
A General Continuous-Time Formulation of Stochastic ADMM and Its Variants

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Stochastic versions of alternating direction method of multiplier (ADMM) and its
2 variants play a key role in many modern large-scale machine learning problems. In
3 this work, we introduce a unified algorithmic framework called generalized stochas-
4 tic ADMM and investigate it via a continuous-time analysis. The generalized
5 framework widely include many stochastic ADMM variants including standard,
6 linearized and gradient-based ADMM. Our continuous-time analysis provides with
7 new insights on stochastic ADMM and variants. We show that its dynamics is
8 approximated by a stochastic differential equations with small noise parameters.
9 We rigorously proved that under some proper scaling, the trajectory of stochastic
10 ADMM weakly converges to the trajectory of the stochastic differential equation.
11 Our analysis also provides a simple proof of the empirical rule that the relaxation
12 parameter must be between 0 and 2.

1 Introduction

13
14 For modern industrial scale machine learning problems, with the massive amount of data that are
15 not only extremely large but often stored or even collected in a distributed manner, stochastic
16 first-order methods almost become one of the default choices due to its excellent performance for
17 online streaming large-scale dataset. Stochastic version of alternating direction method of multiplier
18 (ADMM) algorithms is a popular approach to handle this distributed setting, especially for the
19 regularized empirical risk minimization. Consider the following stochastic optimization problem:

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad V(x) := f(x) + g(Ax) = \mathbb{E}_\xi f(x, \xi) + g(Ax), \quad (1)$$

20 where $f(x) = \mathbb{E}_\xi f(x, \xi)$ with $f(x, \xi)$ as the loss incurred on a sample ξ , $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$,
21 $g : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$, $A \in \mathbb{R}^{m \times d}$, and both f and g are convex and differentiable. The alternating
22 direction method of multiplier (ADMM) [2] rewrites (1) as a constrained optimization problem

$$\begin{aligned} & \underset{x \in \mathbb{R}^d, z \in \mathbb{R}^m}{\text{minimize}} && \mathbb{E}_\xi f(x, \xi) + g(z) \\ & \text{subject to} && Ax - z = 0. \end{aligned} \quad (2)$$

23 Here and throughout the rest of the paper, we overload f to ease the notation, i.e. we adopt the
24 two-argument $f(\cdot, \xi)$ for the stochastic instance and one-argument $f(\cdot)$ for its expectation. Note
25 the classical setting of linear constraint $Ax + Bz = c$ can be reformulated as $z = Ax$ by a linear
26 transformation when B is invertible. In the batch learning setting, $f(x)$ can be approximated by the
27 empirical risk $f_{emp} = \frac{1}{N} \sum_{i=1}^N f(x, \xi_i)$. However, the computation cost for minimizing f_{emp} with a
28 large amount of samples is significantly high and the efficiently is reasonably limited given the time
29 and resource constraints. In the stochastic setting, at each step x_k is updated based on a small batch
30 of samples (or even one) ξ_k instead of a large batch or full batch.

31 **General Formulation for Stochastic ADMM** Introducing stochasticity to ADMM [2] is in parallel
32 to incorporating noise into gradient descent [55, 52, 42]. At iterate $k + 1$, a sample ξ_{k+1} is randomly
Submitted to 34th Conference on Neural Information Processing Systems (NeurIPS 2020). Do not distribute.

33 drawn as an independent realization of ξ . Analogous to stochastic gradient descent, our **stochastic**
 34 **ADMM (sADMM)** introduced in [52, 42] performs the following updates:

$$x_{k+1} = \operatorname{argmin}_x \left\{ f(x, \xi_{k+1}) + \frac{\rho}{2} \|Ax - z_k + u_k\|_2^2 \right\}, \quad (3a)$$

$$z_{k+1} = \operatorname{argmin}_z \left\{ g(z) + \frac{\rho}{2} \|\alpha Ax_{k+1} + (1 - \alpha)z_k - z + u_k\|_2^2 \right\}, \quad (3b)$$

$$u_{k+1} = u_k + (\alpha Ax_{k+1} + (1 - \alpha)z_k - z_{k+1}). \quad (3c)$$

35 $\rho > 0$ is called the penalty parameter [2]. Here $\alpha \in (0, 2)$ is introduced as a *relaxation parameter*
 36 [2]. The algorithm is splitted into different regimes of α as follows:

- 37 • $\alpha = 1$, corresponding to the **standard stochastic ADMM**;
- 38 • $\alpha > 1$, corresponding to the **over-relaxed stochastic ADMM**;
- 39 • $\alpha < 1$, corresponding to the **under-relaxed stochastic ADMM**.

40 Historically, the above relaxation schemes were introduced and analyzed in [11], and experiments
 41 in [10, 12] suggest that over-relaxed regime of $\alpha > 1$ can improve the convergence rate. The
 42 acceleration phenomenon in over-relaxed ADMM has been discussed from the continuous perspective
 43 in deterministic setting [63, 16], where relaxed ADMM algorithm accelerates by a factor of $\alpha \in (1, 2)$
 44 in its convergence rate.

45 Since the emerging of ADMM, many variants have been introduced recently for solving a variety of
 46 optimization tasks. We focus on two stochastic ADMM variants to cater the need in application:

47 (i) In the linearized ADMM [20, 33, 43], (4a) replaces (3a) where

$$x_{k+1} := \operatorname{argmin}_x \left\{ f(x, \xi_{k+1}) + \frac{\tau}{2} \left\| x - \left(x_k - \frac{\rho}{\tau} A^\top (Ax_k - z_k + u_k) \right) \right\|_2^2 \right\}. \quad (4a)$$

48 In words, the augmented Lagrangian function is approximated by linearizing the quadratic term
 49 of x in (3a) plus the addition of a proximal term $\frac{\tau}{2} \|x - x_k\|_2^2$.

50 (ii) In the gradient-based ADMM [6, 54, 7, 32], (5a) replaces (3a) where

$$x_{k+1} := x_k - \frac{1}{\tau} \left(f'(x_k, \xi_{k+1}) + \rho A^\top (Ax_k - z_k + u_k) \right). \quad (5a)$$

51 In words, in lieu to solving x -subproblem (3a) accurately, we apply one step of gradient descent
 52 with stepsize $1/\tau$.

