

---

# Differential Inclusions for Modeling Nonsmooth ADMM Variants: A Continuous Limit Theory

---

Anonymous Authors<sup>1</sup>

## Abstract

Recently, there has been a great deal of research attention on understanding the convergence behavior of first-order methods. One line of this research focuses on analyzing the convergence behavior of first-order methods using tools from continuous dynamical systems such as ordinary differential equation and differential inclusions. These research results shed lights on better understanding first-order methods from a non-optimization point of view. The alternating direction method of multipliers (ADMM) is a widely used first-order method for solving optimization problems arising from machine learning and statistics, and it is important to investigate its behavior using these new techniques from dynamical systems. Existing works along this line have been mainly focusing on problems with smooth objective functions, which excludes many important applications that are traditionally solved by ADMM variants. In this paper, we analyze some well-known and widely used ADMM variants for nonsmooth optimization problems using the tools of differential inclusions. In particular, we analyze the convergence behavior of linearized ADMM and gradient-based ADMM for nonsmooth problems and show their connections with dynamical systems. We anticipate that these results will provide new insights on understanding ADMM for solving nonsmooth problems.

## 1. Introduction

Recently, there has been tremendous interests in using continuous-time dynamical system tools to analyze first-order optimization algorithms such as Nesterov's accelerated gradient method (AGM) (Nesterov, 1983) and its

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

variants. In the seminal work Su et al. (2016), the authors designed a differential equation for modeling AGM, and analyzed the connection between the solution of the differential equation and the continuous limit of the iterates of AGM. This work provided new insights on understanding the convergence behavior of AGM. Later investigations along this line mainly focused on analyzing AGM and its variants such as FISTA and heavy ball method using the tools of differential equation, differential inclusion, and more generally, continuous dynamical systems (see, e.g., Wibisono et al. (2016); Wilson et al. (2016); Krichene et al. (2015); Shi et al. (2018)). Very recently, França et al. (2018a;b) made a significant step towards understanding the alternating direction method of multipliers (ADMM) using the tools from continuous dynamical systems. ADMM is now a widely used algorithm for solving problems with separable structures, which include a lot of important applications arising from image processing, signal processing, machine learning, statistics etc. It has a close connection with some classical operator-splitting methods in numerical PDEs such as Douglas-Rachford (Douglas & Rachford, 1956) and Peaceman-Rachford (Peaceman & Rachford, 1955) operator-splitting methods that dated back to the 1950s. These operator-splitting methods were later studied in Gabay & Mercier (1975); Glowinski & Marroco (1975); Gabay (1983); Fortin & Glowinski (1983); Glowinski & Le Tallec (1989); Eckstein & Bertsekas (1992). But the renaissance of ADMM was due to several works in 2007-2008 that introduced this algorithm to solving signal processing and image processing problems (Combettes & Pesquet, 2007; Goldstein & Osher, 2009; Wang et al., 2008). Since then, ADMM was successfully used for solving important applications in many areas in science and engineering. The popularity and importance of ADMM has been partly demonstrated by the recognition of the highly influential survey paper Boyd et al. (2011). As a result, the works of França et al. (França et al., 2018a;b) are very timely and important as they provided new tools for further understanding the convergence behavior of this influential algorithm.

However, one major drawback of França et al. (2018a;b) is that they assume that the objective function is smooth, which in fact rules out most of the applications solved by ADMM and its variants. More specifically, França et al.

(2018a;b) consider the following problem:

$$\min_{x \in \mathbb{R}^d} f(x) + g(Ax), \quad (1)$$

where  $A \in \mathbb{R}^{m \times d}$ ,  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $g: \mathbb{R}^m \rightarrow \mathbb{R}$ , and both  $f$  and  $g$  are convex and differentiable. We need to point out that under such assumption, there exist many other efficient algorithms for solving (1) and ADMM may not be a good choice.

In this paper, we allow  $f$  and  $g$  to be nonsmooth functions in (1). To apply ADMM, one standard technique is to rewrite (1) as

$$\begin{aligned} \min_{x, z} f(x) + g(z) \\ \text{s.t. } Ax - z = 0. \end{aligned} \quad (2)$$

One typical iteration of ADMM for solving (2) is

$$x^{k+1} := \operatorname{argmin}_x \left[ f(x) + \frac{\rho}{2} \|Ax - z^k + u^k\|^2 \right], \quad (3a)$$

$$z^{k+1} := \operatorname{argmin}_z \left[ g(z) + \frac{\rho}{2} \|Ax^{k+1} - z + u^k\|^2 \right], \quad (3b)$$

$$u^{k+1} := u^k + Ax^{k+1} - z^{k+1}, \quad (3c)$$

where  $u$  is the (scaled) Lagrange multiplier and  $\rho > 0$  is a penalty parameter in the augmented Lagrangian function:

$$\mathcal{L}_\rho(x, z, u) = f(x) + g(z) + \rho \langle u, Ax - z \rangle + \frac{\rho}{2} \|Ax - z\|^2. \quad (4)$$

Throughout this paper, we assume that both  $f$  and  $g$  are proper closed convex functions.

**Our Contributions.** We extend the analysis of França et al. (2018a;b) to nonsmooth cases using a completely different tool: differential inclusion, which is motivated by the analysis of FISTA by Vassilis et al. (2018). More specifically, we analyze the convergence rate of continuous limit of two widely used ADMM variants for nonsmooth problems: linearized ADMM and gradient-based ADMM. We anticipate that these results will provide new insights on understanding ADMM for solving nonsmooth problems.

## 2. Continuous Limit of Linearized ADMM and Gradient-Based ADMM

We need to point out that the ADMM given in (3) is rarely used in practice, because for most applications, the  $x$ -subproblem does not have closed-form solution and an iterative solver is still needed to solve it. Note that although the  $z$ -subproblem in (3) corresponds to the proximal mapping of function  $g$ , the  $x$ -subproblem does not correspond to the proximal mapping of  $f$  because of the presence of matrix  $A$ . Moreover, it is possible that in some applications,  $f$  does not have an easy proximal mapping. Two most commonly

used nonsmooth ADMM variants in practice are linearized ADMM and gradient-based ADMM, and they are suitable for the following two cases, which cover most applications of ADMM:

- Case (i): *Linearized ADMM* is suitable for the case where  $f$  is nonsmooth with easy proximal mappings; one representative application in this case is the Lasso problem where  $f(x) = \|x\|_1$  and  $g(\cdot) = \frac{1}{2} \|\cdot - b\|^2$  (Tibshirani, 1996a).
- Case (ii): *Gradient-based ADMM* is suitable for the case where  $g$  is nonsmooth with easy proximal mapping,  $f$  is differentiable but does not have an easy proximal mapping; one representative application in this case is the sparse logistic regression problem where  $g$  is the  $\ell_1$  norm and  $f$  is the logistic loss function (Liu et al., 2009). Note that  $A = I$  in this particular application.

We now provide more details about the applicability of linearized ADMM and gradient-based ADMM.

In Case (i), where both  $f$  and  $g$  are nonsmooth with easy proximal mappings, the  $z$ -subproblem (3b) corresponds to the proximal mapping of  $g$  and is thus easy to solve; while the  $x$ -subproblem (3a) does not correspond to the proximal mapping of  $f$  due to the presence of matrix  $A$ . Linearized ADMM addresses this issue by adding a suitably chosen proximal term  $\frac{1}{2} \|x - x^k\|_{\tau_L I - \rho A^\top A}$  to the objective function of (3a), which results in the following subproblem whose solution corresponds to the proximal mapping of  $f$ :

$$\min_x f(x) + \frac{\tau_L}{2} \left\| x - \left( x_k - \frac{\rho}{\tau_L} A^\top (Ax_k - z_k + u_k) \right) \right\|^2,$$

where  $1/\tau_L$  can be viewed as the step size of the gradient step of the quadratic penalty. Combining this subproblem with (3b) and (3c) yields the linearized ADMM. Here we consider a slightly more general version of linearized ADMM by adding a relaxation term to the intermediate residual  $Ax_{k+1} - z_k$ , which is summarized in (5):

$$x_{k+1} = \operatorname{argmin}_x \left\{ f(x) + \frac{\tau_L}{2} \left\| x - \left( x_k - \frac{\rho}{\tau_L} A^\top (Ax_k - z_k + u_k) \right) \right\|^2 \right\}, \quad (5a)$$

$$z_{k+1} = \operatorname{argmin}_z \left\{ g(z) + \frac{\rho}{2} \left\| \alpha Ax_{k+1} + (1 - \alpha)z_k - z + u_k \right\|^2 \right\}, \quad (5b)$$

$$u_{k+1} = u_k + (\alpha Ax_{k+1} + (1 - \alpha)z_k - z_{k+1}), \quad (5c)$$

where  $\alpha \in (0, 2)$  is a relaxation parameter, and when  $\alpha = 1$ , it reduces to the classical linearized ADMM. As a widely used nonsmooth ADMM variant, linearized ADMM

has been studied by many researchers, (see, e.g., Chen & Teboulle (1994); Eckstein (1994); He et al. (2002); Zhang et al. (2010); Yang & Zhang (2011); Lin et al. (2011); Ma (2016); Xu (2015); Yang & Yuan (2013); Ouyang et al. (2015)). The difference between (3b)-(3c) and (5b)-(5c) is that  $Ax_{k+1}$  is replaced by  $\alpha Ax_{k+1} + (1 - \alpha)z_k$ . This is called relaxation, which has been suggested in many papers (see, e.g., Eckstein & Bertsekas (1992)) to provide more flexibility and potentially improve the convergence speed of the algorithm.

We now use the total variation minimization problem Rudin et al. (1992) as an example to show how (5) works for a particular problem. The total variation minimization problem can be casted as the following form after variable splitting,

$$\begin{aligned} & \text{minimize} && \frac{1}{2}\|x - b\|_2^2 + \lambda\|z\|_1 \\ & \text{subject to} && z = Dx, \end{aligned}$$

which is in the form of (1) with  $A = D$ ,  $f(\cdot) = \frac{1}{2}\|\cdot - b\|_2^2$  and  $g(\cdot) = \lambda\|\cdot\|_1$ . When linearized ADMM (5) is applied to solve (2), the two subproblems (5a) and (5b) respectively correspond to the proximal mappings of  $\|\cdot\|_1$  and  $\frac{1}{2}\|\cdot - b\|_2^2$ , which are both very easy to compute.

In Case (ii), gradient-based ADMM is suitable for the case where  $f$  does not have an easy proximal mapping. In this case, the linearized ADMM (5) is not a good choice, because the  $x$ -subproblem (5a) is still not easy to solve. As a result, the gradient-based ADMM is proposed to address this issue. A typical iteration of gradient-based ADMM is as follows.

$$x_{k+1} = x_k - (1/\tau_G)(\nabla f(x_k) + \rho A^\top (Ax_k - z_k + u_k)), \quad (6a)$$

$$z_{k+1} = \underset{z}{\operatorname{argmin}} \left\{ g(z) + \frac{\rho}{2} \|\alpha Ax_{k+1} + (1 - \alpha)z_k - z + u_k\|_2^2 \right\}, \quad (6b)$$

$$u_{k+1} = u_k + (\alpha Ax_{k+1} + (1 - \alpha)z_k - z_{k+1}), \quad (6c)$$

where we again used a relaxation term to make the algorithm more general, and  $1/\tau_G$  is the step size of the negative gradient step taken on the augmented Lagrangian function  $\mathcal{L}_\rho(x^k, z^k, u^k)$ . Note that since we assume that  $f$  is differentiable in this case, the objective function of the  $x$ -subproblem in (3a) becomes a differentiable function. As a result, gradient-based ADMM suggests that a gradient step is taken instead of minimizing the augmented Lagrangian function directly, which results in the new  $x$ -subproblem in (6a). This gradient-based ADMM has been studied in the literature extensively, (see, e.g., Condat (2013); Vu (2013); Davis & Yin (2017); Lin et al. (2017)).

