

The proof of Theorem 2 is detailed in Section C in the supplementary materials ¹. The third theorem states that our proposed ADMM can achieve a sublinear convergence rate of o(1/k) under Assumption 1, despite the nonconvex and complex nature of Problem 1. Such a rate is state-of-the-art even compared to those methods for simpler convex problems. The theorem is shown as follows:

Theorem 3 (Convergence Rate). If $\rho > 2H$, for a bounded sequence $(x_1^k, \dots, x_n^k, z^k, y^k)$, define $u_k = \min_{0 \le l \le k} (\|z^{l+1} - z^l\|_2^2 + \sum_{i=1}^n \|A_i(x_i^{l+1} - x_i^l)\|_2^2)$, then the convergence rate of u_k is o(1/k).

The proof of this theorem is in Section C in the supplementary materials ¹. The o(1/k) convergence rate of the proposed ADMM is consistent with much existing work analyzing the convex ADMM, including [12, 17, 22]. Our contribution in term of convergence rate is that we extend the guarantee of o(1/k) into the multi-convex Problem 1.

Our proposed ADMM is more general than some influential works in terms of formulation. The relations between our proposed ADMM and previous works are summarized as follows:

1. Generalization of Block Coordinate Descent (BCD) for multi-convex problems. When the linear constraint $\sum_{i=1}^{n} A_i x_i = z$ is removed in Problem 1, then the proposed ADMM is reduced to the Block Coordinate Descent [39].

2. Generalization of multi-block ADMM. When $f(x_1, \dots, x_n) = 0$, the proposed ADMM is reduced to the convex multi-block ADMM [27], i.e. the ADMM with no less than three variables.

Apart from general formulations, the convergence guarantees of our proposed ADMM cover more applications than previous literature. For example, [37] requires the coupled objective $f(x_1, \dots, x_n)$ to be Lipschitz differentiable. However, some important applications such as weakly-constrained multi-task learning (Section 4.1) and learning with signed-network constraints (Section 4.2) do not satisfy this condition. But they are covered by our convergence guarantees of the multi-convex ADMM to a Nash point.

4 Applications

In this section, we apply our proposed ADMM to two real-world applications, both of which conform to Problem 1 and benefit from the convergence properties of the proposed ADMM.

4.1 Weakly-constrained Multi-task Learning

In multi-task learning problems, multiple tasks are learned jointly to achieve a better performance compared with learning tasks independently [32]. Most work on multi-task learning has tended to enforce the assumption of similarity among the feature weight values across tasks [2, 11, 36, 32, 43] because this makes it possible to use convex regularization terms like $\ell_{2,1}$ norms [36] and Graph Laplacians [43]. However, this assumption is usually too strong and is seldom satisfied by the real-world data. Instead of requiring feature weights to be similar in magnitude, a more conservative but probably more reasonable assumption is that multiple tasks share similar polarities for the same feature, which means that if a feature is positively relevant to the output of a task, then its weight will also be positive for other related tasks. This assumption is appropriate for many applications. For example, the feature 'number of clinic visits' will be positively related to flu outbreaks, while the feature 'popularity of vaccination' will be negatively related to them, even though their feature weights can vary dramatically for different countries (namely tasks here). This is achieved by enforcing the requirement for every pair of tasks with neighboring indices to have the same weight sign. This optimization objective is shown as follows:

$$\min_{w_1, \cdots, w_n} \sum_{i=1}^n (Loss_i(w_i) + \Omega_i(w_i))$$
s.t. $w_{i,j} w_{i+1,j} \ge 0 \ (i = 1, 2, \cdots, n-1, j = 1, 2, \cdots, m).$
(5)

where *n* and *m* denote the number of tasks and features, respectively, $w_{i,j}$ is the weight of the *j*-th feature in the *i*-th task, w_i is the weight of the *i*-th task, and $Loss_i(w_i)$ and $\Omega_i(w_i)$ are the loss function and the regularization term of the *i*-th task, respectively. The inequality constraint implies that the *i*-th and the *i* + 1-th tasks share the same sign for their weights. Equation (5) is rewritten in the following form to fit in our proposed ADMM framework:

$$\min_{w_1, \cdots, w_n, z} \sum_{i=1}^n (Loss_i(w_i) + \Omega_i(z_i)) + \lambda_1 \sum_{i=1}^{n-1} \sum_{j=1}^m c_1(w_{i,j}w_{i+1,j})$$
(6)
s.t. $z_i = w_i \ (i = 1, 2, \cdots, n).$

where $z = [z_1; \dots; z_n]$ is an auxiliary variable, and $\lambda_1 > 0$ is a tuning parameter. Notice that the inequality constraint $w_{i,j}w_{i+1,y} \ge 0$ is transformed to a quadratic penalty $c_1(x)$ such that $c_1(x) = \begin{cases} x^2 & x < 0 \\ 0 & x \ge 0 \end{cases}$ which makes the formulation consistent with Problem 1. The proposed ADMM algorithm for this case is shown in Appendix D.1 in the supplementary materials ¹.

4.2 Learning with Signed-Network Constraints

The application of network models for social network analysis has attracted the attention of a large number of researchers [7]. For example, influential societal events often spread across many social networking sites and are expressed in different languages. Such multi-lingual indicators usually transmit similar semantic information through networks and have thus been utilized to facilitate social event forecasting [41]. The problem with network constraints is formulated as follows:

$$\min_{\beta_1, \dots, \beta_n} Loss(\beta_1, \dots, \beta_n) + \sum_{i=1}^n \omega_i(\beta_i)$$

s.t. $\exists (\beta_i, \beta_j) \in E_s, \exists (\beta_p, \beta_q) \in E_d \ (1 \le i, j, p, q \le n)$

where β_i is the weight of the *i*-th node. $Loss(\beta_1, \dots, \beta_n)$ is a loss function and $\omega_i(\beta_i)$ is a regularization term for the *i*-th node. $E_s = \{(\beta_i, \beta_j) | \beta_i \beta_j \ge 0\}$ and $E_d = \{(\beta_p, \beta_q) | \beta_p \beta_q \le 0\}$ are two edge sets to represent two opposite relationships: $(\beta_i, \beta_j) \in \{(\beta_i, \beta_j) | \beta_i \beta_j \ge 0\}$

 E_s means that $\beta_i\beta_j \ge 0$, while $(\beta_p, \beta_q) \in E_d$ means that $\beta_p\beta_q \le 0$. The constraint means that some pair (β_i, β_j) satisfies the edge set E_s , and that some pair (β_p, β_q) satisfies the edge set E_d . For example, in the problem of social event forecasting with French and English, E_s and E_d are edge sets of synonyms and antonyms between French and English, and the weight pair of the French word "bien" and the English word "good" belongs to E_s . The problem with network constraints can be reformulated approximately to the following:

$$\min_{\beta_1, \cdots, \beta_n, z} Loss(\beta_1, \cdots, \beta_n) + \sum_{i=1}^n \omega_i(z_i) + \lambda_2(\sum_{(\beta_i, \beta_j) \in E_s} c_2(\beta_i, \beta_j) + \sum_{(\beta_p, \beta_q) \in E_d} c_3(\beta_p, \beta_q)) s.t. \ z_i = \beta_i \ (i = 1, 2, \cdots, n)$$
(7)

where $z = [z_1; \dots; z_n]$ is an auxiliary variable, and $\lambda_2 > 0$ is a tuning parameter. The constraint $(\beta_i, \beta_j) \in E_s$ and $(\beta_p, \beta_q) \in E_d (1 \le i, j, p, q \le n)$ are transformed to two quadratic penalties $c_2(\beta_i, \beta_j)$ and $c_3(\beta_p, \beta_q)$ as follows:

$$c_2(\beta_i,\beta_j) = \begin{cases} (\beta_i\beta_j)^2 & (\beta_i,\beta_j) \notin E_s \\ 0 & (\beta_i,\beta_j) \in E_s \end{cases}, \\ c_3(\beta_p,\beta_q) = \begin{cases} (\beta_p\beta_q)^2 & (\beta_p,\beta_q) \notin E_d \\ 0 & (\beta_p,\beta_q) \in E_d \end{cases}$$

The proposed ADMM for this case is also shown in Appendix D.2 in the supplementary materials ¹.