53 In this paper we formulate a general scheme, called **Generalized stochastic ADMM (G-sADMM)**,
 54 to accommodate and unify all these variants of stochastic ADMM:

$$x_{k+1} = \operatorname{argmin}_x \hat{\mathcal{L}}_{k+1}(x, z_k, u_k), \quad (6a)$$

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

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

55 where the objective $\hat{\mathcal{L}}_{k+1}$ in (6a) for x -subproblem is defined as

$$\begin{aligned} \hat{\mathcal{L}}_{k+1}(x, z_k, u_k) &= (1 - \omega_1) f(x, \xi_{k+1}) + \omega_1 f'(x_k, \xi_{k+1})(x - x_k) \\ &\quad + (1 - \omega) \cdot \frac{\rho}{2} \|Ax - z_k + u_k\|_2^2 + \omega \rho A^\top (Ax_k - z_k + u_k)(x - x_k) + \frac{\tau}{2} \|x - x_k\|_2^2, \end{aligned} \quad (7)$$

56 with parameter $\rho, \tau \geq 0$, $\alpha \in (0, 2)$ and explicitness parameters $\omega_1, \omega \in [0, 1]$.

57 **Main Results** We study the continuous-time limit of the stochastic sequence $\{x_k\}$ defined in (6a)
 58 as $\rho \rightarrow \infty$. Following the recent line of work França et al. [14, 15] and Yuan et al. [63], we adopt
 59 a (sped-up) time rescaling of ρ^{-1} , namely the stepsize for (stochastic) ADMM, and hence we are
 60 in the small-stepsize regime as in (stochastic) gradient descent. Also, the proximal parameter τ in
 61 (7) grows linearly in ρ i.e. $\tau/\rho \rightarrow c$. In our setting of continuous limit theory, we focus on a fixed
 62 interval by T (referred to as "time"), corresponding to $\lfloor \rho T \rfloor$ steps (referred to as "step") in discrete
 63 time. Our main results can be concluded as follows:

64 (i) As $\rho \rightarrow \infty$ the sequence $\{x_{\lfloor \rho t \rfloor}\}$ for $1 \leq t \leq T$ admits a continuous limit $\{X_t : t \in [0, T]\}$, a
 65 stochastic process that solves the following stochastic differential equation (SDE):

$$\left[cI + \left(\frac{1}{\alpha} - \omega \right) A^\top A \right] dX_t = -\nabla V(X_t) dt + \sqrt{\frac{1}{\rho}} \sigma(X_t) dW_t, \quad (8)$$

66 where we recall from (1) that $V(x) = f(x) + g(Ax)$. The right hand of (8) consists of a
 67 deterministic part and a stochastic part, where $\sigma(x)$ is some diffusion matrix defined later
 68 in (9) and W_t is the standard Brownian motion in \mathbb{R}^d . The continuous-limit approximation
 69 is rigorously characterized in weak convergence in the sense that applying a test function φ
 70 of a certain class, the expectation of $\varphi(x_k) - \varphi(X_{k/\rho})$ is uniformly convergent to zero in
 71 $k \in [0, \lfloor \rho T \rfloor]$ when ρ tends to infinity as stated in Theorem 1.2.3. The above SDE carries the
 72 small parameter $\sqrt{1/\rho}$ in its diffusion term and is also called *stochastic modified equation*
 73 (SME) [31], due to the historical reason in numerical analysis [38, 25].

74 (ii) We provide a continuous-time explanation of the effect of relaxation parameter $\alpha \in (0, 2)$ in
 75 stochastic ADMM, which is the first among stochastic ADMM work to our best knowledge. In
 76 light of our weak convergence of the discrete algorithm x_k , the residual $r_k = Ax_k - z_k$ roughly
 77 satisfies $r_{k+1} \approx (1 - \alpha)r_k$ and converges to 0 at a geometric rate $|1 - \alpha| < 1$. The rigorous
 78 statement is in Corollary 2.

79 **Contributions** Our contribution to this paper is the first continuous-time analysis of stochastic
 80 ADMM. First, we present a unified stochastic differential equation as the continuous-time model
 81 of variants of stochastic ADMM (standard ADMM, linearized ADMM, gradient-based ADMM)
 82 in the regime of large ρ under weak convergence. Second, we characterize the time evolution of
 83 $\text{std}(x_k)$ and $\text{std}(z_k)$ in their continuous-time counterparts and explicitly show that the standard
 84 deviation of stochastic ADMM variants has the scaling $\rho^{-1/2}$ in stochastic ADMM. Finally, our
 85 theoretical analysis of continuous model provides a simple justification of effect of relaxation
 86 parameter $\alpha \in (0, 2)$, and further sheds light on the principled choices of parameters c, α, ω .

87 2 Related Work

88 **ADMM:** ADMM is a widely used algorithm for solving problems with separable structures in
 89 machine learning, statistics, control etc. ADMM has a close connection with operator-splitting
 90 methods [8, 44, 18, 19]. But ADMM comes back to popularity due to several works like [5, 21, 56]
 91 and the highly influential survey paper [2]. Numerous modern machine learning applications are
 92 inspired by ADMM, for example, [40, 35, 62, 36, 51, 50, 45].

93 **Variants of ADMM** Linearized ADMM and gradient-based ADMM are widely used variants of
 94 ADMM. Linearized ADMM has been studied extensively, for example, in [3, 9, 22, 33, 34, 60–
 95 62, 65, 43]; gradient-based ADMM has been studied extensively too, for example, in [6, 54, 7, 32].
 96 In this work, we present a general formulation for relaxed, linearized and gradient-based ADMM,
 97 and extended all of them to stochastic optimization.

98 **Relaxation scheme for ADMM** There are several important works studying the relaxation scheme
 99 of ADMM [10, 9, 11, 12] and propose to choose $\alpha \in (0, 2)$ with empirical suggestion for over
 100 relaxation $1 < \alpha < 2$. In this work, we give a simple and elegant proof of the theoretical reason that
 101 $\alpha \in (0, 2)$.

102 **Continuous model for ADMM:** The recent work in [14, 16] establishes the first deterministic
 103 continuous-time model for standard ADMM in the form of ordinary differential equation (ODE) for
 104 the smooth ADMM and [15, 63] extends its work to non-smooth ADMM, using the tool of differential
 105 inclusion, which is motivated by [53, 41].

106 **Stochastic and online ADMM:** The use of stochastic and online techniques for ADMM has recently
 107 drawn a lot of interest. [55] first proposed the online ADMM, which learns from only one sample
 108 (or a small mini-batch) at a time. [42, 52] proposed the variants of stochastic ADMM to attack the
 109 difficult nonlinear optimization problem inherent in $f(x, \xi)$ by linearization. Very recently, further
 110 accelerated algorithms for the stochastic ADMM have been developed in [67, 23].