We now use sparse logistic regression Koh et al. (2007) as an example to show how (6) works for a particular problem.

The sparse logistic regression problem can be casted as

$$\min_{x,v} \frac{1}{p} \sum_{i=1}^p \log(1 + \exp(-b_i(a_i^\top x + v))) + \lambda\|x\|_1, \quad (7)$$

which is in the form of (1) with  $f$  being the logistic loss function,  $g$  being the  $\ell_1$  norm and  $A = I$ . Note that the logistic loss function  $f$  does not adopt an easy proximal mapping, but it is differential and thus a gradient step can be taken for the  $x$ -subproblem. When the gradient-based ADMM (6) is applied to solve (7), the two subproblems (6a) and (6b) are both easy to be implemented.

### 3. Main Results

In this section, we present the main results on the convergence of the continuous limit of the iterates of linearized ADMM (5) and gradient-based ADMM (6). We focus on the continuous approximation results when  $t = \rho^{-1}k$  and  $\rho \rightarrow \infty$ , where  $t$  denotes the time and  $k$  is the iteration counter. The following assumption is needed for the results in this section.

**Assumption 1.** We assume that  $F(x) := f(x) + g(Ax)$  is a lower semi-continuous, sub-differentiable function properly defined over  $\mathbb{R}^d$ , where

(i)  $F$  is lower semi-continuous means that

$$\liminf_{y \rightarrow x} F(y) \geq F(x), \forall x \in \mathbb{R}^d.$$

(ii)  $F$  is sub-differentiable means that,  $\forall x, y \in \mathbb{R}^d$ ,  $\exists g \in \mathbb{R}^d$  satisfying

$$F(x) \geq F(y) + g^\top(x - y).$$

(iii) We assume that matrix  $A$  has full column rank and singular values  $\sigma_1 \geq \dots \geq \sigma_d > 0$ .

For linearized ADMM, we have the following approximation theorem.

**Theorem 2.** Let Assumption 1 hold, and the relaxation parameter  $\alpha \in (0, 2)$ . Rescaling the time by setting  $t = \rho^{-1}k$ . Assume  $\tau_L$  is dependent on  $\rho$ , and  $\tau_L/\rho \rightarrow c$  when  $\rho \rightarrow \infty$ , where  $c$  is a constant. The continuous-time limit of the iterates  $\{x_k\}$  of linearized ADMM (5) as  $\rho \rightarrow \infty$  is given by the differential inclusion

$$0 \in \partial F(X(t)) + \left( cI + \frac{1 - \alpha}{\alpha} A^\top A \right) \dot{X}(t), \quad (8)$$

with initial value  $X(0) = x_0$ .

Similarly, we have the following Theorem 3 for the continuous limit of iterates of gradient-based ADMM (6).

**Theorem 3.** Let Assumption 1 hold, and the relaxation parameter  $\alpha \in (0, 2)$ . In addition, we assume that  $f$  is smooth. Rescaling the time by setting  $t = \rho^{-1}k$ . Assume  $\tau_G$  is dependent on  $\rho$ , and  $\tau_G/\rho \rightarrow c$  when  $\rho \rightarrow \infty$ , where  $c$  is a constant. continuous-time limit of the iterates  $\{x_k\}$  of gradient-based ADMM (6) as  $\rho \rightarrow \infty$  is given by the same differential inclusion (8) with initial value  $X(0) = x_0$  as linearized ADMM.

**Remark 4.** Here we assume mildly that the solution of the differential inclusion (8) exists and is unique. For conditions that guarantee the existence and uniqueness of differential inclusion, we refer the readers to Adly et al. (2006); Attouch et al. (2002); Paoli (2000).

The next theorem shows the convergence rate of the continuous-time limit of the iterates  $\{x_k\}$  generated by linearized ADMM (5) and gradient-based ADMM (6) in the convex case. We use  $x^*$  to denote the minimizer of  $F$ . We recall that  $\sigma_1$  is the largest singular value of  $A$  and  $\sigma_d$  is the smallest singular value of  $A$ . For simple of notations, we define  $\kappa_1 = \sqrt{c + \frac{1-\alpha}{\alpha}\sigma_1^2}$  and  $\kappa_d = \sqrt{c + \frac{1-\alpha}{\alpha}\sigma_d^2}$  to be the largest and smallest singular value of  $cI + \frac{1-\alpha}{\alpha}A^T A$  correspondingly.

**Theorem 5.** Let Assumption 1 hold. Assume that  $c$  and  $\alpha$  are chosen such that the matrix  $(cI + \frac{1-\alpha}{\alpha}A^T A)$  is positive definite, with largest and smallest eigenvalues being  $\kappa_1^2, \kappa_d^2$ . The shock solution  $X(t)$  of the differential inclusion (8) with initial value  $X(0) = x_0$  obtained in Theorems 2 and 3 has bounded trajectory and  $\mathcal{O}(t^{-1})$  convergence rate almost everywhere, i.e., for a.e.  $t > 0$  it holds that

$$\|X(t) - x^*\| \leq \frac{\kappa_1}{\kappa_d} \|x_0 - x^*\|,$$

$$F(X(t)) - F(x^*) \leq \frac{\kappa_1^2 \|x_0 - x^*\|^2}{2t}.$$

Moreover,

$$\int_0^{+\infty} [F(X(t)) - F(x^*)] \leq \frac{\kappa_1^2}{2} \|x_0 - x^*\|^2,$$

$$\int_0^{+\infty} t \|\dot{X}(t)\|^2 \leq \frac{\kappa_1^2}{2\kappa_d^2} \|x_0 - x^*\|^2.$$

**Remark 6.** The concept of shock solution is given in Attouch et al. (2002); Paoli (2000); Vassilis et al. (2018). We include its definition and existence result in supplementary materials for completeness.

**Remark 7.** Theorem 5 proves an  $\mathcal{O}(t^{-1})$  convergence rate in function value gap when the objective function  $F$  is convex.

## 4. Proofs of Main Results

In this section, we prove the main results provided in Section 3. Section 4.1 and Section 4.2 prove Theorems 2 and 3,

respectively. We defer to Appendix A the proof of Theorem 5 due to limitation of space.

### 4.1. Proof of Theorem 2

Due to the strong convexity of the optimization subproblems (5a) and (5b), it is easy to verify that the sequence  $\{x_k, z_k, u_k\}$  is unique. We have from the first-order optimality conditions of (5a) and (5b) that

$$0 \in \partial f(x_{k+1}) + \tau_L \left[ x_{k+1} - \left( x_k - \frac{\rho}{\tau_L} A^T (Ax_k - z_k + u_k) \right) \right], \quad (9a)$$

$$0 \in \frac{1}{\rho} \partial g(z_{k+1}) - (\alpha Ax_{k+1} + (1 - \alpha)z_k - z_{k+1} + u_k). \quad (9b)$$

We detail the proof in the following:

- (i) Adding up (9b) and (5c) eliminates the common term  $(\alpha Ax_{k+1} + (1 - \alpha)z_k - z_{k+1} + u_k)$  and reduces to a simple  $u$ -update:

$$u_{k+1} \in \frac{1}{\rho} \partial g(z_{k+1}). \quad (10)$$

Taking the continuous limit  $\rho \rightarrow \infty$  gives  $U(t) = 0$ , and hence  $\dot{U}(t) = 0$ .<sup>1</sup>

- (ii) Reorganize (9a) into the following form:

$$0 \in \partial f(x_{k+1}) + \tau_L(x_{k+1} - x_k) + \rho A^T (Ax_k - z_k + u_k). \quad (11)$$

Bringing (10) into (11) leads to:

$$0 \in \partial f(x_{k+1}) + A^T \partial g(z_k) + \tau_L(x_{k+1} - x_k) + \rho A^T (Ax_k - z_k), \quad (12)$$

where again from (5c),

$$u_{k+1} - u_k = \alpha Ax_{k+1} + (1 - \alpha)z_k - z_{k+1} = \alpha A(x_{k+1} - x_k) - (z_{k+1} - z_k) + \alpha(Ax_k - z_k),$$

and hence

$$Ax_k - z_k = \frac{1}{\alpha} [(u_{k+1} - u_k) + (z_{k+1} - z_k)] - A(x_{k+1} - x_k). \quad (13)$$

<sup>1</sup>Although the continuous version of  $U(t)$  is constantly zero, it is different with  $u_k = 0$ . One may regard  $u_k$  as an infinitesimal number that dynamically changes in the system.



Plugging (13) into (12) gives

$$0 \in \partial f(x_{k+1}) + A^\top \partial g(z_k) + \tau_L(x_{k+1} - x_k) + \rho A^\top \left( \frac{1}{\alpha} [(u_{k+1} - u_k) + (z_{k+1} - z_k)] - A(x_{k+1} - x_k) \right). \quad (14)$$

Taking the limit  $\rho \rightarrow \infty$  and letting  $\tau_L/\rho \rightarrow c$ , using the fact that  $\dot{U}(t) = 0$ , (14) reduces to

$$0 \in \partial f(X(t)) + A^\top \partial g(Z(t)) + (cI - A^\top A) \dot{X}(t) + \frac{1}{\alpha} A^\top \dot{Z}(t). \quad (15)$$

(iii) We directly take the  $\rho \rightarrow \infty$  limit in (5c) and conclude

$$Z(t) = AX(t), \quad \dot{Z}(t) = A\dot{X}(t).$$

Combining the above and (15) concludes

$$0 \in \partial F(X(t)) + \left( cI + \frac{1-\alpha}{\alpha} A^\top A \right) \dot{X}(t),$$

which completes the proof.

## 4.2. Proof of Theorem 3

Again the sequence  $\{x_k, z_k, u_k\}$  is unique due to the strong convexity of the optimization subproblem (5a) and (5b). It follows from the optimality conditions that

$$0 = \nabla f(x_k) + \rho A^\top (Ax_k - z_k + u_k) + \tau_G(x_{k+1} - x_k), \quad (16a)$$

$$0 \in \partial g(z_{k+1}) - \rho(\alpha Ax_{k+1} + (1-\alpha)z_k - z_{k+1} + u_k), \quad (16b)$$

$$u_{k+1} = u_k + (\alpha Ax_{k+1} + (1-\alpha)z_k - z_{k+1}). \quad (16c)$$

Seeing  $\tau_L$  in the place of  $\tau_G$ , (16b) and (16c) are identical to (9b) and (5c), while (16a) is identical to (9a) with  $\partial f(x_{k+1})$  replaced by  $\nabla f(x_k)$ .

Carrying out the proof of Theorem 2 in §4.1 gives (14) with  $\partial f(x_{k+1})$  replaced by  $\nabla f(x_k)$ , and hence taking corresponding limits gives DI (15) with  $\partial f(X(t))$  replaced by  $\nabla f(X(t))$ . The rest of the proof follows in the same fashion as Part (iii) in the proof of Theorem 2.