5 Experiments

In this section, we test our proposed ADMM using ten real-world datasets on two applications detailed in Section 4. Scalability, effectiveness, and convergence properties are compared with several existing state-of-the-art methods on ten real datasets. All experiments were conducted on a 64-bit machine with Intel(R) Core(TM) processor (i7-6820HQ CPU@ 2.70GHZ) and 16.0GB memory.



5.1 Experiment I: Weak-constrained Multi-task Learning

To evaluate the effectiveness of our method on the application of weak-constrained multi-task learning described in Equation (6), a real-world school dataset is used to evaluate the effectiveness of our proposed ADMM. It consists of the examination scores in three years of 15,362 students from 139 secondary schools, which are treated as tasks for examination scores prediction based on 27 input features such as year of the examination, school-specific features, and student-specific features. The dataset is publicly

available and the detailed description can be found in the original paper [21]. ρ was set to 1000. Here we chose two kinds of λ_1 : (1) $\lambda_1^k = 10^5$; (2) $\lambda_1^{k+1} = \lambda_1^k + 10$ with $\lambda_1^k = 1$. $\lambda_1(1)$ and $\lambda_1(2)$ are the first and the second choice of λ_1 , respectively.

Metrics. In this experiment, five metrics were utilized to evaluate model performance. Mean Squared Error (MSE) measures the average of the squares of the difference between observation and estimation. Different from MSE, Mean Squared Logarithmic Error (MSLE) measures the ratio of observation to estimation. Mean Absolute Error (MAE) is also an error measurement but computed in the absolute value. The less the above three metrics are, the better a regression model is. Explained Variance (EV) computes the ratio of the variance of the error to that of observation. The coefficient of determination or R2 score is the proportion of the variance in the dependent variable that is predictable from the independent variable. The higher score of EV and R2 are, the better a regression model is.

Baselines. To validate the effectiveness of the proposed ADMM, five benchmark multitask learning models served as comparison methods. Loss functions were set to least square errors. The number of iterations was set to 5,000. The regularization parameter α was set based on 5-fold cross-validation on the training set.

1. multi-task learning with Joint Feature Selection (JFS) [2, 43]. JFS is one of the most commonly used strategies in multi-task learning. It captures the relatedness of multiple tasks by a constraint of a weight matrix to share a common set of features. α was set to 100.

2. Clustered Multi-Task Learning (CMTL) [42, 43]. CMTL assumes that multiple tasks are clustered into several groups. Tasks in the same group are similar to each other. α was set to 1.

3. multi-task Lasso (mtLasso) [43]. mtLasso extends the classic Lasso model to the multi-task learning setting. α was set to 10.

4. a convex relaxation of Alternating Structure Optimization (cASO) [43, 1]. cASO decomposes each task into two components: task-specific feature mapping and task-shared feature mapping. α was set to 0.01.

5. Block Coordinate Descent (BCD) [39]. BCD is an intuitive method to solve multi-convex problems, which optimizes each variable alternately. α was set to 10.

Performance. As discussed in Section 4.1, the convergence of our proposed ADMM is guaranteed based on our theoretical framework. To verify this, Figures 1(a) and 1(b) illustrate the residual and objective values in different iterations, which demonstrates the convergence of the proposed ADMM on this nonconvex problem. Then the performance of examination score prediction on this dataset is illustrated in Table 1. Table 1 shows the mean and the standard deviation of all methods, which were repeated 10 times by initializing parameters randomly, to make experimental evaluation robust. It

Lap.	le	1:	Perf	orma	ance	in	Ez	xpei	rim	ent	t I
------	----	----	------	------	------	----	----	------	-----	-----	-----

Mean									
Method	MSE	MSLE	MAE	EV	R2				
JFS	114.1052	0.4531	8.4349	0.2948	0.2948				
CMTL	114.9892	0.4647	8.4756	0.2876	0.2875				
mtLasso	115.3143	0.4625	8.4725	0.2873	0.2873				
cASO	137.8336	0.5204	9.3450	0.1606	0.1605				
BCD	149.2313	0.5577	9.8057	0.1299	0.0777				
$ADMM(\lambda_1(1))$	113.6975	0.4423	8.4024	0.2950	0.2960				
$ADMM(\lambda_1(2))$	113.2400	0.4428	8.3943	0.3002	0.3002				
Standard Deviation									
Method	MSE	MSLE	MAE	EV	R2				
JFS	2.02	0.02	0.06	0.02	0.02				
CMTL	1.85	0.02	0.05	0.01	0.01				
mtLasso	1.77	0.02	0.05	0.01	0.01				
cASO	7.26	0.01	0.06	0.01	0.01				
BCD	1.41	0.01	0.06	0.15	0.01				
$ADMM(\lambda_1(1))$	0.83	0.005	0.03	0.01	0.01				
$\text{ADMM}(\lambda_1(2))$	0.95	0.01	0.04	0.02	0.02				

shows that $\lambda_1(2)$ outperforms $\lambda_1(1)$ in four out of five metrics for the proposed ADMM. In addition, the proposed ADMM achieves the best performance in all the metrics, compared to all comparison methods. Moreover, the standard deviation of the proposed ADMM is about 30% smaller than any other comparison method. This is because our method only enforces that the sign of the feature weight across different tasks is the same, while comparison methods typically perform too aggressive assumptions on the similarity among tasks. For example, CMTL enforces that the correlated tasks need to have similar feature weights using squared regularization on the difference between feature weights. JFS and mtLasso still tend to enforce similar weights on features in different tasks by $\ell_{2,1}$ norm. Because their enforcement is weaker than CMTL, their performance is better. cASO gets relatively weak performance because it optimizes an approximation of a nonconvex problem, and thus the approximate solution may be distant from the true solution to the original problem. Finally, the BCD performs the worst among all methods, even though it shares the same formulation with our proposed ADMM. This reflects the advantage of our proposed ADMM algorithm: dual information in one iteration can be passed to its following iteration by dual variables, which yields better performance.

Scalability. To investigate the scalability of the proposed ADMM compared with all baselines in Experiment I, we measured the training time of them in the school dataset when the number of features varies. The training time was averaged by running 20 times. Figure 2 shows the training time of all methods when the number of features ranges from 10 to 28. The training time of all methods increased linearly concerning the number of features. cASO was the most efficient of all methods, while the proposed ADMM was ranked second.



Fig. 2: The training time of all methods in Experiment I.

mtLasso and JFS also trained a model within 5 seconds on average. CMTL was timeconsuming for training, which spent more than 10 seconds.

5.2 Experiment II: Event Forecasting with Multi-lingual Indicators

Datasets. To evaluate the performance of our proposed ADMM on the application in Section 4.2, extensive experiments on nine real-world datasets have been performed. The dataset is obtained by randomly sampling 10% (by volume) of the Twitter data from Jan 2013 to Dec 2014. The data in the first and second years are used and training and test set, respectively. For the topic (i.e., social unrest) of interest, 1,806 keywords in the three major languages in Latin America, namely English, Spanish, and Portuguese, were provided by the paper [41]. Their translation relationships have also been labeled as semantic links among them, such as "protest" in English, "protesta" in Spanish, and "protesto" in Portuguese. The event forecasting results were validated against a labeled event set, known as the Gold Standard Report (GSR), which is publicly available [25]. **Metric and Baselines.** The metric used to evaluate the performance is Area Under the receiver operating characteristic Curve (AUC). Five comparison methods including the state-of-the-art Multi-Task learning (MTL), Multi-Resolution Event Forecast-

ing (MREF), and Distant-supervision of Heterogeneous Multitask Learning (DHML) as well as classic methods Logistic Regression (LogReg) and Lasso. ρ was set to 10. Here we chose two kinds of λ_2 : (1) $\lambda_2^k = 10^5$; (2) $\lambda_2^{k+1} = \lambda_2^k + 10$ with $\lambda_2^k = 1$. $\lambda_2(1)$ and $\lambda_2(2)$ are the first and the second choice of λ_2 , respectively. All the hyperparameters were tuned by 5-fold cross-validation.