¹We refer the readers to Appendix A for backgrounds of SDE and the notion of weak convergence rate.

²We warn the reader that this is towards an orthogonal direction to the convergence of x_k to the optimizer x_* of the objective function as $k \rightarrow \infty$ for a fixed ρ ,

³Appendix A is a summary of the background of the stochastic differential equation and the weak approximation, and the theory on the stochastic modified equation for interested readers.

111 **Continuous models and analysis for stochastic optimization:** The work in [49] is one seminal
112 work of using continuous-time dynamical system to analyze discrete algorithms for optimization
113 such as Nesterov’s accelerated gradient method, with important extension to high resolution and
114 symplectic structure in [47, 48, 24], to variational perspective in [59, 58]. The mathematical analysis
115 of continuous formulation for stochastic optimization algorithms has become an important trend in
116 recent years. We select an under-represented list out of numerous works: analyzing SGD from the
117 perspective of Langevin MCMC [4, 66], analyzing stochastic mirror-descent [68], analyzing the SGD
118 for the online algorithm of specific statistical models [28, 26, 30].

119 **Stochastic modified equation:** The mathematical connection between the stochastic gradient descent
120 (SGD) and stochastic modified equation (SME) has been insightfully discovered in [29, 31]. This
121 SME technique, originally arising from the numerical analysis of SDE [38, 25], has become the
122 major mathematical tool for stochastic or online algorithms. The idea of using optimal control for
123 stochastic continuous models to find adaptive parameters like step-size for stochastic optimization is
124 used in [29] to controlled adaptive step size and momentum, and in [1] to improve the control on
125 momentum, in [27] to choose batch size.

126 3 Weak Approximation to Stochastic ADMM

127 In this section, we show the weak approximation to the generalized family of stochastic ADMM
128 variants (**G-sADMM**) (6). In the following sections, we define $\epsilon := \rho^{-1}$ to serve as an analog of
129 stepsize for (stochastic) ADMM. The **G-sADMM** scheme (6) contains the relaxation parameter α ,
130 the parameter $c = \tau/\rho$ and the implicitness parameters ω, ω_1 . This scheme (6) is very general and
131 includes many existing variants as follows. It recovers all deterministic variants of ADMM for the
132 setting $f(x, \xi) \equiv f(x)$. If $\omega_1 = \omega = \tau = 0$, then (6) is the standard stochastic ADMM (**sADMM**)
133 or online ADMM [42, 67]. For the setting of parameters $\omega_1 = 0, \omega = 1$ and $c > 0$ it becomes
134 stochastic version of the linearized ADMM for acceleration [9, 22, 65, 43]. If $\omega_1 = \omega = 1$ and $c > 0$,
135 this scheme is the stochastic version of the gradient-based ADMM [6, 54, 7, 32]. In addition, the
136 stochastic ADMM considered in [42] for non-smooth function corresponds to the setting here with
137 $\alpha = 1, \omega_1 = 1, \omega = 0$ and $\tau = \tau_k \propto \sqrt{k}$.

138 **Notations and Assumptions** We use $\|\cdot\|$ to denote the Euclidean two norm if the subscript is not
139 specified. And all vectors are referred to as column vectors. $f'(x, \xi), g'(z)$ and $f''(x, \xi), g''(z)$ refer
140 to the first (gradient) and second (Hessian) derivatives w.r.t. x . The first assumption is **Assumption I:**
141 $f(x), g$ and for each $\xi, f(x, \xi)$, are closed proper convex functions; A has full column rank.

142 Let \mathcal{F} be the set of functions of at most polynomial growth. To apply the SME theory, we need the
143 following **Assumptions II** on the smoothness:

- 144 (i) The second order derivative f'', g'' are uniformly bounded in x , and almost surely in ξ for
145 $f''(x, \xi)$. $\mathbb{E} \|f'(x, \xi)\|_2^2$ is uniformly bounded in x .
- 146 (ii) $f(x), f(x, \xi), g(x)$ and the partial derivatives up to order 5 belong to \mathcal{F} .
- 147 (iii) $f'(x)$ and $f'(x, \xi)$ satisfy a uniform growth condition: $\|f'(x)\| + \|f'(x, \xi)\| \leq C(1 + \|x\|)$
148 for a constant C independent of ξ .

149 **Main Theorem** Given the noisy gradient $f'(x, \xi)$ and its expectation $f'(x) = \mathbb{E}_\xi f'(x, \xi)$, we
150 define the following matrix $\sigma(x) \in \mathbb{R}^{d \times d}$ by

$$\Sigma(x) = \sigma(x)\sigma(x)^\top = \mathbb{E}_\xi \left[(f'(x, \xi) - f'(x))(f'(x, \xi) - f'(x))^\top \right]. \quad (9)$$

151 **Theorem 1** (SME for G-sADMM). *Let $\alpha \in (0, 2)$, $\omega_1, \omega \in \{0, 1\}$ and $c = \tau/\rho \geq 0$. Let*
152 *$\epsilon = \rho^{-1} \in (0, 1)$. $\{x_k\}$ denote the sequence of stochastic ADMM (6) with the initial choice*
153 *$z_0 = Ax_0$. Define X_t as a stochastic process satisfying the SDE*

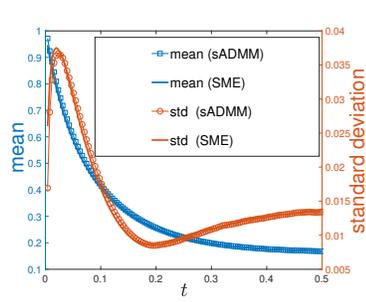
$$\widehat{M}dX_t = -\nabla V(X_t)dt + \sqrt{\epsilon}\sigma(X_t)dW_t \quad (10)$$

154 *where the matrix*

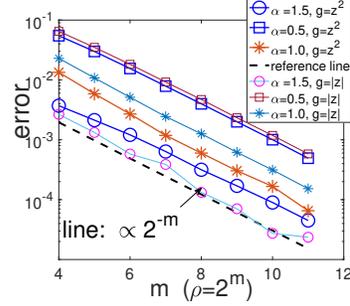
$$\widehat{M} := cI + \left(\frac{1}{\alpha} - \omega \right) A^\top A. \quad (11)$$

155 *Then we have $x_k \rightarrow X_{k\epsilon}$ in weak convergence of order 1, with the following precise meaning: For*
156 *any time interval $T > 0$ and for any test function φ such that φ and its partial derivatives up to order*
157 *4 belong to \mathcal{F} , there exists a constant C such that*

$$|\mathbb{E} \varphi(X_{k\epsilon}) - \mathbb{E} \varphi(x_k)| \leq C\epsilon, \quad \forall k \leq \lfloor T/\epsilon \rfloor. \quad (12)$$



(a) The expectation (left axis) and standard deviation (right axis) of x_k and X_t . $\epsilon = 2^{-7}$. $\alpha = 1.5$.