## 5. Continuous Limit of Generalized ADMM

In this section, we study the continuous limit of the generalized ADMM (G-ADMM) proposed by Eckstein & Bertsekas (1992). We point out that this has been studied by França et al. (2018a;b) for the cases where  $f$  and  $g$  are both smooth. We now extend the analysis to problems with nonsmooth

$f$  and  $g$ . G-ADMM allows more flexibility of ADMM by introducing a new relaxation parameter  $\alpha \in (0, 2)$ , and it updates the iterates as

$$x_{k+1} = \operatorname{argmin}_x \left\{ f(x) + \frac{\rho}{2} \|Ax - z_k + u_k\|^2 \right\}, \quad (17a)$$

$$z_{k+1} = \operatorname{argmin}_z \left\{ g(z) + \frac{\rho}{2} \|\alpha Ax_{k+1} + (1-\alpha)z_k - z + u_k\|^2 \right\}, \quad (17b)$$

$$u_{k+1} = u_k + (\alpha Ax_{k+1} + (1-\alpha)z_k - z_{k+1}). \quad (17c)$$

### 5.1. Differential Inclusion for G-ADMM

Following França et al. (2018a), we rescale the time by a factor of  $\rho^{-1}$ , i.e. let  $t = \rho^{-1}k$ , and obtain the continuous-time approximation. One can see from (17a) that larger parameter  $\rho > 0$  results in smaller-pace updates of  $x_k$ . We study the limit of updates in the regime  $\rho \rightarrow \infty$ .

Our main result on the differential inclusion approximation of G-ADMM is as follows.

**Theorem 8.** *Let Assumption 1 hold, and the relaxation parameter  $\alpha \in (0, 2)$ . Rescale the time by setting  $t = \rho^{-1}k$ . The continuous limit of iterates of  $\{x_k\}$  in Algorithm (17) as  $\rho \rightarrow \infty$  is given by the following differential inclusion:*

$$\frac{1}{\alpha} \dot{X}(t) + (A^\top A)^{-1} \partial F(X(t)) \ni 0, \quad (18)$$

with  $X(0) = x_0$ .

We move forward and analyze the convergence property of differential inclusion (18).

We recall that  $\sigma_1$  is the largest singular value of  $A$  and  $\sigma_d$  is the smallest singular value of  $A$ .

**Theorem 9.** *When the Assumption 1 holds, the shock solution  $X(t)$  of differential inclusion (18) has bounded trajectory and  $\mathcal{O}(t^{-1})$  convergence rate almost everywhere, i.e., for a.e.  $t \geq 0$ ,*

$$\|X(t) - x^*\| \leq \frac{\sigma_1}{\sigma_d} \|x_0 - x^*\|,$$

$$F(X(t)) - F(x^*) \leq \frac{\sigma_1^2 \|x_0 - x^*\|^2}{2\alpha t}.$$

Moreover,

$$\int_0^{+\infty} [F(X(t)) - F(x^*)] dt \leq \frac{\sigma_1^2}{2\alpha} \|x_0 - x^*\|^2,$$

$$\int_0^{+\infty} t \|\dot{X}(t)\|^2 dt \leq \frac{\sigma_1^2}{2\sigma_d^2} \|x_0 - x^*\|^2.$$

Here, the key idea of the convergence analysis of differential inclusion (18) is to use a sequence of approximating differential equations (ADE) that approaches the differential inclusion.

## 5.2. Differential Inclusion for Accelerated G-ADMM

Goldstein et al. (2014) proposed an accelerated ADMM by incorporating Nesterov's extrapolation technique. This method is generalized by França et al. (2018a;b) which jointly consider relaxation and acceleration. The accelerated G-ADMM considered in França et al. (2018a;b) is described as follows.

$$x_{k+1} = \operatorname{argmin}_x \left\{ f(x) + \frac{\rho}{2} \|Ax - \hat{z}_k + \hat{u}_k\|^2 \right\}, \quad (19a)$$

$$z_{k+1} = \operatorname{argmin}_z \left\{ g(z) + \frac{\rho}{2} \|\alpha Ax_{k+1} + (1 - \alpha)\hat{z}_k - z + \hat{u}_k\|^2 \right\}, \quad (19b)$$

$$u_{k+1} = \hat{u}_k + (\alpha Ax_{k+1} + (1 - \alpha)\hat{z}_k - z_{k+1}), \quad (19c)$$

$$\hat{u}_{k+1} = u_{k+1} + \gamma_{k+1}(u_{k+1} - u_k), \quad (19d)$$

$$\hat{z}_{k+1} = z_{k+1} + \gamma_{k+1}(z_{k+1} - z_k), \quad (19e)$$

where  $\gamma_{k+1} = \frac{k}{k+r}$  with  $r > 0$ .

Here we extend the results in França et al. (2018a;b) to the case where  $f$  and  $g$  are nonsmooth functions. The main results are in the following theorems.

**Theorem 10.** *Let Assumption 1 hold, and the relaxation parameter  $\alpha \in (0, 2)$ . Rescale the time by setting  $t = \rho^{-1/2}k$ , the continuous-time approximation of Algorithm (19) iteration as  $\rho \rightarrow \infty$  is given by the differential inclusion*

$$\frac{1}{\alpha} \left( \ddot{X}(t) + \frac{r}{t} \dot{X}(t) \right) + (A^\top A)^{-1} \partial F(X(t)) \ni 0, \quad (20)$$

with  $X(t_0) = x_0$  and  $\dot{X}(t_0) = 0$ . Here  $t_0$  is an arbitrary positive starting time.

Unlike all previous time rescaling scheme where  $k = \rho t$ , for accelerated G-ADMM the time rescaling is  $k = \rho^{1/2}t$ , which is in accordance with the idea of "acceleration".  $r$  is the damping constant of differential inclusion (20). We refer to the case  $r \geq 3$  as high friction case,  $0 < r < 3$  as low friction case and provide two separate convergence theorems under these two cases. Moreover, we specifically point out that the convergence rate can be sharper when  $r > 3$ .

In the following, we define  $\Delta_0^2 = \max\{t_0^2(F(x_0) - F(x^*)), \|x_0 - x^*\|^2\}$ . We also remark that all the factors  $C$  used later are functions of the damping constant  $r$ , the relaxation parameter  $\alpha$ , and singular values  $\sigma_1$  and  $\sigma_d$ .

Firstly, we show that  $X(t)$  has almost surely bounded trajectory for  $t \geq t_0$ , and the convergence rate is  $\mathcal{O}(t^{-2})$  with regard to both  $F(X(t)) - F(x^*)$  and  $\|\dot{X}(t)\|^2$ .

**Theorem 11 (High Friction).** *When  $r \geq 3$ , the shock solution  $X(t)$  of Differential Inclusion (20) has bounded trajectory and  $\mathcal{O}(t^{-2})$  convergence rate almost everywhere, i.e.*

there exists positive factors  $C_1, C_2, C_3$  such that, for a.e.  $t \geq t_0$ ,

$$\|X(t) - x^*\| \leq C_1 \Delta_0,$$

$$F(X(t)) - F(x^*) \leq \frac{C_2 \Delta_0^2}{t^2}, \quad \|\dot{X}(t)\| \leq \frac{C_3 \Delta_0}{t}.$$

When  $r > 3$ , there exist positive factors  $C_4, C_5$  such that

$$\int_{t_0}^{\infty} t(F(X(t)) - F(x^*)) dt \leq C_4 \Delta_0^2,$$

$$\int_{t_0}^{\infty} t \|\dot{X}(t)\|^2 dt \leq C_5 \Delta_0^2.$$

Using the bounded integration result in Theorem 11, we could show that, when  $r > 3$ , the convergence rate is actually  $\mathcal{O}(t^{-2})$ .

**Remark 12.** *The convergence rate in Theorem 11 is not sharp when  $r > 3$ , in the sense that  $F(X(t)) - F(x^*) = \mathcal{O}(t^{-2})$  and  $\|\dot{X}(t)\| = \mathcal{O}(t^{-1})$  as  $t \rightarrow \infty$ , i.e. for all  $r > 3$ ,*

$$\lim_{t \rightarrow \infty} t^2(F(X(t)) - F(x^*)) = 0, \quad \lim_{t \rightarrow \infty} t \|\dot{X}(t)\| = 0.$$

The damping constant  $r$  has a magic number 3, as discussed in Su et al. (2016), in the sense that a  $\mathcal{O}(t^{-2})$  convergence rate could be guaranteed for AGM in high friction case  $r \geq 3$ , but this result does not apply to low friction case  $0 < r < 3$ . Attouch et al. (2017), Attouch et al. (2018), Vasilis et al. (2018), França et al. (2018a) derive the  $\mathcal{O}(t^{-2r/3})$  convergence rate for AGM under low friction case and extend  $\mathcal{O}(t^{-2})$  convergence rate from AGM to accelerated G-ADMM in high friction case, respectively.

For nonsmooth and low friction case, we show a  $\mathcal{O}(t^{-2r/3})$  convergence rate for  $F(X(t)) - F(x^*)$ , and for  $\|\dot{X}(t)\|^2$  if the trajectory is almost surely bounded.

**Theorem 13 (Low Friction).** *When  $0 < r < 3$ , the shock solution  $X(t)$  of Differential Inclusion (20) has  $\mathcal{O}(t^{-2r/3})$  convergence rate almost everywhere, i.e. there exists positive factor  $C_6$  such that, for a.e.  $t \geq t_0$ ,*

$$F(X(t)) - F(x^*) \leq \frac{C_6 t_0^{-2(3-r)/3} \Delta_0^2}{t^{2r/3}}.$$

If in addition the trajectory  $\{X(t)\}_{t \geq t_0}$  is bounded almost everywhere for  $t \geq t_0$ , then there also exists some positive factor  $C_7$  such that for a.e.  $t \geq t_0$ ,

$$\|\dot{X}(t)\| \leq \frac{C_7 t_0^{-(3-r)/3} \Delta_0}{t^{r/3}}.$$

Theorems 11 and 13 provide convergence results for accelerated G-ADMM in continuous-time scheme with the second-order differential inclusion and accompanying tools, which sheds light on discrete-time scheme (Su et al., 2016).

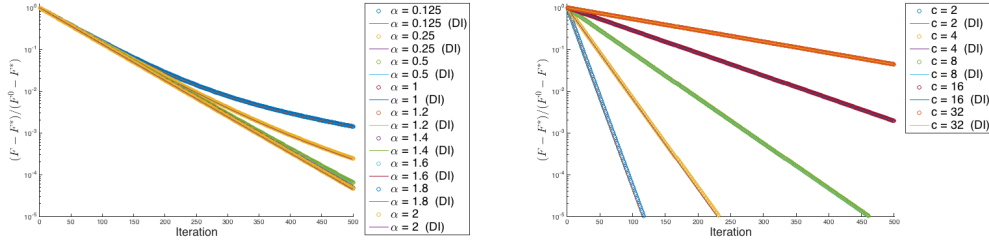


Figure 1. On total variation minimization problem, the plots are the trajectory of linearized ADMM with  $\rho = 10$  and the corresponding differential inclusion, the first plot is for different  $\alpha$  from  $2^{-3}$  to 2 when  $c = 10$ , second plot is for different  $c$  from 1 to 32 when  $\alpha = 1.6$ .

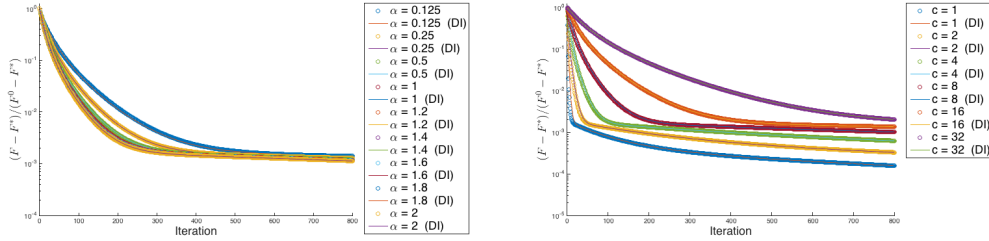


Figure 2. On sparse logistic regression, the plots are gradient ADMM and the differential inclusion when  $\rho = 10$ , first plot is for different  $\alpha$  from  $2^{-3}$  to 2 when  $c = 10$ , second plot is for different  $c$  from 1 to 32 when  $\alpha = 1.6$ .