Table 2: Event forecasting performance in AUC in each of the 9 datasets.

Performance. As shown in Table 2, $\lambda_2(2)$ outperforms $\lambda_2(1)$ marginally in seven out of nine datasets for the proposed ADMM, and they generally perform the best among all the methods, with DHML and BCD the

	BR	CL	CO	EC	EL	MX	PY	UY	VE
LogReg	0.686	0.677	0.644	0.599	0.618	0.661	0.616	0.628	0.667
LASSO	0.685	0.677	0.648	0.603	0.636	0.665	0.615	0.666	0.669
MTL	0.722	0.669	0.810	0.617	0.772	0.795	0.600	0.811	0.771
MREF	0.714	0.563	0.515	0.784	0.612	0.693	0.658	0.681	0.588
DHML	0.845	0.683	0.846	0.839	0.780	0.793	0.737	0.835	0.835
BCD	0.847	0.668	0.850	0.830	0.773	0.800	0.736	0.835	0.856
ADMM $(\lambda_2(1))$	0.864	0.699	0.870	0.848	0.794	0.820	0.746	0.850	0.867
ADMM $(\lambda_2(2))$	0.867	0.701	0.872	0.851	0.798	0.823	0.747	0.847	0.865

second-best performer. They all outperform the others typically by at least 5%-10%. This is because they leverage the multilingual correlation among the features to boost up the model's generalizability. Thanks to the framework of multi-task learning, MTL and MREF obtained a competitive performance with AUC typically over 0.7, which outperforms simple methods like LogReg and LASSO by 5% on average.

Efficiency. In Experiment II, we also compared the training time of the proposed ADMM in comparison with all baselines on 9 datasets. The training time was averaged by running 5 times. The training time was shown in Table 3. We do not show BCD because its training time is similar to the proposed ADMM. Overall, the proposed ADMM was the most efficient of all methods for all datasets. It consumed no more than 30 seconds on all datasets. MTL

	LogReg	LASSO	MTL	MREF	DHML	ADMM
BR	30193	1535	233	25889	332	14
CL	2981	242	35	6521	852	11
СО	8060	780	108	14714	87	31
EC	312	295	17	4332	46	25
EL	551	261	17	4669	33	3
MX	17712	2043	853	31349	175	29
PY	7297	527	40	9495	242	5
UY	748	336	20	5305	82	3
VE	5563	1008	49	5769	179	28

Table 3: Comparison of running time (in seconds) on 9 datasets in Experiment II.

ranked second, but it spent hundreds of seconds on some datasets, like BR and MX. As the most time-consuming baselines, LogReg and MREF trained a model in thousands of seconds or more.

6 Conclusions

We propose an ADMM framework for multi-convex problems with multiple coupled variables. It not only inherits the merits of general ADMMs but also provides advantageous theoretical properties on convergence conditions and properties under mild conditions. Besides, several machine learning applications of recent interest are discussed as special cases of our proposed ADMM. Extensive experiments have been conducted on ten real-world datasets, and demonstrate the effectiveness, scalability, and convergence properties of our proposed ADMM.

Acknowledgement

This work was supported by the National Science Foundation (NSF) Grant No. 1755850, No. 1841520, No. 2007716, No. 2007976, No. 1942594, No. 1907805, a Jeffress Memorial Trust Award, Amazon Research Award, NVIDIA GPU Grant, and Design Knowledge Company (subcontract No: 10827.002.120.04).

Bibliography

- Rie Kubota Ando and Tong Zhang. A framework for learning predictive structures from multiple tasks and unlabeled data. Journal of Machine Learning Research, 6(Nov):1817–1853, 2005.
- Andreas Argyriou, Theodoros Evgeniou, and Massimiliano Pontil. Multi-task feature learning. Advances in neural information processing systems, pages 41-48, 2007 [2] Janchao Bai, Jicheng Li, Fenguini Xu, and Hongchao Zhang. Generalized symmetric ADMM for separable convex optimization. Computational optimization and applications, 70(1):129–170, 2018. Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. SIAM journal on imaging sciences, 2(1):183–202, 2009. [3]
- [4] [5]
- Radu loan Boş and Dang-Khoa Nguyen. The proximal alternating direction method of multipliers in the nonconvex setting: convergence analysis and rates. *Mathematics of Operations Research*, 45(2):682–712, 2020. Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. Distributed optimization and statistical learning via the alternating direction method of [6] multipliers. Foundations and Trends® in Machine Learning, 3(1):1-122, 2011.

- Pater J Carrington, John Sout, Jan Arnave M. Martine Learning, 14(1):1122-011. Peter J Carrington, John Sout, and Stanley Wasserman. Models and methods in social network analysis, volume 28. Cambridge university press, 2005. Fabio Cavalletti and Tapio Rajala. Tangent lines and lipschitz differentiability spaces. Analysis and Geometry in Metric Spaces, 4(1):85–103, 2016. MT Chao, Y Zhang, and BJ Jian. An inertial proximal alternating direction method of multipliers for nonconvex optimization. International Journal of Computer Mathematics, 98(6):1199–1217, 2021. [9] [10]
- Mathematics, 98(6):1199–1217, 2021.
 Rick Chartrand and Brendt Wohlberg. A nonconvex ADMM algorithm for group sparsity with sparse groups. Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on, pages 6009–6013, 2013.
 Jianhui Chen, Jiayu Zhou, and Jieping Ye. Integrating low-rank and group-sparse structures for robust multi-task learning. Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 42–50, 2011.
 Wei Deng, Ming-Jun Lai, Zhimin Peng, and Wotao Yin. Parallel multi-block ADMM with 0 (1/k) convergence. Journal of Scientific Computing, 71(2):712–736, 2017.
 Wenb Gao, Donald Goldfarth, and Frank E Curtis. ADMM for multiaffine constrained optimization. Optimization Methods and Software, 52(2):2573–303, 2020.
 Tom Goldstein, Brendan O'Donoghue, Simon Setzer, and Richard Baraniuk. Fast alternating direction optimization methods. SIAM Journal on Imaging Sciences, 7(3):1588–1672, 2014. [11]
- [12]
- [14] 7(3):1588-1623, 2014.
- [15] Yan Gu, Bo Jiang, and Deren Han. A semi-proximal-based strictly contractive peaceman-rachford splitting method. arXiv preprint arXiv:1506.02221, pages 1–20, 2015 Arash Hassibi, Jonathan How, and Stephen Boyd. A path-following method for solving bmi problems in control. American Control Conference, 1999. Proceedings of
- [16] the 1999, 2:1385-1389, 1999. [17]
- Bingsheng He and Xiaoming Yuan. On the o(1/n) convergence rate of the douglas-rachford alternating direction method. SIAM Journal on Numerical Analysis, 50(2):700–709, 2012. [18] Moitaba Kadkhodaie, Konstantina Christakopoulou, Maziar Saniabi, and Arindam Baneriee. Accelerated alternating direction method of multipliers. Proceedings of
- Wojado Radkibada Kakibada K [19]
- [20]
- 2015 [21] Ya Li, Xinmei Tian, Tongliang Liu, and Dacheng Tao. Multi-task model and feature joint learning. International Joint Conference on Artificial Intelligence, pages
- Adda, 3649, 2015.
 Tian-Yi Lin, Shi-Qian Ma, and Shu-Zhong Zhang. On the sublinear convergence rate of multi-block ADMM. Journal of the Operations Research Society of China, [22] 3(3):251-274, 2015
- Nelson Merentes and Kazimierz Nikodem. Remarks on strongly convex functions. Aequationes mathematicae, 80(1-2):193–199, 2010. Hua Ouyang, Niao He, Long Tran, and Alexander Gray. Stochastic alternating direction method of multipliers. International Conference on Machine Learning, pages [24] 80-88, 2013 [25] Terry Reed. Open source indicators project: https://doi.org/10.7910/DVN/EN8FUW, 2017.
- Kiryuc Shen, Steven Diamond, Madeleine Udell, Yuantao Gu, and Stephen Boyd. Disciplined multi-convex programming. Control And Decision Conference (CCDC), 2017 29th Chinese, pages 895–900, 2017. [26]
- [27] Min Tao and Xiaoming Yuan. Convergence analysis of the direct extension of ADMM for multiple-block separable convex minimization. Advances in Computational Mathematics, 44(3):773-813, 2018.
- Paul Tseng. Dail coordinate ascent methods for non-strictly convex minimization. Mathematical programming, 59(1-3):231–247, 1993.
 Paul Tseng. Convergence of a block coordinate descent method for nondifferentiable minimization. Journal of optimization theory and applications, 109(3):475–494, [28] [29]
- 2001. [30] Huahua Wang and Arindam Banerjee. Bregman alternating direction method of multipliers. Advances in Neural Information Processing Systems, pages 2816–2824,
- 2014 [31] Junxiang Wang, Zheng Chai, Yue Cheng, and Liang Zhao. Toward model parallelism for deep neural network based on gradient-free ADMM framework. 2020 IEEE
- International Conference on Data Mining (ICDM), pages 591-600, 2020. [32]
- International Conference on Data Mining (ICDM), pages 591–600, 2020. Junxiang Wang, Tyang Gao, Andreas Zulf, Jingyuan Yang, and Liang Zhao. Incomplete label uncertainty estimation for petition victory prediction with dynamic features. 2018 IEEE International Conference on Data Mining (ICDM), pages 537–546, 2018. Junxiang Wang, Hongyi Li, Zheng Chai, Yongchao Wang, Yue Cheng, and Liang Zhao. Towards quantized model parallelism for graph-augmented mlps based on gradient-free admm framework. arXiv preprint arXiv:2105.09837, 2021. Junxiang Wang, Fuxun Yu, Xiang Chen, and Liang Zhao. ADMM for efficient deep learning with global convergence. Proceedings of the 25th ACM SIGKDD International Conference on Roweledge Discovery & Data Mining, pages 111–119, 2019. [33]
- [34]
- [35]
- Control and Optimization, page 100009, 2021. Lu Wang, Yan Li, Jiayu Zhou, Dongxiao Zhu, and Jieping Ye. Multi-task survival analysis. 2017 IEEE International Conference on Data Mining (ICDM), pages 485–494, 2017. [36] [371 Yu Wang, Wotao Yin, and Jinshan Zeng. Global convergence of ADMM in nonconvex nonsmooth optimization. Journal of Scientific Computing, 78(1):29-63, 2019.
- 1381
- The training rotation in the maximum complete on transformed in the measurement of the maximum complete on participation of the Society for Industrial and Applied Mathematics, 11(3):588–593, 1963.
 Yangyang Xu and Wotao Yin. A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion. *SIAM Journal on imaging sciences*, 6(3):1758–1789, 2013. [39] Zheng Xu, Soham De, Mario Figueiredo, Christoph Studer, and Tom Goldstein. An empirical study of ADMM for nonconvex problems. NIPS 2016 Workshop on
- [40] [41]
- Energy RA, Johann Jegen Kallo Tigen Kallo Tigen Kallo Timor Marking and Kallo [42] Jiayu Zhou, Jianhui Chen, and Jieping Ye. Clustered multi-task learning via alternating structure optimization. Advances in neural information processing systems.
- ges 702-710, 2011
- Jiayu Zhou, Jianhui Chen, and Jieping Ye. Malsar: Multi-task learning via structural regularization. Arizona State University, 21, 2011. [43]