(b) The weak convergence error err_m versus m for various α and ℓ_2, ℓ_1 regularization g .

Figure 1: The match between the stochastic ADMM and the SME and the verification of the first order weak approximation. The result is based on the average of 10^5 independent runs. The step size $\epsilon = 1/\rho = 2^{-m}T$ and $T = 0.5$ in (b). The details can be referred to in Appendix C

158 The proof is in the next section. Here we highlight one useful corollary resulted from our proof of the
159 theorem.

160 **Corollary 2.** For generalized stochastic ADMM, when $\rho \rightarrow \infty$, the residual r_k satisfied the relation:
161 $r_{k+1} = (1 - \alpha)r_k$. From this relation, we can derive the condition on the relaxation parameter α
162 for the convergence $r_k \rightarrow 0$ as $k \rightarrow \infty$: $|1 - \alpha| < 1$, which matches the empirical range $\alpha \in (0, 2)$
163 proposed for standard ADMM with relaxation Boyd et al. [2].

164 As an illustration of Theorem I on the first order weak convergence, Figure 1b verifies this rate for
165 various α and even different g . To show that the SME does not only provide the expectation of the
166 solution, but also provides the fluctuation of the numerical solution x_k for any given ϵ , Figure 1a
167 plots the match of mean and standard deviation of X_t from the SME versus the x_k from the sADMM.
168 One corollary of Theorem I is that the standard deviation of the stochastic ADMM x_k is $O(\sqrt{\epsilon})$ and
169 the rescaled two standard deviations $\epsilon^{-1/2} \text{std}(x_k)$ and $\epsilon^{-1/2} \text{std}(X_{k\epsilon})$ are close as the functions
170 of the time $t_k = k\epsilon$. Other heuristics from Theorem I is about the continuous dynamics of the z_k
171 variable and the residual r_k , which is discussed in Appendix B

172 4 Proof of Main Results

173 *Proof of Theorem 7* The ADMM scheme is the iteration of the triplet (x, z, λ) where $\lambda := \epsilon u$. In our
174 proof, we shall show this triplet iteration can be effectively reduced to the iteration of x variable only
175 in the form of $x_{k+1} = x_k + \epsilon \mathcal{A}(\epsilon, x_k, \xi_{k+1})$. The main approach for this reduction is the detailed
176 analysis of the residual $r_k = x_k - Az_k$ below. This *one step difference* $x_{k+1} - x_k$ will play a central
177 role in the theory of weak convergence, Milstein's theorem (Theorem 5 in Appendix A [37]).

178 For notational ease, we drop the random variable ξ_{k+1} in the scheme (6); the readers bear in mind
179 that f and its derivatives do involve ξ and all conclusions hold almost surely for ξ .

180 The optimality conditions for the scheme (6) are

$$\begin{aligned} \omega_1 \epsilon f'(x_k) + (1 - \omega_1) \epsilon f'(x_{k+1}) + \epsilon A^\top \lambda_k \\ + A^\top (\omega Ax_k + (1 - \omega)Ax_{k+1} - z_k) + c(x_{k+1} - x_k) = 0 \end{aligned} \quad (13a)$$

$$\epsilon g'(z_{k+1}) = \epsilon \lambda_k + \alpha Ax_{k+1} + (1 - \alpha)z_k - z_{k+1} \quad (13b)$$

$$\epsilon \lambda_{k+1} = \epsilon \lambda_k + \alpha Ax_{k+1} + (1 - \alpha)z_k - z_{k+1} \quad (13c)$$

181 Note that due to (13b) and (13c), the last condition (13c) can be replaced by $\lambda_{k+1} = g'(z_{k+1})$.
182 So, without loss of generality, one can assume that $\lambda_k \equiv g'(z_k)$ for any integer k . The optimality
183 conditions (13) now can be written only in the pair (x, z) :

$$\begin{aligned} \omega_1 \epsilon f'(x_k) + (1 - \omega_1) \epsilon f'(x_{k+1}) + \epsilon A^\top g'(y_k) \\ + A^\top (\omega Ax_k + (1 - \omega)Ax_{k+1} - z_k) + c(x_{k+1} - x_k) = 0 \end{aligned} \quad (14a)$$

$$\epsilon g'(z_{k+1}) - \epsilon g'(z_k) = \alpha Ax_{k+1} + (1 - \alpha)z_k - z_{k+1} \quad (14b)$$

184 We think of (x_{k+1}, z_{k+1}) in (14) as the function of x_k, z_k and ϵ and seek the asymptotic expansion
185 of this function in the regime of small ϵ , by using Dominant Balance technique (see, e.g. [57]) with

186 the Taylor expansion for $f'(x_{k+1})$ and $g'(z_{k+1})$ around x_k and z_k , respectively. The expansion of
 187 the one step difference up to the order ϵ is

$$x_{k+1} - x_k = -M^{-1}A^\top r_k - \epsilon M^{-1} [f'(x_k) + A^\top g'(z_k)] - \epsilon N_k A^\top r_k + \epsilon^2 c_k, \quad (15a)$$

$$z_{k+1} - z_k = (1 - \epsilon g''(z_k))\alpha(A(x_{k+1} - x_k) + r_k) + \epsilon^2 c'_k, \quad (15b)$$

188 where r_k is the *residual*

$$r_k := Ax_k - z_k \quad (16)$$

189 and the matrix M is

$$M := cI + (1 - \omega)A^\top A. \quad (17)$$

190 I is the identity matrix. N_k is a matrix depending on the hessian $f''(x_k)$ and M . c_k and c'_k are
 191 uniformly bounded in k (and in random element ξ) due to the bound of functions f'' and g'' in
 192 **Assumption II**. Throughout the rest of the paper, we shall use the notation $\mathcal{O}(\epsilon^p)$ to denote these
 193 terms bounded by a (nonrandom) constant multiplier of ϵ^p , for $p = 1, 2, \dots$