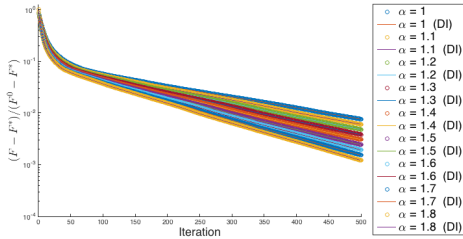


Figure 3. On LASSO problem, the plots are the trajectory of G-ADMM with  $\rho = 50$  and the corresponding differential inclusion.

## 6. Numerical experiments

According to the convergence theorems in previous sections, both the linearized ADMM and the gradient ADMM share the same differential inclusion (8). In the following two examples, we show that for  $\rho = 10$  the trajectory of the linearized and gradient ADMM algorithm is very close to the trajectory of the differential inclusion.

To compute the trajectory of the differential inclusion numerically, we use entropic approximation [Teboulle \(1992\)](#) of the nonsmooth objective to compute the trajectory of the differential inclusion, we remark that entropic approximation is parallel in theory to Moreau-Yosida approximation, which use quadratic approximation. Entropic approximation would provide a smooth approximation of the sub-gradient for  $\|z\|_1$  as  $\tanh \beta z$ , where we choose  $\beta = 10^6$  in the following experiments, so that the numerical approximation

error is of the order  $10^{-6}$ .

### 6.1. Total variation minimization with linearized ADMM

Consider the total variation minimization problem (2), see [Rudin et al. \(1992\)](#). Using variable splitting, we can write the problem as

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|x - b\|_2^2 + \lambda \|z\|_1 \\ & \text{subject to} && z = Dx, \end{aligned}$$

where  $D$  is the finite difference matrix,  $x, z \in \mathbb{R}^n$ . This problem fits to the general framework of ADMM with  $A = D$ ,  $f(x) = \frac{1}{2} \|x - b\|_2^2$  and  $g(z) = \lambda \|z\|_1$ , since  $f$  is quadratic with easy proximal and  $g$  is nonsmooth with easy proximal, namely, its proximal operator is the soft thresholding operator  $S_\lambda(z) = \text{sign}(z)(|z| - \lambda)_+$ , we could use linearized ADMM to solve this problem.

We generated a problem instance, where the true signal  $x^0$  is a piece-wise constant signal, the observation is  $b = x^0 + n$ ,  $n \sim \mathcal{N}(0, I)$ . We solve the total variation minimization with  $\lambda = 0.01$  to recover the true signal  $x^0$  from  $b$ , following the MATLAB examples of the paper [Boyd et al. \(2011\)](#).

In the following, we run both the linearized ADMM algorithm with  $\rho = 10$  and the differential inclusion. We demonstrate the trajectory of linearized ADMM algorithm and the differential inclusion for different parameter configuration  $\alpha$  and  $c$ . First of all, we fix  $c = 10$ , and vary  $\alpha$  from  $2^{-3}$  to 2; then, we fix  $\alpha = 1.6$ , and vary  $c$  from 1 to

32. Figures 1 are the trajectory of linearized ADMM with  $\rho = 10$  and the corresponding differential inclusion, with the  $x$  axis being the iteration count, the  $y$  axis being the relative error  $\frac{F(x_t) - F(x^*)}{F(x_0) - F(x^*)}$ , where  $F(x^*)$  is the function value at optimal solution,  $F(x_0)$  is the function value at initialization. We can see that the differential inclusion matches the trajectory of the linearized ADMM algorithm very closely for all parameter settings.

## 6.2. Sparse logistic regression with gradient ADMM

Consider the sparse logistic regression problem (7) (see Koh et al. (2007) for more details). Using variable splitting, we can write the problem as

$$\begin{aligned} & \text{minimize} && \frac{1}{p} \sum_i \log(1 + \exp(-b_i(a_i^\top w + v))) + \lambda \|z\|_1 \\ & \text{subject to} && z = w, \end{aligned}$$

with variable  $x \in \mathbb{R}^n, v \in \mathbb{R}$ .

This problem fits to the general framework with variable  $\bar{x} = (x, v)$ ,  $A = I$ ,  $f(\bar{x}) = \log(1 + \exp(-b_i(a_i^\top x + v)))$  and  $g(\bar{x}) = \lambda \|\bar{x}_{1:n}\|_1$ , since  $f$  is differentiable but does not have an easy proximal mapping and  $g$  is nonsmooth with easy proximal, we could use gradient ADMM to solve this problem.

We generated a problem instance following the MATLAB examples of the paper Boyd et al. (2011). More specifically, we chose a true weight vector  $x^{\text{true}}$  sampled from Bernoulli-Gaussian distribution with mean 0, variance 1 and sparsity level 0.1, along with the true intercept  $v^{\text{true}}$  sampled from standard normal. Each feature vector  $a_i$  was generated from Bernoulli-Gaussian distribution at sparsity level 0.2. The labels were then generated using  $b_i = \text{sign}(a_i^\top x^{\text{true}} + v^{\text{true}} + \nu_i)$ , where  $\nu_i \sim \mathcal{N}(0, 0.1)$ . The regularization parameter is set to  $\lambda = 0.1\lambda_{\max}$  according to Koh et al. (2007), where  $\lambda_{\max} = \|A^\top \tilde{b}\|_\infty$  is the critical value above which the solution of the problem is  $x = 0$ , where  $\tilde{b}$  is defined in page 93 of Boyd et al. (2011).

We run both the gradient ADMM algorithm with  $\rho = 10$  and the differential inclusion. We demonstrate the trajectory of gradient ADMM algorithm and the differential inclusion for different parameter configuration  $\alpha$  and  $c$ . First of all, we fix  $c = 10$ , and vary  $\alpha$  from  $2^{-3}$  to 2; then, we fix  $\alpha = 1.6$ , and vary  $c$  from 1 to 32. Figures 2 are the trajectory of gradient ADMM with  $\rho = 10$  and the corresponding differential inclusion, with the  $x$  axis being the iteration count, the  $y$  axis being the relative error  $\frac{F(x_t) - F(x^*)}{F(x_0) - F(x^*)}$ , where  $F(x^*)$  is the function value at optimal solution,  $F(x_0)$  is the function value at initialization. We can see that the differential inclusion matches the trajectory of the gradient ADMM algorithm very closely for all parameter settings.

## 6.3. LASSO with G-ADMM

Additionally, we show that the trajectory of G-ADMM algorithm with different  $\alpha$  is close to the trajectory of the differential inclusion in LASSO example Tibshirani (1996b).

The Lasso problem can be casted as

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|Ax - b\|^2 + \lambda \|z\|_1 \\ & \text{subject to} && z = x, \end{aligned}$$

which fits to the general framework with  $f(\cdot) = \frac{1}{2} \|A \cdot - b\|^2$  and  $g(\cdot) = \lambda \|\cdot\|_1$ . When G-ADMM (17) is applied to solve this problem, the two subproblems (17a) and (17b) respectively correspond to the proximal mappings of  $\frac{1}{2} \|A \cdot - b\|^2$  and  $\|\cdot\|_1$ , which are both very easy to compute.

We generated a problem instance following the MATLAB examples of the paper Boyd et al. (2011). Specifically, we sample true sparse signal  $x_0$  from Bernoulli-Gaussian distribution with mean 0, variance 1 and sparsity level  $p = 0.05$ ,  $A$  is sampled from Gaussian random matrix of size 100 by 400 with columns norm normalized to one,  $b = Ax_0 + v$ , where  $v \sim \mathcal{N}(0, 0.001)$ . The regularization parameter is set to  $\lambda = 0.1\lambda_{\max}$  according to Koh et al. (2007), where  $\lambda_{\max} = \|A^\top b\|_\infty$  is the critical value above which the solution of the problem is  $x = 0$ .

We run both the G-ADMM algorithm with  $\rho = 50$  and the differential inclusion. We vary  $\alpha$  from 1 to 1.8 as the range that people typically use for G-ADMM Boyd et al. (2011). Figure 3 is the trajectory of G-ADMM with  $\rho = 50$  and the corresponding differential inclusion, with the  $x$  axis being the iteration count, the  $y$  axis being the relative error  $\frac{F(x_t) - F(x^*)}{F(x_0) - F(x^*)}$ , where  $F(x^*)$  is the function value at optimal solution,  $F(x_0)$  is the function value at initialization. We can see that the differential inclusion matches the trajectory of the linearized ADMM algorithm very closely for all parameter settings.

## 7. Conclusions

In this paper, we analyzed the convergence behavior of the continuous limits of some widely used nonsmooth ADMM variants: linearized ADMM, gradient-based ADMM, as well as G-ADMM and its Nesterov's acceleration. Such continuous limits are characterized by the tool of differential inclusion and promote understandings of these ADMM variants from the angles of dynamical systems. Our novel continuous-time convergence theorems characterize these ADMM variants, which is further supported by experimental results. The differential inclusion for linearized and gradient ADMM variants suggests that we could choose the algorithmic parameters via a principled approach that uses the condition number of the matrix, which serves as a practical guidance learned from theoretical insights.