Appendix

A Preliminary Lemmas

In this section, we give preliminary lemmas which are also used in the proofs of Lemmas 1 and 2. While Lemmas 4 and 5 depend on the optimality conditions of subproblems, Lemmas 3, 6 and 7 require Assumption 1.

Lemma 3. It holds that $\forall z_1, z_2 \in \mathbb{R}^q$,

$$h(z_1) \leq h(z_2) + \nabla h(z_2)^T (z_1 - z_2) + (H/2) \|z_1 - z_2\|^2, \quad -h(z_1) \leq -h(z_2) - \nabla h(z_2)^T (z_1 - z_2) + (H/2) \|z_1 - z_2\|^2.$$

Proof. Because h(z) is Lipschitz differentiable by Assumption 1, so is -h(z). Therefore, this lemma is proven exactly as same as Lemma 2.1 in [4].

Lemma 4. It holds that $y^k = \nabla h(z^k)$ for all $k \in \mathbb{N}$.

Proof. The optimality condition for the problem with regard to z^k gives rise to

$$\nabla h(z^k) - y^{k-1} - \rho(\sum_{i=1}^n A_i x_i^k - z^k) = 0$$

Because $y^k = y^{k-1} + \rho(\sum A_i x_i^k - z^k)$, we have $y^k = \nabla h(z^k)$.

Lemma 5. It holds that for $\forall k \in \mathbb{N}$,

$$L_{\rho}(\cdots, x_{i-1}^{k+1}, x_i^k, \cdots) - L_{\rho}(\cdots, x_i^{k+1}, x_{i+1}^k, \cdots) \ge (\rho/2) \|A_i x_i^k - A_i x_i^{k+1}\|_2^2.$$
(8)

Proof.

$$\begin{split} &L_{\rho}(\cdots, x_{i-1}^{k+1}, x_{i}^{k}, \cdots) - L_{\rho}(\cdots, x_{i}^{k+1}, x_{i+1}^{k}, \cdots) \\ &= f(\cdots, x_{i-1}^{k+1}, x_{i}^{k}, \cdots) - f(\cdots, x_{i}^{k+1}, x_{i+1}^{k}, \cdots) \\ &+ (y^{k})^{T} (A_{i}x_{i}^{k} - A_{i}x_{i}^{k+1}) + (\rho/2) \|\sum_{j=1}^{i-1} A_{j}x_{j}^{k+1} + \sum_{j=i}^{n} A_{j}x_{j}^{k} - z^{k} \|_{2}^{2} - (\rho/2) \|\sum_{j=1}^{i} A_{j}x_{j}^{k+1} + \sum_{j=i+1}^{n} A_{j}x_{j}^{k} - z^{k} \|_{2}^{2} \\ &= f(\cdots, x_{i-1}^{k+1}, x_{i}^{k}, \cdots) - f(\cdots, x_{i}^{k+1}, x_{i+1}^{k}, \cdots) \\ &+ (y^{k})^{T} (A_{i}x_{i}^{k} - A_{i}x_{i}^{k+1}) + (\rho/2) \|A_{i}x_{i}^{k} - A_{i}x_{i}^{k+1} \|_{2}^{2} + \rho(\sum_{j=1}^{i} A_{j}x_{j}^{k+1} + \sum_{j=i+1}^{n} A_{j}x_{j}^{k} - z^{k})^{T} (A_{i}x_{i}^{k} - A_{i}x_{i}^{k+1}) \\ &= f(\cdots, x_{i-1}^{k+1}, x_{i}^{k}, \cdots) - f(\cdots, x_{i}^{k+1}, x_{i+1}^{k}, \cdots) \\ &+ (A_{i}^{T}y^{k} + \rho A_{i}^{T} (\sum_{j=1}^{i} A_{j}x_{j}^{k+1} + \sum_{j=i+1}^{n} A_{j}x_{j}^{k} - z^{k}))^{T} (x_{i}^{k} - x_{i}^{k+1}) + (\rho/2) \|A_{i}x_{i}^{k} - A_{i}x_{i}^{k+1}\|_{2}^{2}. \end{split}$$

where the second equality follows from the cosine rule: $||b + c||^2 - ||a + c||^2 = ||b - a||^2 + 2(a + c)^T(b - a)$ with $a = A_i x_i^{k+1}$, $b = A_i x_i^k$ and $c = \sum_{j=1}^{i-1} A_j x_j^{k+1} + \sum_{j=i+1}^n A_j x_j^k - z^k$. The optimality condition of x_i^{k+1} leads to

$$\begin{aligned} 0 &\in \partial_{x_{i}} L_{\rho}(\cdots, x_{i}^{k+1}, x_{i+1}^{k}, \cdots) \\ &= \partial_{x_{i}} f(\cdots, x_{i}^{k+1}, x_{i+1}^{k}, \cdots) + A_{i}^{T} y^{k} + \rho A_{i}^{T} (\sum_{j=1}^{i} A_{j} x_{j}^{k+1} + \sum_{j=i+1}^{n} A_{j} x_{j}^{k} - z^{k}) \\ &- A_{i}^{T} y^{k} - \rho A_{i}^{T} (\sum_{j=1}^{i} A_{j} x_{j}^{k+1} + \sum_{j=i+1}^{n} A_{j} x_{j}^{k} - z^{k}) \in \partial_{x_{i}} f(\cdots, x_{i}^{k+1}, x_{i+1}^{k}, \cdots). \end{aligned}$$