194 Equation (15) clearly shows that the dynamics of the pair (x_{k+1}, z_{k+1}) depends on the residual
 195 $r_k = Ax_k - z_k$ in the previous step. In the limit of $\epsilon = 0$, the one step difference (i.e. *truncation error*)
 196 $x_{k+1} - x_k \rightarrow -M^{-1}A^\top r_k$ which does not vanish unless r_k is exact zero. We next turn to the analysis
 197 of the residual first. (15b) together with (15a) implies $z_{k+1} - z_k = -\alpha AM^{-1}A^\top r_k + \alpha r_k + \mathcal{O}(\epsilon)$
 198 by absorbing ϵ order terms into $\mathcal{O}(\epsilon)$. By using (14b) for Ax_{k+1} , we have the recursive relation for
 199 residual r_{k+1} :

$$\begin{aligned} r_{k+1} &= Ax_{k+1} - z_{k+1} = \left(\frac{1}{\alpha} - 1\right)(z_{k+1} - z_k) + \frac{\epsilon}{\alpha}(g'(z_{k+1}) - g'(z_k)) \\ &= (1 - \alpha)(I - AM^{-1}A^\top)r_k + \mathcal{O}(\epsilon). \end{aligned} \quad (18)$$

200 Since $r_0 = 0$ at the initial step by setting $z_0 = Ax_0$, then (18) shows that $r_1 = Ax_1 - y_1$ become
 201 $\mathcal{O}(\epsilon)$ after one iteration. Next, we have the following important property, particularly with assumption
 202 $\alpha = 1$, we have a stronger result than (18).

203 **Proposition 3.** For any integer $k \geq 0$, if $r_k = \mathcal{O}(\epsilon)$, then

$$r_{k+1} = (1 - \alpha + \epsilon\alpha g''(z_k))(r_k + A(x_{k+1} - x_k)) + \mathcal{O}(\epsilon^3). \quad (19)$$

204 If $\alpha = 1$, equation (19) reduces to the second order in ϵ :

$$r_{k+1} = \epsilon\alpha g''(z_k)(r_k + A(x_{k+1} - x_k)) + \mathcal{O}(\epsilon^3) = \mathcal{O}(\epsilon^2). \quad (20)$$

Proof of Proposition 3 Since $r_k = Ax_k - z_k = \mathcal{O}(\epsilon)$, then the one step difference $x_{k+1} - x_k$ and
 $z_{k+1} - z_k$ are both at order $\mathcal{O}(\epsilon)$ because of (15a) and (15b). We solve $\delta z := z_{k+1} - z_k$ from (14b)
 by linearizing the implicit term $g'(z_{k+1})$ and use the assumption that the third order derivative of g
 exits:

$$\epsilon g''(z_k)\delta z + \epsilon\mathcal{O}((\delta z)^2) + \delta z = \alpha(r_k + A\delta x).$$

205 where $\delta x := x_{k+1} - x_k = \mathcal{O}(\epsilon)$. Then since $\mathcal{O}((\delta z)^2) = \mathcal{O}(\epsilon^2)$, the expansion of $\delta z = z_{k+1} - z_k$
 206 in ϵ is

$$z_{k+1} - z_k = \alpha(1 - \epsilon g''(z_k))(r_k + A(x_{k+1} - x_k)) + \mathcal{O}(\epsilon^3) \quad (21)$$

207 Then

$$\begin{aligned} r_{k+1} &= r_k + A(x_{k+1} - x_k) - (z_{k+1} - z_k) \\ &= \left(1 - \alpha + \epsilon\alpha g''(z_k)\right)(r_k + A(x_{k+1} - x_k)) + \mathcal{O}(\epsilon^3) \end{aligned}$$

208 □

209 We now focus on (15a) for $x_{k+1} - x_k$, which is rewritten

$$\begin{aligned} M(x_{k+1} - x_k) &= -A^\top r_k - \epsilon(f'(x_k) + A^\top g'(z_k)) \\ &\quad + \epsilon^2(1 - \omega_1)f''(x_k)M^{-1}\left(f'(x_k) + A^\top g'(z_k) + \frac{1}{\epsilon}A^\top r_k\right) + \mathcal{O}(\epsilon^3). \end{aligned} \quad (22)$$

210 This expression does not contain the parameter α explicitly, but the residual $r_k = Ax_k - y_k$
 211 significantly depends on α (see Proposition 3). If $\alpha = 1$, then r_k is on the order of ϵ^2 , which leads to

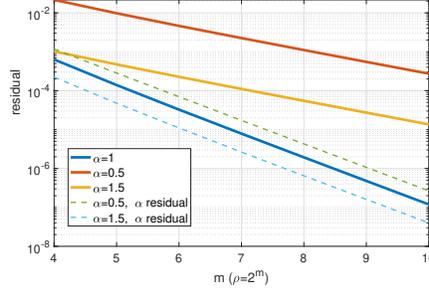


Figure 2: Illustration of the order of the residual r_k and the α -residual r_k^α defined in (23).

212 useful fact that there is no contribution from r_k to the weak approximation of x_k at the order 1. But
 213 for the relaxation case where $\alpha \neq 1$, r_k contains the first order term coming from $z_{k+1} - z_k$.
 214 To obtain a second order for some “residual” for the relaxed scheme where $\alpha \neq 1$, we need a new
 215 definition, α -residual, to account for the gap induced by α . Motivated by (13b), we first define

$$r_{k+1}^\alpha := \alpha A x_{k+1} + (1 - \alpha) z_k - z_{k+1}. \quad (23)$$

216 It is connected to the original residual r_{k+1} and r_k since it is easy to check that

$$r_{k+1}^\alpha = \alpha r_{k+1} + (\alpha - 1)(z_{k+1} - z_k) = \alpha r_k + \alpha A(x_{k+1} - x_k) - (z_{k+1} - z_k). \quad (24)$$

217 Obviously, when $\alpha = 1$, this α -residual r^α is the original residual $r = Ax - y$. In our proof, we need
 218 a modified α -residual, denoted by

$$\widehat{r}_{k+1}^\alpha := \alpha r_k + (\alpha - 1)(z_{k+1} - z_k) \quad (25)$$

We shall show that both r_{k+1}^α and \widehat{r}_{k+1}^α are as small as $\mathcal{O}(\epsilon^2)$:

$$r_{k+1}^\alpha = \mathcal{O}(\epsilon^2) \quad \text{and} \quad \widehat{r}_{k+1}^\alpha = \mathcal{O}(\epsilon^2).$$

219 In fact, by the second equality of (24), (21) becomes $z_{k+1} - z_k = (1 - \epsilon g''(z_k))(r_{k+1}^\alpha + z_{k+1} -$
 220 $z_k) + \mathcal{O}(\epsilon^3)$, so,

$$r_{k+1}^\alpha = \epsilon(1 + \epsilon g''(z_k))g''(z_k)(z_{k+1} - z_k) + \mathcal{O}(\epsilon^3)$$

221 which is $\mathcal{O}(\epsilon^2)$ since $z_{k+1} - z_k = \mathcal{O}(\epsilon)$.