## References

- Adly, S., Attouch, H., and Cabot, A. Finite time stabilization of nonlinear oscillators subject to dry friction. In *Nonsmooth mechanics and analysis*, pp. 289–304. Springer, 2006.
- Attouch, H., Cabot, A., and Redont, P. The dynamics of elastic shocks via epigraphical regularization of a differential inclusion barrier and penalty approximations. *Adv. Math. Sci. Appl.*, 12:273–306, 2002.
- Attouch, H., Chbani, Z., and Riahi, H. Rate of convergence of the nesterov accelerated gradient method in the subcritical case  $\alpha \leq 3$ . *arXiv preprint arXiv:1706.05671*, 2017.
- Attouch, H., Chbani, Z., Peypouquet, J., and Redont, P. Fast convergence of inertial dynamics and algorithms with asymptotic vanishing viscosity. *Mathematical Programming*, 168(1-2):123–175, 2018.
- Boyd, S., Parikh, N., Chu, E., Peleato, B., Eckstein, J., et al. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3(1):1–122, 2011.
- Chen, G. and Teboulle, M. A proximal-based decomposition method for convex minimization problems. *Mathematical Programming*, 64:81–101, 1994.
- Combettes, P. L. and Pesquet, J.-C. A Douglas-Rachford splitting approach to nonsmooth convex variational signal recovery. *IEEE Journal of Selected Topics in Signal Processing*, 1(4):564–574, 2007.
- Condat, L. A primal-dual splitting method for convex optimization involving Lipschitzian, proximable and linear composite terms. *J. Optimization Theory and Applications*, 158(2):460–479, 2013.
- Davis, D. and Yin, W. A three-operator splitting scheme and its optimization applications. *Set-Valued and Variational Analysis*, 25(4):829–858, 2017.
- Douglas, J. and Rachford, H. H. On the numerical solution of the heat conduction problem in 2 and 3 space variables. *Transactions of the American Mathematical Society*, 82:421–439, 1956.
- Eckstein, J. Some saddle-function splitting methods for convex programming. *Optimization Methods and Software*, 4(1):75–83, 1994.
- Eckstein, J. and Bertsekas, D. P. On the Douglas–Rachford splitting method and the proximal point algorithm for maximal monotone operators. *Mathematical Programming*, 55(1-3):293–318, 1992.
- Fortin, M. and Glowinski, R. *Augmented Lagrangian methods: applications to the numerical solution of boundary-value problems*. North-Holland Pub. Co., 1983.
- França, G., Robinson, D. P., and Vidal, R. ADMM and accelerated ADMM as continuous dynamical systems. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 1559–1567, 2018a.
- França, G., Robinson, D. P., and Vidal, R. Relax, and accelerate: A continuous perspective on ADMM. *arXiv preprint arXiv:1808.04048*, 2018b.
- Gabay, D. Applications of the method of multipliers to variational inequalities. In Fortin, M. and Glowinski, R. (eds.), *Augmented Lagrangian Methods: Applications to the Solution of Boundary Value Problems*. North-Holland, Amsterdam, 1983.
- Gabay, D. and Mercier, B. *A dual algorithm for the solution of non linear variational problems via finite element approximation*. Institut de recherche d’informatique et d’automatique, 1975.
- Glowinski, R. and Le Tallec, P. *Augmented Lagrangian and Operator-Splitting Methods in Nonlinear Mechanics*. SIAM, Philadelphia, Pennsylvania, 1989.
- Glowinski, R. and Marroco, A. Sur l’approximation, par éléments finis d’ordre un, et la résolution, par pénalisation-dualité d’une classe de problèmes de dirichlet non linéaires. *Revue française d’automatique, informatique, recherche opérationnelle. Analyse numérique*, 9(R2):41–76, 1975.
- Goldstein, T. and Osher, S. The split Bregman method for L1-regularized problems. *SIAM J. Imaging Sci.*, 2:323–343, 2009.
- Goldstein, T., O’Donoghue, B., Setzer, S., and Baraniuk, R. Fast alternating direction optimization methods. *SIAM Journal on Imaging Sciences*, 7(3):1588–1623, 2014.
- He, B. S., Liao, L., Han, D., and Yang, H. A new inexact alternating direction method for monotone variational inequalities. *Mathematical Programming*, 92:103–118, 2002.
- Koh, K., Kim, S.-J., and Boyd, S. An interior-point method for large-scale l1-regularized logistic regression. *Journal of Machine learning research*, 8(Jul):1519–1555, 2007.
- Krichene, W., Bayen, A., and Bartlett, P. L. Accelerated mirror descent in continuous and discrete time. In *NIPS*, 2015.
- Lin, T., Ma, S., and Zhang, S. An extragradient-based alternating direction method for convex minimization.

- 495 *Foundations of Computational Mathematics*, 17(1):35–  
496 59, 2017.
- 497 Lin, Z., Liu, R., and Su, Z. Linearized alternating direction  
498 method with adaptive penalty for low-rank representation.  
499 In *Advances in neural information processing systems*,  
500 pp. 612–620, 2011.
- 501 Liu, J., Chen, J., and Ye, J. Large-scale sparse logistic  
502 regression. In *SIGKDD*, 2009.
- 503 Ma, S. Alternating proximal gradient method for convex  
504 minimization. *Journal of Scientific Computing*, 68(2):  
505 546–572, 2016.
- 506 Nesterov, Y. A method for unconstrained convex minimiza-  
507 tion problem with the rate of convergence  $O(1/k^2)$ . In  
508 *Soviet Mathematics Doklady*, volume 27, pp. 372–376,  
509 1983.
- 510 Ouyang, Y., Chen, Y., Lan, G., and Pasiliao Jr, E. An  
511 accelerated linearized alternating direction method of  
512 multipliers. *SIAM Journal on Imaging Sciences*, 8(1):  
513 644–681, 2015.
- 514 Paoli, L. A. An existence result for vibrations with unilateral  
515 constraints: case of a nonsmooth set of constraints. *Math.*  
516 *Models Methods Appl. Sci.*, 10:815–831, 2000.
- 517 Peaceman, D. H. and Rachford, H. H. The numerical so-  
518 lution of parabolic elliptic differential equations. *SIAM*  
519 *Journal on Applied Mathematics*, 3:28–41, 1955.
- 520 Rudin, L. I., Osher, S., and Fatemi, E. Nonlinear total  
521 variation based noise removal algorithms. *Physica D:*  
522 *nonlinear phenomena*, 60(1-4):259–268, 1992.
- 523 Shi, B., Du, S. S., Jordan, M. I., and Su, W. J. Understanding  
524 the acceleration phenomenon via high-resolution differ-  
525 ential equations. *arXiv preprint arXiv:1810.08907*, 2018.
- 526 Su, W., Boyd, S., and Candes, E. J. A differential equation  
527 for modeling Nesterov’s accelerated gradient method: the-  
528 ory and insights. *Journal of Machine Learning Research*,  
529 17(153):1–43, 2016.
- 530 Teboulle, M. Entropic proximal mappings with applications  
531 to nonlinear programming. *Mathematics of Operations*  
532 *Research*, 17(3):670–690, 1992.
- 533 Tibshirani, R. Regression shrinkage and selection via the  
534 lasso. *Journal Royal Statistical Society B*, 58:267–288,  
535 1996a.
- 536 Tibshirani, R. Regression shrinkage and selection via the  
537 lasso. *Journal of the Royal Statistical Society. Series B*  
538 *(Methodological)*, pp. 267–288, 1996b.
- 539 Vassilis, A., Jean-François, A., and Charles, D. The differen-  
540 tial inclusion modeling FISTA algorithm and optimality  
541 of convergence rate in the case  $b \leq 3$ . *SIAM Journal on*  
542 *Optimization*, 28(1):551–574, 2018.
- 543 Vu, B. C. A splitting algorithm for dual monotone inclusions  
544 involving cocoercive operators. *Advances in Computa-*  
545 *tional Mathematics*, 38(3):667–681, 2013.
- 546 Wang, Y., Yang, J., Yin, W., and Zhang, Y. A new alter-  
547 nating minimization algorithm for total variation image  
548 reconstruction. *SIAM Journal on Imaging Sciences*, 1(3):  
549 248–272, 2008.
- 550 Wibisono, A., Wilson, A. C., and Jordan, M. I. A varia-  
551 tional perspective on accelerated methods in optimization.  
552 *Proceedings of the National Academy of Sciences*, pp.  
553 201614734, 2016.
- 554 Wilson, A. C., Recht, B., and Jordan, M. I. A lyapunov  
555 analysis of momentum methods in optimization. *arXiv*  
556 *preprint arXiv:1611.02635*, 2016.
- 557 Xu, Y. Alternating proximal gradient method for sparse  
558 nonnegative Tucker decomposition. *Mathematical Pro-*  
559 *gramming Computation*, 7(1):39–70, 2015.
- 560 Yang, J. and Yuan, X. Linearized augmented lagrangian and  
561 alternating direction methods for nuclear norm minimiza-  
562 tion. *Mathematics of Computation*, 82(281):301–329,  
563 2013.
- 564 Yang, J. and Zhang, Y. Alternating direction algorithms for  
565  $\ell_1$ -problems in compressive sensing. *SIAM journal on*  
566 *scientific computing*, 33(1):250–278, 2011.
- 567 Zhang, X., Burger, M., Bresson, X., and Osher, S. Breg-  
568 manized nonlocal regularization for deconvolution and  
569 sparse reconstruction. *SIAM Journal on Imaging Science*,  
570 3:253–276, 2010.

## A. Proof of Theorem 5

*Proof of Theorem 5.* For notation simplicity, we choose a matrix  $B$  such that  $B^\top B = cI + \frac{1-\alpha}{\alpha}A^\top A$ . Recall that the largest and smallest singular value of  $B$  are  $\kappa_1$  and  $\kappa_d$ . Note that  $\kappa_1 = \sqrt{c + \frac{1-\alpha}{\alpha}\sigma_1^2}$  and  $\kappa_d = \sqrt{c + \frac{1-\alpha}{\alpha}\sigma_d^2}$  where  $\sigma_1, \sigma_d$  are singular value of matrix  $A$ . Then the original DI becomes  $0 \in \partial F(X(t)) + (B^\top B)\dot{X}(t)$ . Because Moreau-Yosida approximation  $F_\mu(X_\mu(t))$  is a continuously differentiable, convex function for all  $\mu > 0$ , we denote its minimizer as  $x_\mu^*$ .

For each  $\mu > 0$ , consider the energy functional of Moreau-Yosida approximation defined as

$$\mathcal{E}_\mu(t) = t(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) + \frac{\lambda}{2}\|B(X_\mu(t) - x_\mu^*)\|^2, \quad (21)$$

where  $\lambda$  is an arbitrary constant chosen as  $\lambda \geq 1$ . Because  $F_\mu$  is a continuously differentiable function, we could write the time derivative of  $\mathcal{E}_\mu(t)$  as

$$\dot{\mathcal{E}}_\mu(t) = (F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) + t\langle \nabla F_\mu(X_\mu(t)), \dot{X}_\mu(t) \rangle + \lambda\langle B^\top B(X_\mu(t) - x_\mu^*), \dot{X}_\mu(t) \rangle. \quad (22)$$

By substituting  $B^\top B\dot{X}_\mu(t)$  with  $-\nabla F_\mu(X_\mu(t))$  and vice versa, we have

$$\dot{\mathcal{E}}_\mu(t) = -t\|B\dot{X}_\mu(t)\|^2 + (F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) - \lambda\langle (X_\mu(t) - x_\mu^*), \nabla F_\mu(X_\mu(t)) \rangle \leq 0, \quad (23)$$

where we used the convexity of  $F_\mu$  and nonnegativity of  $(F_\mu(X_\mu) - F_\mu(x_\mu^*)), \|B\dot{X}_\mu\|$  in the last inequality.

Similar to  $\mathcal{E}_\mu(t)$ , we define the energy functional for  $F(X(t))$  as

$$\mathcal{E}(t) = t(F(X(t)) - F(x^*)) + \frac{\lambda}{2}\|B(X(t) - x^*)\|^2. \quad (24)$$

At time  $t = 0$ , there is an upper bound on  $\mathcal{E}(0)$  as

$$\mathcal{E}(0) = \frac{\lambda}{2}\|B(x_0 - x^*)\|^2 \leq \frac{\lambda\kappa_1^2}{2}\|x_0 - x^*\|^2. \quad (25)$$

By applying the approximation scheme (AS) as  $\mu \rightarrow 0$ , we have for a.e.  $t \geq 0$  that  $\mathcal{E}(t) \leq \mathcal{E}(0)$ .