We have the following result according to the definition of subgradient

$$\begin{split} f(\cdots, x_{i-1}^{k+1}, x_i^k, \cdots) \\ &\geq f(\cdots, x_i^{k+1}, x_{i+1}^k, \cdots) + (-A_i^T y^k - \rho A_i^T (\sum_{j=1}^i A_j x_j^{k+1} + \sum_{j=i+1}^n A_j x_j^k - z^k))^T (x_i^{k+1} - x_i^k) \\ &= f(\cdots, x_i^{k+1}, x_{i+1}^k, \cdots) + (A_i^T y^k + \rho A_i^T (\sum_{j=1}^i A_j x_j^{k+1} + \sum_{j=i+1}^n A_j x_j^k - z^k))^T (x_i^k - x_i^{k+1}). \end{split}$$

Therefore, the lemma is proved.

Lemma 6. If $\rho > 2H$ so that $C_1 = \rho/2 - H/2 - H^2/\rho > 0$, then it holds that

 $L_{\rho}(\cdots, x_n^{k+1}, z^k, y^k) - L_{\rho}(\cdots, x_n^{k+1}, z^{k+1}, y^{k+1}) \ge C_1 \|z^{k+1} - z^k\|_2^2.$

$$\begin{split} & \text{Proof.} \\ & L_{\rho}(x_{1}^{k+1},\cdots,x_{n}^{k+1},z^{k},y^{k}) - L_{\rho}(x_{1}^{k+1},\cdots,x_{n}^{k+1},z^{k+1},y^{k+1}) \\ &= h(z^{k}) + (y^{k})^{T}(\sum_{i=1}^{n}A_{i}x_{i}^{k+1}-z^{k}) + (\rho/2) \|\sum_{i=1}^{n}A_{i}x_{i}^{k+1}-z^{k}\|_{2}^{2} \\ &- h(z^{k+1}) - (y^{k+1})^{T}(\sum_{i=1}^{n}A_{i}x_{i}^{k+1}-z^{k+1}) - (\rho/2) \|\sum_{i=1}^{n}A_{i}x_{i}^{k+1}-z^{k+1}\|_{2}^{2} \\ &= h(z^{k}) - h(z^{k+1}) + (y^{k}-y^{k+1})^{T}\sum_{i=1}^{n}A_{i}x_{i}^{k+1} + (y^{k+1})^{T}z^{k+1} - (y^{k})^{T}z^{k} + (\rho/2) \|\sum_{i=1}^{n}A_{i}x_{i}^{k+1}-z^{k}\|_{2}^{2} \\ &- (\rho/2) \|\sum_{i=1}^{n}A_{i}x_{i}^{k+1}-z^{k+1}\|_{2}^{2} \\ &= h(z^{k}) - h(z^{k+1}) + (y^{k}-y^{k+1})^{T}\sum_{i=1}^{n}A_{i}x_{i}^{k+1} + (y^{k+1})^{T}z^{k+1} - (y^{k})^{T}z^{k} \\ &+ (\rho/2) \|z^{k+1}-z^{k}\|_{2}^{2} + \rho(z^{k+1}-\sum_{i=1}^{n}A_{i}x_{i}^{k+1})^{T}(z^{k}-z^{k+1}) \\ &(\text{cosine rule } \|a-b\|^{2} - \|a-c\|^{2} = \|c-b\|^{2} + 2(c-a)^{T}(b-c) \text{ where } a = \sum_{i=1}^{n}A_{i}x_{i}^{k+1}, b = z^{k} \text{ and } c = z^{k+1}.) \\ &= h(z^{k}) - h(z^{k+1}) + (y^{k}-y^{k+1})^{T}\sum_{i=1}^{n}A_{i}x_{i}^{k+1} + (y^{k+1})^{T}z^{k+1} - (y^{k})^{T}z^{k} + (\rho/2) \|z^{k+1}-z^{k}\|_{2}^{2} \\ &= h(z^{k}) - h(z^{k+1}) + (y^{k}-y^{k+1})^{T}(\sum_{i=1}^{n}A_{i}x_{i}^{k+1}+z^{k}-z^{k+1}) + (y^{k}-y^{k+1})^{T}z^{k+1} - (y^{k})^{T}z^{k} + (\rho/2) \|z^{k+1}-z^{k}\|_{2}^{2} \\ &= h(z^{k}) - h(z^{k+1}) + (y^{k}-y^{k+1})^{T}(\sum_{i=1}^{n}A_{i}x_{i}^{k+1}+z^{k}-z^{k+1}) + (y^{k}-y^{k+1})^{T}z^{k+1} - (y^{k})^{T}z^{k} + (\rho/2) \|z^{k+1}-z^{k}\|_{2}^{2} \\ &= h(z^{k}) - h(z^{k+1}) + (y^{k}-y^{k+1})^{T}(y^{k}-y^{k+1}) - (y^{k+1})^{T}(z^{k}-z^{k+1}) + (y^{k}-y^{k+1})^{T}z^{k+1} - (y^{k})^{T}z^{k} + (\rho/2) \|z^{k+1}-z^{k}\|_{2}^{2} \\ &= h(z^{k}) - h(z^{k+1}) - (1/\rho)(y^{k}-y^{k+1})^{T}(y^{k}-y^{k+1}) - (y^{k+1})^{T}(z^{k}-z^{k+1}) + (p^{k}-y^{k+1})^{T}z^{k+1} - (y^{k})^{T}z^{k} + (\rho/2) \|z^{k+1}-z^{k}\|_{2}^{2} \\ &= h(z^{k}) - h(z^{k+1}) - (y^{k+1})^{T}(z^{k}-z^{k+1}) + (\rho/2) \|z^{k+1}-z^{k}\|_{2}^{2} - (1/\rho) \|y^{k+1}-y^{k}\|_{2}^{2} \\ &= h(z^{k}) - h(z^{k+1}) - (y^{k+1})^{T}(z^{k}-z^{k+1}) + (\rho/2) \|z^{k+1}-z^{k}\|_{2}^{2} - (1/\rho) \|y^{k+1}-y^{k}\|_{2}^{2} \\ &= h(z^{$$

(9)

We choose $\rho > 2H$ to make $C_1 > 0$.

Lemma 7. $\forall k \in \mathbb{N}$, we have $||y^{k+1} - y^k|| \le H ||z^{k+1} - z^k||$.

Proof.

$$\|y^{k+1} - y^k\| = \|\nabla h(z^{k+1}) - \nabla h(z^k)\|$$
 (Lemma 4) $\leq H\|z^{k+1} - z^k\|$ (Assumption 1).

B Proofs of Lemmas 1-2

Proof (Proof of Lemma 1). This follows directly from Lemmas 5 and 6.