222 The difference between $(z_{k+1} - z_k)$ and $(z_{k+2} - z_{k+1})$, is at the order ϵ^2 due to truncation error of
 223 the central difference scheme. Then we have the conclusion

$$\widehat{r}_{k+1}^\alpha = \alpha r_k + (\alpha - 1)(z_{k+1} - z_k) = \mathcal{O}(\epsilon^2) \quad (26)$$

224 by shifting the subscript k by one.

225 By (25) and the above proposition, we have $r_k = (\frac{1}{\alpha} - 1)(z_{k+1} - z_k) + \mathcal{O}(\epsilon^2)$. Furthermore, due to
 226 (21), $r_k = (\frac{1}{\alpha} - 1)(z_{k+1} - z_k) + \mathcal{O}(\epsilon^2) = (1 - \alpha)(r_k + A(x_{k+1} - x_k)) + \mathcal{O}(\epsilon^2)$ which gives

$$r_k = \left(\frac{1}{\alpha} - 1\right)(z_{k+1} - z_k) + \mathcal{O}(\epsilon^2) = \left(\frac{1}{\alpha} - 1\right)A(x_{k+1} - x_k) + \mathcal{O}(\epsilon^2) \quad (27)$$

227 and it follows $z_{k+1} - z_k = A(x_{k+1} - x_k) + \mathcal{O}(\epsilon^2)$.

228 **Remark 1.** As an illustration, the figure above for a toy example (see numerical example part in
 229 Appendix C) shows that $r_k \sim \mathcal{O}(\epsilon)$ for $\alpha \neq 1$ while $r_k \sim \mathcal{O}(\epsilon^2)$ for $\alpha = 1$ and the α -residual r_k^α is
 230 $\mathcal{O}(\epsilon^2)$ regardless of α (the solid lines are r_k and the dashed lines are r_k^α).

231 With the above preparations for the residual analysis, we now apply Theorem 5 to show the main
 232 conclusion of our theorem. Combining equation (22) and (27), and noting the Taylor expansion of
 233 $g'(z_k)$: $g'(y_k) = g'(Ax_k - r_k) = g'(Ax_k) + \mathcal{O}(\epsilon)$ since $r_k = \mathcal{O}(\epsilon)$, and now putting back the
 234 random element ξ into f' , we have

$$M(x_{k+1} - x_k) = -\epsilon (f'(x_k, \xi_{k+1}) + A^\top g'(Ax_k)) - \left(\frac{1}{\alpha} - 1\right) A^\top A(x_{k+1} - x_k) + \mathcal{O}(\epsilon^2) \quad (28)$$

235 For convenience, introduce the matrix

$$\widehat{M} := M + \frac{1-\alpha}{\alpha} A^\top A = c + \left(\frac{1}{\alpha} - \omega\right) A^\top A. \quad (29)$$

and let $\widehat{x}_k := \widehat{M}x_k$, and $\delta\widehat{x}_{k+1} = \widehat{M}(x_{k+1} - x_k)$ Then

$$\delta\widehat{x} = -\epsilon V'(x, \xi) + \epsilon^2 ((1 - \omega_1) f'' M^{-1} V'(x) - A^\top \theta) + \mathcal{O}(\epsilon^3).$$

236 The final step is to compute the momentums in the Milstein's theorem Theorem 5 as follows.

(i)

$$\mathbb{E}[\delta\widehat{x}] = -\epsilon \mathbb{E} V'(x, \xi) + \mathcal{O}(\epsilon^2) = -\epsilon V'(x) + \mathcal{O}(\epsilon^2) \quad (30)$$

(ii)

$$\begin{aligned} \mathbb{E}[\delta\widehat{x} \delta\widehat{x}^\top] &= \epsilon^2 \mathbb{E} \left([f'(x, \xi) + A^\top g'(x)] [f'(x, \xi)^\top + g'(x)^\top A] \right) + \mathcal{O}(\epsilon^3) \\ &= \epsilon^2 (V'(x) V'(x)^\top) - \epsilon^2 (f'(x) + A^\top g'(x)) (f'(x)^\top + g'(x)^\top A) \\ &\quad + \epsilon^2 \mathbb{E} \left([f'(x, \xi) + A^\top g'(x)] [f'(x, \xi)^\top + g'(x)^\top A] \right) + \mathcal{O}(\epsilon^3) \\ &= \epsilon^2 (V'(x) V'(x)^\top) + \epsilon^2 \mathbb{E} \left[(f'(x, \xi) - f'(x)) (f'(x, \xi) - f'(x))^\top \right] + \mathcal{O}(\epsilon^3) \end{aligned}$$

237 (iii) It is trivial that $\mathbb{E}[\prod_{j=1}^s \delta x_{i_j}] = \mathcal{O}(\epsilon^3)$ for $s \geq 3$ and $i_j = 1, \dots, d$.

238 So, Theorem 1 is proved. □

239

240 *Proof of Corollary 2* Based on Proposition 3 $r_{k+1} = (1-\alpha)(r_k + A(x_{k+1} - x_k)) + \epsilon \alpha g''(z_k)(r_k +$
 241 $A(x_{k+1} - x_k)) + \mathcal{O}(\epsilon^3)$. Since g'' are bounded and $x_{k+1} - x_k = \mathcal{O}(\epsilon)$, then $r_{k+1} = (1-\alpha)r_k + \mathcal{O}(\epsilon)$,
 242 with the leading term $(1-\alpha)r_k$. Therefore, when $\epsilon \rightarrow 0$, we need $|1-\alpha| < 1$ for $|r_k|$ to converge
 243 to zero when $k \rightarrow \infty$. □

244 5 Conclusion and Discussion

245 In this paper, we have developed the stochastic continuous dynamics in the form of stochastic
 246 modified equation (SME) to analyze the dynamics of a general family of stochastic ADMM, including
 247 the standard, linearized and gradient-based ADMM with relaxation α . Our continuous model
 248 (10) provides a unified framework to describe the dynamics of stochastic ADMM algorithms and
 249 particularly is able to quantify the fluctuation effect for a large penalty parameter ρ . Our analysis
 250 here generalize the existing works [14, 63] in the form of ordinary differential equation, but our
 251 results intrinsically encodes the impact of the noise in the stochastic ADMM. As the first order
 252 approximation to the stochastic ADMM trajectory, the solution to the continuous model can precisely
 253 describe the mean and the standard deviation (fluctuation) in stochastic trajectory.