By non-negativity of  $F(X) - F(x^*)$  in (24), we find

$$\frac{\lambda\kappa_d^2}{2}\|X(t) - x^*\|^2 \leq \mathcal{E}(0). \quad (26)$$

Combining with upper bound of  $\mathcal{E}(0)$  in (25), we derive for a.e.  $t \geq 0$  that

$$\|X(t) - x^*\| \leq \frac{\kappa_1}{\kappa_d}\|x_0 - x^*\|. \quad (27)$$

Using the nonnegativity of all terms in (24) and monotonicity of  $\mathcal{E}(t)$  on a.e.  $t \geq 0$ , we have, for a.e.  $t \geq 0$ ,

$$t(F(X(t)) - F(x^*)) \leq \mathcal{E}(t) \leq \mathcal{E}(0) \leq \frac{\lambda\kappa_1^2}{2}\|x_0 - x^*\|^2 \quad (28)$$

Choosing  $\lambda = 1$ , we have the following result, for a.e.  $t \geq 0$ ,

$$F(X(t)) - F(x^*) \leq \mathcal{E}(t) \leq \mathcal{E}(0) \leq \frac{\kappa_1^2}{2t}\|x_0 - x^*\|^2 \quad (29)$$

By applying convexity of  $F_\mu$  to (23), we have

$$\dot{\mathcal{E}}_\mu(t) \leq (1 - \lambda)(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) - t\|B\dot{X}_\mu(t)\|^2. \quad (30)$$

Notice that the two terms in (30) are all negative, we find

$$F_\mu(X_\mu(t)) - F_\mu(x_\mu^*) \leq \frac{-\mathcal{E}_\mu(t)}{\lambda - 1} \quad \text{and} \quad t \|\dot{X}_\mu(t)\|^2 \leq -\frac{\mathcal{E}_\mu(t)}{\kappa_d^2}. \quad (31)$$

By integrating over  $(0, T)$ , the inequalities above give for all  $T > 0$  that

$$\int_0^T (F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) dt \leq \frac{\mathcal{E}_\mu(0)}{\lambda - 1}, \quad \int_0^T t \|\dot{X}_\mu(t)\|^2 dt \leq \frac{\mathcal{E}_\mu(0)}{\kappa_d^2}. \quad (32)$$

By applying approximation scheme  $(\mathcal{AS})$ , taking limit  $T \rightarrow \infty$ , choosing  $\lambda \rightarrow \infty$  and  $\lambda = 1$  respectively, and plugging in (25), we have

$$\int_0^\infty (F(X_\mu(t)) - F(x^*)) dt \leq \frac{\kappa_1^2}{2} \|x_0 - x^*\|^2, \quad \int_0^\infty t \|\dot{X}(t)\|^2 dt \leq \frac{\kappa_1^2}{2\kappa_d^2} \|x_0 - x^*\|^2. \quad (33)$$

□

## B. Preliminaries of Differential Inclusion

Recall that we denote  $F(x) = f(x) + g(Ax)$ , and Assumption 1 holds. To transit from the smooth case to the nonsmooth case, we use the tool of differential inclusion (DI) to build the connection between subdifferentiable  $F$  and differentiable functions. One basic example of a differential inclusion takes the form of:

$$\dot{x}(t) \in \partial F(x(t))$$

To bridge the gap between differentiable objective functions and nondifferentiable objective functions, we follow Vassilis et al. (2018) and consider the Moreau-Yosida Approximation, which is a standard tool in convex analysis.

**Definition 14** (Moreau-Yosida Approximation). *Moreau-Yosida Approximation of a convex function  $F$  with parameter  $\mu > 0$  is defined as*

$$F_\mu(x) := \inf_y \left\{ F(y) + \frac{1}{2\mu} \|y - x\|^2 \right\}$$

Use  $J_\mu(x)$  to denote the unique point that achieves the infimum above, then  $\nabla F_\mu(x) = \frac{1}{\mu}(x - J_\mu(x))$ . For any  $\mu > 0$ ,  $F_\mu$  is a convex, continuously differentiable function.

We take the Definition 3.1 in Vassilis et al. (2018) of a shock solution to define a solution of a differential inclusion. The existence of a shock solution are described in Section 3 of Vassilis et al. (2018). More specifically, we can build a sequence  $x_\mu(t)$  such that its subsequence converges, where  $x_\mu(t)$  are the solutions to the Approximate Differential Equation (ADE) defined below:

### Approximate Differential Equation (ADE)

We consider the Moreau-Yosida approximation  $F_\mu(x)$  of the objective  $F(x)$  with  $\mu > 0$ . We consider the following approximating ODE:

$$\begin{cases} \dot{x}_\mu(t) + \nabla F_\mu(x_\mu(t)) = 0 \\ x_\mu(0) = x_0 \end{cases}$$

Here  $\nabla F_\mu$  can approximate  $\partial F$  and  $F_\mu$  is differentiable as is shown in the theory of Moreau-Yosida approximation.

The convergence to a shock solution is described as the Approximation Scheme  $(\mathcal{AS})$ :

### Approximation Scheme $(\mathcal{AS})$

Let  $\{F_\mu\}_{\mu>0}$  be a family of functions such that  $F_\mu$  is the Moreau-Yosida approximation of  $F$  for all  $\mu > 0$ . Then there exists a subsequence  $\{x_\mu\}_{\mu>0}$  of solutions of (ADE) that converges to a shock solution  $x$  of (DI) in the following sense:



- $x_\mu \rightarrow x$  uniformly on  $[0, T]$  for all  $T > 0$  as  $\mu \rightarrow 0$
- $\dot{x}_\mu \rightarrow \dot{x}$  in  $L^p([0, T]; \mathbb{R}^d)$  for all  $p \in [1, \infty)$  for all  $T > 0$  as  $\mu \rightarrow 0$
- $F_\mu(x_\mu) \rightarrow F(x)$  in  $L^p([0, T]; \mathbb{R}^d)$  for all  $p \in [1, \infty)$  for all  $T > 0$  as  $\mu \rightarrow 0$

### C. Proofs of the theorems related to G-ADMM and the A-ADMM

*Proof of Theorem 8.* Due to the strong convexity of the optimization subproblems (17a) and (17b), it is easy to verify that the sequence  $\{x_k, z_k, u_k\}$  is unique. Together with (17c), we have from the first-order optimality conditions of (17a) and (17b) that

$$\partial f(x_{k+1}) + \rho A^T (Ax_{k+1} - z_k + u_k) \ni 0, \quad (34a)$$

$$\frac{1}{\rho} \partial g(z_{k+1}) - (\alpha Ax_{k+1} + (1 - \alpha)z_k - z_{k+1} + u_k) \ni 0, \quad (34b)$$

$$u_{k+1} - (\alpha Ax_{k+1} + (1 - \alpha)z_k - z_{k+1} + u_k) = 0. \quad (34c)$$

We detail the proof in the following:

- (i) Adding up (34b) and (34c) eliminates the common term  $(\alpha Ax_{k+1} + (1 - \alpha)z_k - z_{k+1} + u_k)$  and reduces to a simple  $u$ -update:

$$u_{k+1} \in \frac{1}{\rho} \partial g(z_{k+1}). \quad (35)$$

Taking the continuous limit  $\rho \rightarrow \infty$  gives  $U(t) = 0$ , and hence  $\dot{U}(t) = 0$ .<sup>2</sup>

- (ii) Bringing (35) into (34a) leads to:

$$0 \in \partial f(x_{k+1}) + A^\top \partial g(z_k) + \rho A^\top (Ax_{k+1} - z_k), \quad (36)$$

where again from (34c),

$$u_{k+1} - u_k = \alpha (Ax_{k+1} - z_k) - (z_{k+1} - z_k),$$

and hence

$$Ax_{k+1} - z_k = \frac{1}{\alpha} [(u_{k+1} - u_k) + (z_{k+1} - z_k)] \quad (37)$$

Plugging (37) into (36) gives

$$0 \in \partial f(x_{k+1}) + A^\top \partial g(z_k) + \rho A^\top \left( \frac{1}{\alpha} [(u_{k+1} - u_k) + (z_{k+1} - z_k)] \right). \quad (38)$$

Taking the limit  $\rho \rightarrow \infty$ , using the fact that  $\dot{U}(t) = 0$ , (38) reduces to

$$0 \in \partial f(X(t)) + A^\top \partial g(Z(t)) + \frac{1}{\alpha} A^\top (\dot{Z}(t)). \quad (39)$$

- (iii) We directly take the  $\rho \rightarrow \infty$  limit in (34c) and conclude

$$Z(t) = AX(t), \quad \dot{Z}(t) = A\dot{X}(t).$$

Combining the above and (39) concludes

$$0 \in \partial F(X(t)) + \left( \frac{1}{\alpha} A^\top A \right) \dot{X}(t),$$

which completes the proof.

<sup>2</sup>Although the continuous version of  $U(t)$  is constantly zero, it is different with  $u_k = 0$ . One may regard  $u_k$  as an infinitesimal number that dynamically changes in the system.

□

*Proof of Theorem 9.* For each  $\mu > 0$ , consider the energy functional of Moreau-Yosida approximation defined as

$$\mathcal{E}_\mu(t) = \alpha t(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) + \frac{\lambda}{2} \|A(X_\mu(t) - x_\mu^*)\|^2, \quad (40)$$

where  $\lambda$  is an arbitrary constant chosen as  $\lambda \geq 1$  and  $x_\mu^*$  denotes the minimizer of  $F_\mu$ . Because  $F_\mu$  is a continuously differentiable function, we could write the time derivative of  $\mathcal{E}_\mu(t)$  as

$$\dot{\mathcal{E}}_\mu(t) = \alpha(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) + \alpha t \langle \nabla F_\mu(X_\mu(t)), \dot{X}_\mu(t) \rangle + \lambda \langle A^\top A(X_\mu(t) - x_\mu^*), \dot{X}_\mu(t) \rangle \quad (41)$$

By substituting  $A^\top A \dot{X}_\mu(t)$  with  $-\alpha \nabla F_\mu(X_\mu(t))$  and vice versa, we have

$$\dot{\mathcal{E}}_\mu(t) = -t \|A \dot{X}_\mu(t)\|^2 + \alpha(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) - \lambda \alpha \langle (X_\mu(t) - x_\mu^*), \nabla F_\mu(X_\mu(t)) \rangle \leq 0, \quad (42)$$

where we used the convexity of  $F_\mu$  and nonnegativity of  $(F_\mu(X_\mu(t)) - F(x_\mu^*))$ ,  $\|A \dot{X}_\mu\|$  in the last inequality.

Similar to  $\mathcal{E}_\mu(t)$ , we define the energy functional for  $F(X(t))$  as

$$\mathcal{E}(t) = \alpha t(F(X(t)) - F(x^*)) + \frac{\lambda}{2} \|A(X(t) - x^*)\|^2. \quad (43)$$

At time 0, there is an upper bound on  $\mathcal{E}(0)$  as

$$\mathcal{E}(0) = \frac{\lambda}{2} \|A(X(0) - x^*)\|^2 \leq \frac{\lambda \sigma_1^2}{2} \|x_0 - x^*\|^2. \quad (44)$$

By applying the approximation scheme ( $\mathcal{AS}$ ) as  $\mu \rightarrow 0$  to equation (42), we have for a.e.  $t \geq 0$ ,  $\dot{\mathcal{E}}(t) \leq 0$  and that  $\mathcal{E}(t) \leq \mathcal{E}(0)$ .

In (43), by non-negativity of  $F(X(t)) - F(x^*)$  and  $\|X(t) - x^*\|^2$ , we find

$$\frac{\lambda}{2} \|A(X(t) - x^*)\|^2 \leq \mathcal{E}(0). \quad (45)$$

Combining with upper bound of  $\mathcal{E}(0)$  in (44), and by taking  $\lambda = 1$ , we derive for a.e.  $t \geq 0$  that

$$\|X(t) - x^*\| \leq \frac{\sigma_1}{\sigma_d} \|x_0 - x^*\|. \quad (46)$$

Using the nonnegativity of all terms in (43) and monotonicity of  $\mathcal{E}(t)$  on a.e.  $t \geq 0$ , we have

$$\alpha t(F(X(t)) - F(x^*)) \leq \mathcal{E}(t) \leq \mathcal{E}(0) \leq \frac{\lambda \sigma_1^2}{2} \|x_0 - x^*\|^2 \quad \text{for a.e. } t, \quad (47)$$

which is given by (44). Thus  $(F(X(t)) - F(x^*)) \leq \frac{\sigma_1^2}{2\alpha t} \|x_0 - x^*\|^2$  by taking  $\lambda = 1$ .