Proof (Proof of Lemma 2). There exists z' such that $\sum_{i=1}^{n} A_i x_i^k - z' = 0$. Therefore, we have

$$f(x_1^k, \cdots, x_n^k) + h(z') \ge \min S > -\infty.$$

where $S = \{f(x_1, \dots, x_n) + h(z) : \sum_{i=1}^{n} A_i x_i - z = 0\}$, which is the objective value of Problem 1, and therefore bounded from below. Then we have

$$\begin{split} &L_{\rho}(x_{1}^{k},\cdots,x_{n}^{k},z^{k},y^{k}) \\ &= f(x_{1}^{k},\cdots,x_{n}^{k}) + h(z^{k}) + (y^{k})^{T}(\sum_{i=1}^{n}A_{i}x_{i}^{k}-z^{k}) + (\rho/2) \|\sum_{i=1}^{n}A_{i}x_{i}^{k}-z^{k}\|^{2} \\ &= f(x_{1}^{k},\cdots,x_{n}^{k}) + h(z^{k}) + (y^{k})^{T}(z^{'}-z^{k}) + (\rho/2) \|\sum_{i=1}^{n}A_{i}x_{i}^{k}-z^{k}\|^{2} \left(\sum_{i=1}^{n}A_{i}x_{i}^{k}-z^{'}=0\right) \\ &= f(x_{1}^{k},\cdots,x_{n}^{k}) + h(z^{k}) + (\nabla h(z^{k}))^{T}(z^{'}-z^{k}) + (\rho/2) \|\sum_{i=1}^{n}A_{i}x_{i}^{k}-z^{k}\|^{2} \text{ (Lemma 4)} \\ &\geq f(x_{1}^{k},\cdots,x_{n}^{k}) + h(z^{'}) + (\rho-H)/2 \|\sum_{i=1}^{n}A_{i}x_{i}^{k}-z^{k}\|_{2}^{2} \text{ (Lemma 3 and 4, h(z) is Lipschitz differentiable)} \\ &\geq \min S + (\rho-H)/2 \|\sum_{i=1}^{n}A_{i}x_{i}^{k}-z^{k}\|_{2}^{2} \geq \min S > -\infty. \end{split}$$

Therefore, $L_{\rho}(x_1^k, \cdots, x_n^k, z^k, y^k)$ is bounded from below.

Proofs of Theorems 1-3 С

Proof (Proof of Theorem 1). We show residual convergence and objective convergence based on Lemmas 1 and 2. From Lemma 1, $L_{\rho}(x_1^k, \dots, x_n^k, z^k, y^k)$ decreases monotonically, and $L_{\rho}(x_1^k, \dots, x_n^k, z^k, y^k)$ is lower bounded by Lemma 2. Therefore, $L_{\rho}(x_1^k, \dots, x_n^k, z^k, y^k)$ is convergent because a monotone bounded sequence converges (Monotone Convergence Theorem). According to the continuity of L_{ρ} , we take $k \to \infty$ on the both sides of Inequality (4) to obtain

$$\begin{split} &\lim_{k \to \infty} \left(L_{\rho}(x_1^k, \cdots, x_n^k, z^k, y^k) - L_{\rho}(x_1^{k+1}, \cdots, x_n^{k+1}, z^{k+1}, y^{k+1}) \right) \\ &\geq \lim_{k \to \infty} C_2(\|z^{k+1} - z^k\|_2^2 + \sum_{i=1}^n \|A_i(x_i^{k+1} - x_i^k)\|_2^2). \end{split}$$

On one hand, $L_{\rho}(x_1, \cdots, x_n, z, y)$ is convergent, so we have

$$\lim_{k \to \infty} C_2(\|z^{k+1} - z^k\|_2^2 + \sum_{i=1}^n \|A_i(x_i^{k+1} - x_i^k)\|_2^2) \le 0.$$

On the other hand, $C_2(\|z^{k+1} - z^k\|_2^2 + \sum_{i=1}^n \|A_i(x_i^{k+1} - x_i^k)\|_2^2$ is nonnegative, so we get

$$\lim_{k \to \infty} C_2(\|z^{k+1} - z^k\|_2^2 + \sum_{i=1}^n \|A_i(x_i^{k+1} - x_i^k)\|_2^2) = 0.$$

This suggests that $\lim_{k\to\infty} (z^{k+1}-z^k) = 0$ and $\lim_{k\to\infty} A_i(x_i^{k+1}-x_i^k) = 0$ $(i = 1, \dots, n)$. Moreover, by Lemma 7, $\lim_{k\to\infty} ||y^{k+1}-y^k|| \leq H \lim_{k\to\infty} ||z^{k+1}-z^k|| = 0$. So we have $\lim_{k\to\infty} (y^{k+1}-y^k) = 0$. a). For residual convergence, by the Line 8 of Algorithm 1, we have

$$\lim_{k \to \infty} r^k = \lim_{k \to \infty} (y^k - y^{k-1})/\rho = 0.$$

b). For objective convergence, since

$$L_{\rho}(x_1^k, \cdots, x_n^k, z^k, y^k) = F(x_1^k, \cdots, x_n^k, z^k, y^k) + (y^k)^T r^k + (\rho/2) \|r^k\|_2^2$$

and $L_{\rho}(x_1^k, \cdots, x_n^k, z^k, y^k)$ is convergent, r^k converges to 0 and y^k is bounded, then $F(x_1^k, \cdots, x_n^k, z^k, y^k)$ is also convergent

Proof (Proof of Theorem 2). Obviously $\lim_{k\to\infty} (z^{k+1} - z^k) = 0$ and $\lim_{k\to\infty} (y^{k+1} - y^k) = 0$ from the proof of Theorem 1. In order to prove this theorem, we firstly prove that $\lim_{k\to\infty} (x_i^{k+1} - x_i^k) = 0 (i = 1, \dots, n)$ if either of two assumptions holds, then prove that any limit point $(x_1^*, \dots, x_n^*, z^*)$ is a feasible Nash point of Problem 1. (a). Suppose $A_i(i = 1, \dots, n)$ have full rank. Because $\lim_{k\to\infty} A_i(x_i^{k+1} - x_i^k) = 0$ from the proof of Theorem 1, then obviously $\lim_{k\to\infty} (x_i^{k+1} - x_i^k) = 0$ [22]. (b). Suppose F is strongly convex with second to z = 0.

(b). Suppose F is strongly convex with regard to x_i . Because $L_\rho(x_1, \dots, x_n, z, y) = F(x_1, \dots, x_n, z) + y^T (\sum A_i x_i - z) + (\rho/2) \|\sum A_i x_i - z\|_2^2, F(x_1, \dots, x_n, z), \text{ and } y^T (\sum A_i x_i - z) + (\rho/2) \|\sum A_i x_i - z\|_2^2$ are strongly convex, L_ρ is also strongly convex regard to x_i [23] with the assumed constant $D_i > 0$. We have

$$\begin{split} L_{\rho}(x_{1}^{k+1},\cdots,x_{i-1}^{k+1},x_{i}^{k},x_{i+1}^{k},\cdots,x_{n}^{k},z^{k},y^{k}) &\geq L_{\rho}(x_{1}^{k+1},\cdots,x_{i-1}^{k+1},x_{i+1}^{k},x_{i+1}^{k},\cdots,x_{n}^{k},z^{k},y^{k}) + (v_{i}^{k+1})^{T}(x_{i}^{k}-x_{i}^{k+1}) \\ &+ (D_{i}/2)\|x_{i}^{k+1}-x_{i}^{k}\|_{2}^{2} \end{split}$$

where $\forall v_i^{k+1} \in \partial_{x_i} L_\rho(x_1^{k+1}, \cdots, x_{i-1}^{k+1}, x_i^{k+1}, x_{i+1}^k, \cdots, x_n^k, z^k, y^k)$. The optimality condition of x_i^{k+1} leads

$$0 \in \partial_{x_i} L_\rho(x_1^{k+1}, \cdots, x_{i-1}^{k+1}, x_i^{k+1}, x_{i+1}^k, \cdots, x_n^k, z^k, y^k).$$
 Therefore, we have

$$L_{\rho}(x_{1}^{k+1}, \cdots, x_{i-1}^{k+1}, x_{i}^{k}, x_{i+1}^{k}, \cdots, x_{n}^{k}, z^{k}, y^{k}) \geq L_{\rho}(x_{1}^{k+1}, \cdots, x_{i-1}^{k+1}, x_{i}^{k+1}, x_{i+1}^{k}, \cdots, x_{n}^{k}, z^{k}, y^{k}) + (D_{i}/2) \|x_{i}^{k+1} - x_{i}^{k}\|_{2}^{2}$$
(10)