254 One distinctive feature between the ADMM and its stochastic or online variants is highly similar to
 255 that between gradient descent and stochastic gradient descent [39, 29, 46]: there exists a transition
 256 time t_* after which the fluctuation (with a typical scale $\rho^{-1/2}$) starts to dominate the “drift” term,
 257 which means the traditional acceleration methods in deterministic case will fail to perform well,
 258 despite of their prominent performance in the early stage of training when the drift suppresses the
 259 noise. In our Appendix D, we briefly discuss a few new adaptive strategies such as for parameters
 260 ρ_t, τ_t , the batch size B_t and even relaxation parameters α_t , via the principle of control theory for
 261 the continuous model we obtained. For example, a large α (over-relaxation) is preferred before the
 262 transition time while a small α (under-relaxation) may be preferred later to help reduce the variance.

263 About the possible further development in theory, we note that our generalized family of ADMM
 264 methods is still restricted to the linearized ADMM and the gradient-based ADMM. Recently, there
 265 emerge many new efficient acceleration methods for stochastic ADMM such as [45, 67, 23, 42]. A
 266 potential future work is to extend our stochastic analysis to these new schemes for the understanding
 267 at the continuous level, even with non-smooth functions.

268 **Broader Impact**

269 Our work is a theoretic exploration of the existing and popular stochastic optimization with the
270 ADMM. The continuous viewpoint to understand many machine learning algorithms is a powerful
271 tool with a unifying conciseness and elegance, to complement the pure discrete-level analysis. It
272 offers a unique bridge in theory between the applied mathematics and machine learning. Such a
273 fusion between discrete world and continuum world can boost the beneficial communication between
274 the algorithmic community and the applied math community. For ethical aspects and future societal
275 consequences, the positive part is that large amount of machine learning algorithms are trained using
276 ADMM and stochastic ADMM, our work would shed light for algorithm designing from theoretical
277 point of view, that might lead to more efficient training algorithm that save the energy consumed in
278 training from parallel search of training hyper-parameters. Due to the theoretical nature of our work,
279 there are not much negative impact, the only related negative impact might come from algorithmic
280 automation that might lead to unemployment of workers.

281 **References**

- 282 [1] Wangpeng An, Haoqian Wang, Qingyun Sun, Jun Xu, Qionghai Dai, and Lei Zhang. A pid
283 controller approach for stochastic optimization of deep networks. In *Proceedings of the IEEE*
284 *Conference on Computer Vision and Pattern Recognition*, pages 8522–8531, 2018.
- 285 [2] Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, et al. Distributed opti-
286 mization and statistical learning via the alternating direction method of multipliers. *Foundations*
287 *and Trends® in Machine learning*, 3(1):1–122, 2011.
- 288 [3] G. Chen and M. Teboulle. A proximal-based decomposition method for convex minimization
289 problems. *Mathematical Programming*, 64:81–101, 1994.
- 290 [4] Xiang Cheng, Niladri S Chatterji, Peter L Bartlett, and Michael I Jordan. Underdamped langevin
291 mcmc: A non-asymptotic analysis. *arXiv preprint arXiv:1707.03663*, 2017.
- 292 [5] P. L. Combettes and Jean-Christophe Pesquet. A Douglas-Rachford splitting approach to
293 nonsmooth convex variational signal recovery. *IEEE Journal of Selected Topics in Signal*
294 *Processing*, 1(4):564–574, 2007.
- 295 [6] L. Condat. A primal-dual splitting method for convex optimization involving Lipschitzian,
296 proximable and linear composite terms. *J. Optimization Theory and Applications*, 158(2):
297 460–479, 2013.
- 298 [7] D. Davis and W. Yin. A three-operator splitting scheme and its optimization applications.
299 *Set-Valued and Variational Analysis*, 25(4):829–858, 2017.
- 300 [8] J. Douglas and H. H. Rachford. On the numerical solution of the heat conduction problem in 2
301 and 3 space variables. *Transactions of the American Mathematical Society*, 82:421–439, 1956.
- 302 [9] J. Eckstein. Some saddle-function splitting methods for convex programming. *Optimization*
303 *Methods and Software*, 4(1):75–83, 1994.
- 304 [10] Jonathan Eckstein. Parallel alternating direction multiplier decomposition of convex programs.
305 *Journal of Optimization Theory and Applications*, 80(1):39–62, 1994.
- 306 [11] Jonathan Eckstein and Dimitri P Bertsekas. On the Douglas–Rachford splitting method and the
307 proximal point algorithm for maximal monotone operators. *Mathematical Programming*, 55
308 (1-3):293–318, 1992.
- 309 [12] Jonathan Eckstein and Michael C Ferris. Operator-splitting methods for monotone affine
310 variational inequalities, with a parallel application to optimal control. *INFORMS Journal on*
311 *Computing*, 10(2):218–235, 1998.
- 312 [13] W H Fleming and R W Rishel. *Deterministic and Stochastic Optimal Control*. Stochastic Mod-
313 elling and Applied Probability. Springer New York, 1975. doi: 10.1007/978-1-4612-6380-7.
- 314 [14] Guilherme França, Daniel P Robinson, and René Vidal. ADMM and accelerated ADMM as
315 continuous dynamical systems. In *Proceedings of the 35th International Conference on Machine*
316 *Learning*, pages 1559–1567, 2018.