From (42) and using the convexity of  $F_\mu$ , we have

$$\dot{\mathcal{E}}_\mu(t) \leq \alpha(1 - \lambda)(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) - t \|A \dot{X}_\mu(t)\|^2. \quad (48)$$

Notice that the two terms in (48) are all negative, we find

$$F_\mu(X_\mu(t)) - F_\mu(x_\mu^*) \leq \frac{-\dot{\mathcal{E}}_\mu(t)}{\alpha(\lambda - 1)} \quad \text{and} \quad t \|A \dot{X}_\mu(t)\|^2 \leq -\dot{\mathcal{E}}_\mu(t). \quad (49)$$

By integrating over  $(0, T)$ , the inequalities above give

$$\int_0^T (F_\mu(X_\mu(t)) - F_\mu(x_\mu^*))dt \leq \frac{\mathcal{E}_\mu(0)}{\alpha(\lambda - 1)}, \quad \int_0^T t \|A\dot{X}_\mu(t)\|^2 dt \leq \mathcal{E}_\mu(0). \quad (50)$$

By applying approximation scheme (AS) and plugging in (44), we have

$$\int_0^T (F_\mu(X_\mu(t)) - F_\mu(x_\mu^*))dt \leq \frac{\lambda\sigma_1^2}{2\alpha(\lambda - 1)} \|x_0 - x^*\|^2, \quad \int_0^T t \|A\dot{X}_\mu(t)\|^2 dt \leq \frac{\lambda\sigma_1^2}{2} \|x_0 - x^*\|^2. \quad (51)$$

Taking the limit when  $\mu \rightarrow 0$ ,  $T \rightarrow \infty$  and choosing  $\lambda \rightarrow \infty$  and  $\lambda = 1$  respectively, we get

$$\int_0^\infty (F(X(t)) - F(x^*))dt \leq \frac{\sigma_1^2}{2\alpha} \|x_0 - x^*\|^2, \quad \int_0^\infty t \|\dot{X}(t)\|^2 dt \leq \frac{\sigma_1^2}{2\sigma_d^2} \|x_0 - x^*\|^2. \quad (52)$$

Which completes our proof.  $\square$

*Proof of Theorem 10.* Due to the strong convexity of the optimization subproblems (19a) and (19b), it is easy to verify that the sequence  $\{x_k, z_k, u_k, \hat{u}_k, \hat{z}_k\}$  is unique. Together with (19c), we have from the first-order optimality conditions of (19a) and (19b) that

$$\partial f(x_{k+1}) + \rho A^T (Ax_{k+1} - \hat{z}_k + \hat{u}_k) \ni 0, \quad (53a)$$

$$\frac{1}{\rho} \partial g(z_{k+1}) - (\alpha Ax_{k+1} + (1 - \alpha)\hat{z}_k - z_{k+1} + \hat{u}_k) \ni 0, \quad (53b)$$

$$u_{k+1} - (\alpha Ax_{k+1} + (1 - \alpha)\hat{z}_k - z_{k+1} + \hat{u}_k) = 0. \quad (53c)$$

- (i) Adding up (53b) and (53c) eliminates the common term  $-(\alpha Ax_{k+1} + (1 - \alpha)\hat{z}_k - z_{k+1} + \hat{u}_k)$  and reduces to a simple  $u$ -update:

$$u_{k+1} \in \frac{1}{\rho} \partial g(z_{k+1}). \quad (54)$$

Taking the continuous limit  $\rho \rightarrow \infty$  gives  $U(t) = 0$ , and hence  $\dot{U}(t) = 0$ ,  $\ddot{U}(t) = 0$ .<sup>3</sup>

- (ii) Bringing (54) and equation (19d) which is the definition of  $\hat{u}$  into (53a) leads to:

$$\partial f(x_{k+1}) + A^T \partial g(z_k) + \rho A^T (Ax_{k+1} - \hat{z}_k) + \rho \gamma_{k+1} A^\top (u_{k+1} - u_k) \ni 0, \quad (55)$$

where again from (53c),

$$(Ax_{k+1} - \hat{z}_k) = \frac{1}{\alpha} [(u_{k+1} - \hat{u}_k) + (z_{k+1} - \hat{z}_k)]. \quad (56)$$

In addition, from equation (19d) and equation (19e), we find that  $u_{k+1} - \hat{u}_k = u_{k+1} - (1 + \gamma_{k+1})u_k + \gamma_{k+1}u_{k-1}$  and  $z_{k+1} - \hat{z}_k = z_{k+1} - (1 + \gamma_{k+1})z_k + \gamma_{k+1}z_{k-1}$ . For  $u_{k+1} - \hat{u}_k$ , we add the term  $u_k - u_k + u_{k-1} - u_{k-1}$  to the right hand side, the resulting equation is a combination of the second order difference and first order difference of the sequence  $\{u_k\}$ :

$$u_{k+1} - \hat{u}_k = (u_{k+1} - 2u_k + u_{k-1}) + (1 - \gamma_{k+1})(u_k - u_{k-1}). \quad (57)$$

Similarly, the equation holds that:

$$z_{k+1} - \hat{z}_k = (z_{k+1} - 2z_k + z_{k-1}) + (1 - \gamma_{k+1})(z_k - z_{k-1}). \quad (58)$$

We note that  $1 - \gamma_k = 1 - \frac{k}{k+r} = \frac{r}{\rho^{1/2}t+r}$ . Taking the limit  $\rho \rightarrow \infty$ , under infinitesimal step sizes, using relationships (56), (57), (58) and the fact that  $\dot{U}(t) = 0$ ,  $\ddot{U}(t) = 0$ , equation (55) becomes:

$$\partial f(X(t)) + A^T \partial g(Z(t)) + \frac{1}{\alpha} A^T \left( \frac{r}{t} \dot{Z}(t) + \ddot{Z}(t) \right) \ni 0. \quad (59)$$

<sup>3</sup>Although the continuous version of  $U(t)$  is constantly zero, it is different with  $u_k = 0$ . One may regard  $u_k$  as an infinitesimal number that dynamically changes in the system.

(iii) We directly take the  $\rho \rightarrow \infty$  limit in (53c) and conclude

$$Z(t) = AX(t), \quad \dot{Z}(t) = A\dot{X}(t), \quad \ddot{Z}(t) = A\ddot{X}(t).$$

Combining the above and (59) concludes

$$0 \in \partial F(X(t)) + \left( \frac{1}{\alpha} A^\top A \right) (\dot{X}(t) + \frac{r}{t} \dot{X}(t)),$$

which completes the proof.  $\square$

*Proof of Theorem 11.* Recall that  $x_\mu^*$  is the minimizer of  $F_\mu$ . For each  $\mu > 0$ , consider the energy functional of Moreau-Yosida approximation defined as

$$\mathcal{E}_\mu(t) = t^2(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) + \frac{1}{2\alpha} \left\| A \left( \lambda(X_\mu(t) - x_\mu^*) + t\dot{X}_\mu(t) \right) \right\|^2 + \frac{\lambda(r - \lambda - 1)}{2\alpha} \|A(X_\mu(t) - x_\mu^*)\|^2 \quad (60)$$

where  $\lambda$  is a constant chosen within  $2 \leq \lambda \leq r - 1$ . Because  $F_\mu$  is a continuously differentiable function, we could write the time derivative of  $\mathcal{E}_\mu(t)$  as

$$\begin{aligned} \dot{\mathcal{E}}_\mu &= 2t(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) + t^2 \nabla F_\mu(X_\mu(t))^\top \dot{X}_\mu + \left( \lambda(X_\mu - x_\mu^*) + t\dot{X}_\mu \right)^\top \left( \frac{1}{\alpha} A^\top A \right) \left( (\lambda + 1)\dot{X}_\mu + t\ddot{X}_\mu \right) \\ &\quad + \lambda(r - \lambda - 1)(X_\mu - x_\mu^*)^\top \left( \frac{1}{\alpha} A^\top A \right) \dot{X}_\mu \end{aligned}$$

By substituting  $tA^\top A\ddot{X}_\mu$  with  $-rA^\top A\dot{X}_\mu - \alpha t \nabla F_\mu(X_\mu(t))$ , we have

$$\dot{\mathcal{E}}_\mu = -\lambda t (F_\mu(x_\mu^*) - F_\mu(X_\mu(t)) - (x_\mu^* - X_\mu)^\top \nabla F_\mu(X_\mu(t))) - (\lambda - 2)t(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) - \frac{(r - 1 - \lambda)t}{\alpha} \|A\dot{X}_\mu\|^2 \leq 0 \quad (61)$$

where we used the convexity of  $F_\mu$  and nonnegativity of  $F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)$ ,  $\|A\dot{X}_\mu\|$  in the last inequality.

Similar to  $\mathcal{E}_\mu(t)$ , we define the energy functional for  $F(X(t))$  as

$$\mathcal{E}(t) = t^2(F(X(t)) - F(x^*)) + \frac{1}{2\alpha} \left\| A \left( \lambda(X(t) - x^*) + t\dot{X}(t) \right) \right\|^2 + \frac{\lambda(r - \lambda - 1)}{2\alpha} \|A(X(t) - x^*)\|^2$$

At time  $t_0$ , there is an upper bound on  $\mathcal{E}(t_0)$  as

$$\mathcal{E}(t_0) = t_0^2(F(x_0) - F(x^*)) + \frac{\lambda(r - 1)}{2\alpha} \|A(x_0 - x^*)\|^2 \leq \frac{2\alpha + \lambda(r - 1)\sigma_1^2}{2\alpha} \Delta_0^2 \quad (62)$$

By non-negativity of  $F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)$ ,  $\|X_\mu - x_\mu^*\|^2$  and  $\|\dot{X}_\mu\|^2$ , we find for all  $r \geq 3$  and  $t \geq t_0$  that

$$\frac{d}{dt} (t\|X_\mu - x_\mu^*\|^2) = \|X_\mu - x_\mu^*\|^2 + 2t(X_\mu - x_\mu^*)^\top \dot{X}_\mu \leq \frac{1}{2} \|2(X_\mu - x_\mu^*) + t\dot{X}_\mu\|^2 \leq \frac{\alpha \mathcal{E}_\mu}{\sigma_d^2} \leq \frac{\alpha \mathcal{E}_\mu(t_0)}{\sigma_d^2}$$

By integrating over  $(t_0, t)$ , this gives us

$$t\|X_\mu - x_\mu^*\|^2 - t_0\|x_0 - x_\mu^*\|^2 \leq \frac{\alpha(t - t_0)}{\sigma_d^2} \mathcal{E}_\mu(t_0)$$

By applying the approximation scheme (AS) as  $\mu \rightarrow 0$ , we have for a.e.  $t \geq t_0$  that

$$\|X - x^*\|^2 \leq \frac{\alpha \mathcal{E}(t_0)}{\sigma_d^2} + \|x_0 - x^*\|^2$$



Combining with upper bound of  $\mathcal{E}(t_0)$  in (62), we derive for a.e.  $t \geq t_0$  that

$$\|X(t) - x^*\| \leq C_1 \Delta_0 \quad (63)$$

with factor  $C_1 = \sqrt{\frac{\alpha + (r-1)\sigma_1^2 + \sigma_d^2}{\sigma_d^2}}$ . Here we choose  $\lambda = 2$  to minimize  $C_1$ .