We sum up Inequality (10) from $i = 1, \dots, n$ and Inequality (9) to obtain

$$L_{\rho}(x_1^k, \cdots, x_n^k, z^k, y^k) - L_{\rho}(x_1^{k+1}, \cdots, x_n^{k+1}, z^{k+1}, y^{k+1}) \ge \sum_{i=1}^n (D_i/2) \|x_i^{k+1} - x_i^k\|_2^2 + C_1 \|z^{k+1} - z^k\|_2^2$$
(11)

where $C_1 > 0$ by Lemma 6 if $\rho > 2H$. According to the continuity of L_{ρ} , we take $k \to \infty$ on the both sides of Inequality (11) to obtain

$$\lim_{k \to \infty} \left(L_{\rho}(x_1^k, \cdots, x_n^k, z^k, y^k) - L_{\rho}(x_1^{k+1}, \cdots, x_n^{k+1}, z^{k+1}, y^{k+1}) \right) \ge \lim_{k \to \infty} \left(\sum_{i=1}^n (D_i/2) \|x_i^{k+1} - x_i^k\|_2^2 + C_1 \|z^{k+1} - z^k\|_2^2 \right)$$

On one hand, $L_{\rho}(x_1, \cdots, x_n, z, y)$ is convergent, so we have

$$\lim_{k \to \infty} \left(\sum_{i=1}^{n} (D_i/2) \| x_i^{k+1} - x_i^k \|_2^2 + C_1 \| z^{k+1} - z^k \|_2^2 \right) \le 0$$

On the other hand, $\sum_{i=1}^{n} (D_i/2) \|x_i^{k+1} - x_i^k\|_2^2 + C_1 \|z^{k+1} - z^k\|_2^2$ is nonnegative, so we get

$$\lim_{k \to \infty} \left(\sum_{i=1}^{n} (D_i/2) \| x_i^{k+1} - x_i^k \|_2^2 + C_1 \| z^{k+1} - z^k \|_2^2 \right) = 0$$

This suggests that $\lim_{k\to\infty} (x_i^{k+1} - x_i^k) = 0$ $(i = 1, \dots, n)$ and $\lim_{k\to\infty} (z^{k+1} - z^k) = 0$. Therefore, $\lim_{k\to\infty} (x_i^{k+1} - x_i^k) = 0$ $(i = 1, \dots, n)$ if either of two assumptions holds. Because $(x_1^k, \dots, x_n^k, z^k, y^k)$ is bounded, there exists a subsequence $(x_1^s, \dots, x_n^s, z^s, y^s)$ such that $(x_1^s, \dots, x_n^s, z^s, y^s) \to (x_1^s, \dots, x_n^s, z^s, y^s)$ where $(x_1^s, \dots, x_n^s, z^s, y^s)$ is a limit point. Because $\lim_{s\to\infty} (x_i^{s+1} - x_i^s) = 0$ $(i = 1, \dots, n)$, $\lim_{s\to\infty} (z^{s+1} - z^s) = 0$ and $\lim_{s\to\infty} (y^{s+1} - y^s) = 0$, we have $(x_1^{s+1}, \dots, x_n^{s+1}, z^{s+1}, y^{s+1}) \to (x_1^s, \dots, x_n^s, z^s, y^s)$. Now we prove that the limit point $(x_1^s, \dots, x_n^s, z^s)$ is a fassible Nash point of Problem 1. For feasibility, since $\lim_{k\to\infty} r^k = \lim_{k\to\infty} \sum_{i=1}^n A_i x_i^k - z^k = 0$, so for the subsequence $(x_1^s, \dots, x_n^s, z^s, y^s) \to (x_1^s, \dots, x_n^s, z^s, y^s)$. Now For the Nash point, we obtain the following according to the optimality conditions of $x_i^{s+1}(i = 1, \dots, n)$ and z^{s+1} in Equations (2) and (3), respectively.

in Equations (2) and (3), respectively.

$$\begin{split} & L_{\rho}(x_{1}^{s+1},\cdots,x_{i-1}^{s+1},x_{i}^{s+1},x_{i+1}^{s},\cdots,x_{n}^{s},z^{s},y^{s}) \leq L_{\rho}(x_{1}^{s+1},\cdots,x_{i-1}^{s+1},x_{i},x_{i+1}^{s},\cdots,x_{n}^{s},z^{s},y^{s}), \\ & \forall (x_{1}^{s+1},\cdots,x_{i-1}^{s+1},x_{i},x_{i+1}^{s},\cdots,x_{n}^{s},z^{s}) \in dom(F) \\ & L_{\rho}(x_{1}^{s+1},\cdots,x_{n}^{s+1},z^{s+1},y^{s}) \leq L_{\rho}(x_{1}^{s+1},\cdots,x_{n}^{s+1},z,y^{s}), \ \forall (x_{1}^{s+1},\cdots,x_{n}^{s+1},z) \in dom(F) \end{split}$$

According to the continuity of L_{ρ} , we take $s \to \infty$ on the both sides of two inequalities. Because $(x_1^s, \cdots, x_n^s, z^s, y^s) \to 0$ $(x_1^*, \cdots, x_n^*, z^*, y^*)$ and $(x_1^{s+1}, \cdots, x_n^{s+1}, z^{s+1}, y^{s+1}) \to (x_1^*, \cdots, x_n^*, z^*, y^*)$, we have

$$\begin{split} & L_{\rho}(x_{1}^{*},\cdots,x_{n}^{*},z^{*},y^{*}) \leq L_{\rho}(x_{1}^{*},\cdots,x_{i-1}^{*},x_{i},x_{i+1}^{*},\cdots,x_{n}^{*},z^{*},y^{*}), \ \forall (x_{1}^{*},\cdots,x_{i-1}^{*},x_{i},x_{i+1}^{*},\cdots,x_{n}^{*},z^{*}) \in dom(F) \\ & L_{\rho}(x_{1}^{*},\cdots,x_{n}^{*},z,y^{*}) \leq L_{\rho}(x_{1}^{*},\cdots,x_{n}^{*},z^{*},y^{*}), \ \forall (x_{1}^{*},\cdots,x_{n}^{*},z) \in dom(F) \end{split}$$

Here $\forall (x_1^*, \cdots, x_{i-1}^*, x_i, x_{i+1}^*, \cdots, x_n^*, z^*) \in dom(F)$ and $\forall (x_1^*, \cdots, x_n^*, z) \in dom(F)$ mean $\forall x_i \ s.t. \sum_{j=1, j \neq i}^n A_j x_j^* + A_i x_i - z^* = 0$ and $\forall z \ s.t. \sum_{j=1}^n A_j x_j^* - z = 0$, respectively. Using the fact that $(x_1^*, \cdots, x_n^*, z^*)$ is feasible in Problem 1, we obtain $L_\rho(x_1^*, \cdots, z^*, y^*) = F(x_1^*, \cdots, z^*)$, $L_\rho(x_1^*, \cdots, x_{i-1}^*, x_i, x_{i+1}^*, \cdots, x_n^*, z^*, y^*) = F(x_1^*, \cdots, x_n^*, z, y^*)$ is a feasible Nash point of F defined in Problem 1.

Proof (*Proof of Theorem 3*). To prove this theorem, we will first show that u_k satisfies two conditions: (1). $u_k \ge u_{k+1}$. (2). $\sum_{k=0}^{\infty} u_k$ is bounded. We then conclude the convergence rate of o(1/k) based on these two conditions. Specifically, first, we have

$$u_{k} = \min_{0 \le l \le k} (\|z^{l+1} - z^{l}\|_{2}^{2} + \sum_{i=1}^{n} \|A_{i}(x_{i}^{l+1} - x_{i}^{l})\|_{2}^{2})$$

$$\geq \min_{0 \le l \le k+1} (\|z^{l+1} - z^{l}\|_{2}^{2} + \sum_{i=1}^{n} \|A_{i}(x_{i}^{l+1} - x_{i}^{l})\|_{2}^{2})$$

$$= u_{k+1}$$

Therefore u_k satisfies the first condition. Second,

$$\begin{split} &\sum_{k=0}^{\infty} u_k \\ &= \sum_{k=0}^{\infty} \min_{0 \le l \le k} (\|z^{l+1} - z^l\|_2^2 + \sum_{i=1}^n \|A_i(x_i^{l+1} - x_i^l)\|_2^2) \\ &\le \sum_{k=0}^{\infty} (\|z^{k+1} - z^k\|_2^2 + \sum_{i=1}^n \|A_i(x_i^{k+1} - x_i^k)\|_2^2) \\ &\le (L_{\rho}(x_1^0, \cdots, x_n^0, z^0, y^0) - L_{\rho}^*)/C_2 \text{ (Lemma I)} \end{split}$$

where $L_{\rho}^* = \lim_{k \to \infty} L_{\rho}(x_1^k, \dots, x_n^k, z^k, y^k)$. So $\sum_{k=0}^{\infty} u_k$ is bounded and u_k satisfies the second condition. Finally, it has been proved that the sufficient conditions of convergence rate o(1/k) are: (1) $u_k \ge u_{k+1}$, and (2) $\sum_{k=0}^{\infty} u_k$ is bounded, and (3) $u_k \ge 0$ (Lemma 1.2 in [12]). Since we have proved the first two conditions and the third one $u_k \ge 0$ is obvious, the convergence rate of o(1/k) is proven.

D Algorithms for Applications

D.1 Weakly-constrained Multi-task Learning

Applying the proposed ADMM to solve the problem in Equation (6), we get Algorithm 2. Specifically, Lines 4-9 update primal variables $w_i (i = 1, \dots, n)$ and $z_i (i = 1, \dots, n)$, Line 10 updates the dual variable $y_i (i = 1, \dots, n)$.

Algorithm 2 The Proposed ADMM to Solve Equation (6).

1: Denote $z = [z_1; \dots; z_n], y = [y_1; \dots; y_n].$ 2: Initialize $\rho, k = 0.$ 3: repeat 4: Update w_1^{k+1} by Equation (12). 5: for i=2 to n-1 do 6: Update w_i^{k+1} by Equation (13). 7: end for 8: Update w_n^{k+1} by Equation (14). 9: Update z_i^{k+1} by Equation (15) in parallel. 10: $y_i^{k+1} \leftarrow y_i^k + \rho(w_i^{k+1} - z_i^{k+1})$ $(i = 1, \dots, n)$ in parallel. 11: $k \leftarrow k + 1.$ 12: until convergence. 13: Output $w_i(i = 1, \dots, n), z.$

All subproblems are detailed as follows: **1. Update** w^{k+1}

The w_i^{k+1} $(i = 1, \cdots, n)$ are updated as follows:

$$\begin{split} w_{1}^{k+1} &\leftarrow \arg\min_{w_{1}} Loss_{1}(w_{1}) + \lambda_{1} \sum_{j=1}^{m} c_{1}(w_{1,j}w_{2,j}^{k}) + (\rho/2) \|w_{1} - z_{1}^{k} + y_{1}^{k}/\rho\|_{2}^{2}. \end{split}$$
(12)
$$w_{i}^{k+1} &\leftarrow \arg\min_{w_{i}} Loss_{i}(w_{i}) + \lambda_{1} \sum_{j=1}^{m} c_{1}(w_{i,j}w_{i+1,j}^{k}) + \lambda_{1} \sum_{j=1}^{m} c_{1}(w_{i-1,j}^{k+1}w_{i,j}) + (\rho/2) \|w_{i} - z_{i}^{k} + y_{i}^{k}/\rho\|_{2}^{2}. \end{aligned}$$
(13)
$$w_{n}^{k+1} \leftarrow \arg\min_{w_{n}} Loss_{n}(w_{n}) + \lambda_{1} \sum_{j=1}^{m} c_{1}(w_{n-1,j}^{k+1}w_{n,j}) + (\rho/2) \|w_{n} - z_{n}^{k} + y_{n}^{k}/\rho\|_{2}^{2}. \end{aligned}$$
(14)

They can be solved by the Iterative Soft Thresholding Algorithm (ISTA) [4]. Take w_1^{k+1} as an example, The ISTA leads to

$$w_1^{t+1} \leftarrow \lambda_1 \sum_{j=1}^m c_1(w_{1,j}w_{2,j}^k) + 1/(2\eta) \|w_1 - (w_1^t - \eta \nabla \phi(w_1^t))\|_2^2$$

where w_1^t is the *t*-th iteration in the ISTA, $\eta > 0$ is a learning rate, $\phi(w_1^t) = Loss_1(w_1^t) + \rho/2||w_1^t - z_1^k + y_1^k/\rho||_2^2$. For each entry of $w_{1,j}^{t+1}(j = 1, 2, \cdots, m)$, we have the following closed-form solution as follows: 1). If $w_{1,j}^{t+1}w_{2,j}^k \leq 0$, then $w_{1,j}^{t+1} \leftarrow (w_{1,j}^t - \eta \nabla_j \phi(w_1^t))/(2\eta \lambda_1 w_{2,j}^2 + 1)$. 2). If $w_{1,j}^{t+1}w_{2,j}^k \geq 0$, then $w_{1,j}^{t+1} \leftarrow w_{1,j}^t - \eta \nabla_j \phi(w_1^t)$. where $\nabla_j \phi(w_1^t)$ is the *j*-th entry of $\nabla \phi(w_1^t)$. 2. Update z^{k+1} The z_i^{k+1} ($i = 1, \cdots, n$) are updated as follows:

$$z_i^{k+1} \leftarrow \arg\min_{z_i} \Omega_i(z_i) + (\rho/2) \|w_i^{k+1} - z_i + y_i^k/\rho\|_2^2 (i = 1, \cdots, n).$$
(15)

For ℓ_1 or ℓ_2 regularization, they have closed-form solutions.

D.2 Learning with Signed-Network Constraints

Applying proposed ADMM to solve the problem in Equation (7), we get Algorithm 3. Specifically, Lines 4-7 update primal variables $\beta_i (i = 1, \dots, n)$ and z, Line 8 updates the dual variable y. All subproblems are detailed as follows:

Algorithm 3 The Proposed ADMM to Solve Equation (7).

1: Denote $z = [z_1; \dots; z_n], y = [y_1; \dots; y_n].$ 2: Initialize $\rho, k = 0.$ 3: repeat 4: for i=1 to n do 5: Update β_i^{k+1} by Equation (16). 6: end for 7: Update z_i^{k+1} ($i = 1, \dots, n$) by Equation (17) in parallel. 8: $y_i^{k+1} \leftarrow y_i^k + \rho(\beta_i^{k+1} - z_i^{k+1})$ ($i = 1, \dots, n$) in parallel. 9: $k \leftarrow k+1$. 10: **until** convergence. 11: Output $\beta_i (i = 1, \dots, n), z$.

1. Update β^{k+1} The β_i^{k+1} $(i = 1, \cdots, n)$ are updated as follows: β_i^{k+1}

$$\begin{split} & \overset{k+1}{\leftarrow} \arg\min_{\beta_{i}} Loss(\cdots, \beta_{i-1}^{k+1}, \beta_{i}, \beta_{i+1}^{k}, \cdots) + \lambda_{2} (\sum_{(\beta_{i}, \beta_{j}^{k+1}) \in E_{s,j} < i} c_{2}(\beta_{i}, \beta_{j}^{k+1}) + \sum_{(\beta_{i}, \beta_{j}^{k}) \in E_{s,j} > i} c_{2}(\beta_{i}, \beta_{j}^{k}) \\ & + \sum_{(\beta_{i}, \beta_{q}^{k+1}) \in E_{d}, q < i} c_{3}(\beta_{i}, \beta_{q}^{k+1}) + \sum_{(\beta_{i}, \beta_{q}^{k}) \in E_{d}, q > i} c_{3}(\beta_{i}, \beta_{q}^{k}) + (\rho/2) \|\beta_{i} - z_{i}^{k} + y_{i}^{k}/\rho\|_{2}^{2}). \end{split}$$

$$(16)$$

Similar to updating w_i^{k+1} in Algorithm 2, they can be solved efficiently by the ISTA [4].

2. Update z_i^{k+1} The z_i^{k+1} $(i = 1, \dots, n)$ are updated as follows:

$$z_i^{k+1} \leftarrow \arg\min_{z_i} \omega_i(z_i) + (\rho/2) \|\beta_i^{k+1} - z_i + y_i^k/\rho\|_2^2 (i = 1, \cdots, n).$$
(17)

Similar to updating z_i^{k+1} in Algorithm 2, they usually have closed-form solutions.