From (61), we know that  $\mathcal{E}_\mu(t)$  is nonincreasing for  $t \geq t_0$ , for all  $\mu > 0$ . By applying (AS) we find that  $\mathcal{E}(t)$  is nonincreasing for a.e.  $t \geq t_0$ . Using the nonnegativity of all three terms in (60) and monotonicity of  $\mathcal{E}(t)$  on a.e.  $t \geq t_0$ , we have for a.e.  $t \geq t_0$  that

$$F(X(t)) - F(x^*) \leq \frac{1}{t^2} \mathcal{E}(t) \leq \frac{1}{t^2} \mathcal{E}(t_0) \leq \frac{C_2}{t^2} \Delta_0^2$$

where factor  $C_2 = 1 + (r-1)\sigma_1^2/\alpha$  is given by (62) with  $\lambda = 2$ , and

$$\|\lambda(X(t) - x^*) + t\dot{X}\|^2 \leq \frac{2\alpha}{\sigma_d^2} \mathcal{E}(t) \leq \frac{2\alpha}{\sigma_d^2} \mathcal{E}(t_0) \leq \frac{2\alpha + \lambda(r-1)\sigma_1^2}{\sigma_d^2} \Delta_0^2$$

Therefore, by triangle inequality and (63),

$$\|\dot{X}(t)\| \leq \frac{1}{t} \|\lambda(X(t) - x^*) + t\dot{X}(t)\| + \frac{1}{t} \lambda \|X(t) - x^*\| \leq \frac{C_3}{t} \Delta_0$$

with factor  $C_3 = \sqrt{\frac{2\alpha + 2(r-1)\sigma_1^2}{\sigma_d^2}} + 2C_1$ . Here we choose  $\lambda = 2$  to minimize  $C_3$ .

From (61), we have

$$\dot{\mathcal{E}}_\mu \leq -(\lambda - 2)t(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) - \frac{(r-1-\lambda)t}{\alpha} \|A\dot{X}_\mu\|^2$$

when  $r = 3$ , we could only choose  $\lambda = 2 = r - 1$  and the right hand side of the inequality above is always zero. However, if we further assume  $r > 3$ , then we could choose  $\lambda = r - 1$  and  $\lambda = 2$  respectively, such that

$$t(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) \leq -\frac{1}{r-3} \dot{\mathcal{E}}_\mu \quad \text{and} \quad t\|\dot{X}_\mu\|^2 \leq -\frac{\alpha}{(r-3)\sigma_d^2} \dot{\mathcal{E}}_\mu$$

By integrating over  $(t_0, \infty)$ , the inequalities above give

$$\int_{t_0}^{\infty} t(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) dt \leq \frac{\mathcal{E}_\mu(t_0)}{r-3} \quad \text{and} \quad \int_{t_0}^{\infty} t\|\dot{X}_\mu(t)\|^2 dt \leq \frac{\alpha \mathcal{E}_\mu(t_0)}{(r-3)\sigma_d^2}$$

By applying (AS) and plugging in (62), we have

$$\int_{t_0}^{\infty} t(F(X(t)) - F(x^*)) dt \leq C_4 \Delta_0^2 \quad \text{and} \quad \int_{t_0}^{\infty} t\|\dot{X}(t)\|^2 dt \leq C_5 \Delta_0^2$$

with factors  $C_4 = \frac{2\alpha + (r-1)^2\sigma_1^2}{2(r-3)\alpha}$  and  $C_5 = \frac{\alpha + (r-1)\sigma_1^2}{(r-3)\sigma_d^2}$ . □

*Proof of Remark 12.* Here we consider a simpler energy functional by taking  $\lambda = 0$  in  $\mathcal{E}_\mu(t)$

$$\mathcal{U}_\mu(t) = t^2(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) + \frac{t^2}{2\alpha} \|A\dot{X}_\mu(t)\|^2$$

Then by taking its time derivative, we have

$$\dot{\mathcal{U}}_\mu = 2t(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) + t^2 \nabla F_\mu(X_\mu(t))^\top \dot{X}_\mu + \frac{t}{\alpha} \|A\dot{X}_\mu\|^2 + t^2 \dot{X}_\mu^\top \left(\frac{1}{\alpha} A^\top A\right) \ddot{X}_\mu$$

By substituting  $A^\top A \ddot{X}_\mu$  with  $-\frac{r}{t} A^\top A \dot{X}_\mu - \alpha \nabla F_\mu(X_\mu(t))$ ,

$$\dot{\mathcal{U}}_\mu = 2t(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) - \frac{(r-1)t}{\alpha} \|A\dot{X}_\mu\|^2 \leq 2t(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*))$$

Combining with Theorem 11, we know that function  $\mathcal{U}_\mu(t) - \int_{t_0}^t 2t(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*))dt$  is nonincreasing and bounded, therefore implying that it is convergent. Consequently, by approximation scheme (AS) and bounded integration formula in Theorem 11, we know that  $\mathcal{U}(t)$  is also convergent as  $t \rightarrow \infty$  for all  $r > 3$ .

By taking integration of  $\frac{1}{t}\mathcal{U}_\mu(t)$  over  $(t_0, \infty)$ , we directly have

$$\int_{t_0}^{\infty} \frac{1}{t}\mathcal{U}_\mu(t)dt = \int_{t_0}^{\infty} t(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*))dt + \frac{1}{2} \int_{t_0}^{\infty} t\|\dot{X}_\mu(t)\|^2 dt$$

Applying (AS) and Theorem 11 implies that

$$\int_{t_0}^{\infty} \frac{1}{t}\mathcal{U}(t)dt \leq (C_4 + C_5/2)\Delta_0^2$$

Combining with our knowledge that  $\int_{t_0}^{\infty} \frac{1}{t}dt = \infty$  and  $\mathcal{U}(t)$  is convergent when  $t \rightarrow \infty$ , we can derive  $\lim_{t \rightarrow \infty} \mathcal{U}(t) = 0$ . This proves that  $F(X(t)) - F(x^*) = o(t^{-2})$  and  $\|\dot{X}(t)\| = o(t^{-1})$  as  $t \rightarrow \infty$ .  $\square$

*Proof of Theorem 13.* The energy functional we used in Theorem 11 is no longer applicable, because we can not find  $\lambda$  satisfying  $\lambda - 2 \geq 0$  and  $r - 1 - \lambda \geq 0$  simultaneously when  $0 < r < 3$ . Here we consider a new energy functional for the Moreau-Yosida approximation

$$\mathcal{E}_\mu(t) = t^2(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) + \frac{1}{2\alpha} \left\| \frac{2r}{3}A(X_\mu(t) - x_\mu^*) + tA\dot{X}_\mu(t) \right\|^2 + \frac{r(3-r)}{9\alpha} \|A(X_\mu(t) - x_\mu^*)\|^2 \quad (64)$$

By taking its time derivative, we have

$$\begin{aligned} \dot{\mathcal{E}}_\mu &= 2t(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) + t^2 \nabla F_\mu(X_\mu(t))^\top \dot{X}_\mu + \left( \frac{2r}{3}(X_\mu - x_\mu^*) + t\dot{X}_\mu \right)^\top \left( \frac{1}{\alpha} A^\top A \right) \left( \left( \frac{2r}{3} + 1 \right) \dot{X}_\mu + t\ddot{X}_\mu \right) \\ &\quad + \frac{2r(3-r)}{9} (X_\mu - x_\mu^*)^\top \left( \frac{1}{\alpha} A^\top A \right) \dot{X}_\mu \end{aligned}$$

By substituting  $tA^\top A\ddot{X}_\mu$  with  $-rA^\top A\dot{X}_\mu - \alpha t \nabla F_\mu(X_\mu)$  and applying the convexity of  $F_\mu$ , we have

$$\dot{\mathcal{E}}_\mu \leq \frac{2(3-r)}{3} t(F_\mu(X_\mu(t)) - F_\mu(x_\mu^*)) + \frac{4r(3-r)}{9\alpha} (X_\mu - x_\mu^*)^\top A^\top A\dot{X}_\mu + \frac{3-r}{3\alpha} t \|A\dot{X}_\mu\|^2$$

Although this energy functional does not have nonnegative derivative, there is a special relationship between it and its derivative. We notice that

$$\dot{\mathcal{E}}_\mu - \frac{2(3-r)}{3t} \mathcal{E}_\mu \leq -\frac{2r(3-r)(3+r)}{27\alpha t} \|A(X_\mu - x_\mu^*)\|^2 \leq 0$$

This implies that, for  $\mathcal{H}_\mu(t) := t^{-\frac{2(3-r)}{3}} \mathcal{E}_\mu(t)$ , for all  $t \geq t_0$ ,

$$\dot{\mathcal{H}}_\mu = t^{-\frac{2(3-r)}{3}} \cdot \left( \dot{\mathcal{E}}_\mu - \frac{2(3-r)}{3t} \mathcal{E}_\mu \right) \leq 0$$

Therefore,  $\mathcal{H}_\mu(t)$  is nonincreasing over  $t \geq t_0$ , for all  $\mu > 0$ . By making similar definition as  $\mathcal{H}(t) := t^{-\frac{2(3-r)}{3}} \mathcal{E}(t)$  and applying the approximation scheme, we have that  $\mathcal{H}(t)$  is nonincreasing for a.e.  $t \geq t_0$ . At time  $t_0$ ,

$$\mathcal{H}(t_0) \leq t_0^{-\frac{2(3-r)}{3}} \cdot \left( 1 + \frac{r(3+r)}{9\alpha} \sigma_1^2 \right) \Delta_0^2$$

By the nonnegativity of all terms in (64) and the monotonicity of  $\mathcal{H}(t)$ , we have for a.e.  $t \geq t_0$  that

$$F(X(t)) - F(x^*) \leq \frac{1}{t^{\frac{2r}{3}}} \mathcal{H}(t) \leq \frac{1}{t^{\frac{2r}{3}}} \mathcal{H}(t_0) \leq \frac{C_6 t_0^{-\frac{2(3-r)}{3}} \Delta_0^2}{t^{\frac{2r}{3}}}$$

with factor  $C_6 = 1 + \frac{r(3+r)\sigma_1^2}{9\alpha}$ .

Similarly, we have for a.e.  $t \geq t_0$  that

$$\left\| \frac{2r}{3}(X(t) - x^*) + t\dot{X} \right\|^2 \leq \frac{2\alpha}{\sigma_d^2} t^{\frac{2(3-r)}{3}} \mathcal{H}(t) \leq \frac{2\alpha}{\sigma_d^2} t^{\frac{2(3-r)}{3}} \mathcal{H}(t_0) \leq \frac{2\alpha C_6 t_0^{-\frac{2(3-r)}{3}}}{\sigma_d^2} t^{\frac{2(3-r)}{3}} \Delta_0^2$$

If we also assume the trajectory  $\{X(t)\}_{t \geq t_0}$  is bounded, then by adopting the same interpretation as in Theorem 11, there exists some positive factor  $C_0$  such that, for a.e.  $t \geq t_0$ ,  $\|X(t) - x^*\| \leq C_0 \Delta_0$ . Then triangle inequality gives us, for a.e.  $t \geq t_0$ , that

$$\|\dot{X}\| \leq \frac{1}{t} \left\| \frac{2r}{3}(X(t) - x^*) + t\dot{X} \right\| + \frac{2r}{3t} \|X(t) - x^*\| \leq \frac{C_7 t_0^{-\frac{3-r}{3}} \Delta_0}{t^{\frac{r}{3}}}$$

with factor  $C_7 = \sqrt{\frac{2\alpha C_6}{\sigma_d^2}} + \frac{2r}{3} C_0$ . □

990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044