- 317 [15] Guilherme Franca, Daniel P Robinson, and René Vidal. A dynamical systems perspective on
318 nonsmooth constrained optimization. *arXiv preprint arXiv:1808.04048*, 2018.
- 319 [16] Guilherme França, Daniel P Robinson, and René Vidal. Relax, and accelerate: A continuous
320 perspective on ADMM. *arXiv preprint arXiv:1808.04048*, 2018.
- 321 [17] M. I. Freidlin and A. D. Wentzell. *Random Perturbations of Dynamical Systems*. Grundlehren
322 der mathematischen Wissenschaften. Springer-Verlag, New York, 3 edition, 2012.
- 323 [18] Daniel Gabay and Bertrand Mercier. *A dual algorithm for the solution of non linear vari-*
324 *ational problems via finite element approximation*. Institut de recherche d’informatique et
325 d’automatique, 1975.
- 326 [19] Roland Glowinski and A Marroco. Sur l’approximation, par éléments finis d’ordre un, et la
327 résolution, par pénalisation-dualité d’une classe de problèmes de dirichlet non linéaires. *Revue*
328 *française d’automatique, informatique, recherche opérationnelle. Analyse numérique*, 9(R2):
329 41–76, 1975.
- 330 [20] Donald Goldfarb, Shiqian Ma, and Katya Scheinberg. Fast alternating linearization methods for
331 minimizing the sum of two convex functions. *Mathematical Programming*, 141(1-2):349–382,
332 2013.
- 333 [21] T. Goldstein and S. Osher. The split Bregman method for L1-regularized problems. *SIAM J.*
334 *Imaging Sci.*, 2:323–343, 2009.
- 335 [22] B. S. He, L. Liao, D. Han, and H. Yang. A new inexact alternating direction method for
336 monotone variational inequalities. *Mathematical Programming*, 92:103–118, 2002.
- 337 [23] Feihu Huang, Songcan Chen, and Heng Huang. Faster stochastic alternating direction method
338 of multipliers for nonconvex optimization. In Kamalika Chaudhuri and Ruslan Salakhutdinov,
339 editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of
340 *Proceedings of Machine Learning Research*, pages 2839–2848, Long Beach, California, USA,
341 09–15 Jun 2019. PMLR. URL <http://proceedings.mlr.press/v97/huang19a.html>.
- 342 [24] Michael I Jordan. Dynamical, symplectic and stochastic perspectives on gradient-based opti-
343 mization. *University of California, Berkeley*, 2018.
- 344 [25] P.E. Kloeden and E. Platen. *Numerical Solution of Stochastic Differential Equations*. Stochastic
345 Modelling and Applied Probability. Springer, New York, corrected edition, 2011. ISBN
346 9783662126165. URL <https://books.google.com.hk/books?id=r9r6CAAAQBAJ>.
- 347 [26] Chris Junchi Li, Zhaoran Wang, and Han Liu. Online ICA: Understanding global dynamics of
348 nonconvex optimization via diffusion processes. In *Advances in Neural Information Processing*
349 *System*, pages 4967–4975, 2016.
- 350 [27] Chris Junchi Li, Lei Li, Junyang Qian, and Jian-Guo Liu. Batch size matters: A diffusion approx-
351 imation framework on nonconvex stochastic gradient descent. *arXiv preprint arXiv:1705.07562*,
352 2017.
- 353 [28] Chris Junchi Li, Mengdi Wang, Han Liu, and Tong Zhang. Diffusion approximations for online
354 principal component estimation and global convergence. In *Advances in Neural Information*
355 *Processing System*, 2017.
- 356 [29] Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and adaptive stochastic
357 gradient algorithms. In *34th International Conference on Machine Learning, ICML 2017*, 34th
358 International Conference on Machine Learning, ICML 2017, pages 3306–3340. International
359 Machine Learning Society (IMLS), 1 2017.
- 360 [30] Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and adaptive stochastic
361 gradient algorithms. *arXiv preprint arXiv:1511.06251v3*, 2017.
- 362 [31] Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and dynamics of
363 stochastic gradient algorithms I: Mathematical foundations. *Journal of Machine Learning*
364 *Research*, 20(40):1–47, 2019.

- 365 [32] T. Lin, S. Ma, and S. Zhang. An extragradient-based alternating direction method for convex
366 minimization. *Foundations of Computational Mathematics*, 17(1):35–59, 2017.
- 367 [33] Zhouchen Lin, Risheng Liu, and Zhixun Su. Linearized alternating direction method with
368 adaptive penalty for low-rank representation. In *Advances in neural information processing*
369 *systems*, pages 612–620, 2011.
- 370 [34] S. Ma. Alternating proximal gradient method for convex minimization. *Journal of Scientific*
371 *Computing*, 68(2):546–572, 2016.
- 372 [35] Morteza Mardani, Qingyun Sun, David Donoho, Vardan Papyan, Hafez Monajemi, Shreyas
373 Vasanaawala, and John Pauly. Neural proximal gradient descent for compressive imaging. In
374 *Advances in Neural Information Processing Systems*, pages 9596–9606, 2018.
- 375 [36] Morteza Mardani, Qingyun Sun, Vardan Papyan, Shreyas Vasanaawala, John Pauly, and
376 David Donoho. Degrees of freedom analysis of unrolled neural networks. *arXiv preprint*
377 *arXiv:1906.03742*, 2019.
- 378 [37] G. N. Milstein. Weak approximation of solutions of systems of stochastic differential equations.
379 *Theory of Probability & Its Applications*, 30(4):750–766, 1986. doi: 10.1137/1130095. URL
380 <https://doi.org/10.1137/1130095>.
- 381 [38] G.N. Milstein. *Numerical Integration of Stochastic Differential Equations*, volume 313 of
382 *Mathematics and Its Applications*. Springer, 1995. ISBN 9780792332138. URL [https://](https://books.google.com.hk/books?id=o2y80r_a4W0C)
383 books.google.com.hk/books?id=o2y80r_a4W0C
- 384 [39] Eric Moulines and Francis R. Bach. Non-asymptotic analysis of stochastic approximation
385 algorithms for machine learning. In J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira,
386 and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24*,
387 pages 451–459. Curran Associates, Inc., 2011. URL [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/4316-non-asymptotic-analysis-of-stochastic-approximation-algorithms-for-machine-learning.pdf)
388 [4316-non-asymptotic-analysis-of-stochastic-approximation-algorithms-for-machine-learning.](http://papers.nips.cc/paper/4316-non-asymptotic-analysis-of-stochastic-approximation-algorithms-for-machine-learning.pdf)
389 [pdf](http://papers.nips.cc/paper/4316-non-asymptotic-analysis-of-stochastic-approximation-algorithms-for-machine-learning.pdf)
- 390 [40] Robert Nishihara, Laurent Lessard, Benjamin Recht, Andrew Packard, and Michael I Jordan. A
391 general analysis of the convergence of admm. *arXiv preprint arXiv:1502.02009*, 2015.
- 392 [41] Stanley Osher, Feng Ruan, Jiechao Xiong, Yuan Yao, and Wotao Yin. Sparse recovery via
393 differential inclusions. *Applied and Computational Harmonic Analysis*, 41(2):436–469, 2016.
- 394 [42] Hua Ouyang, Niao He, Long Tran, and Alexander Gray. Stochastic alternating direction
395 method of multipliers. In Sanjoy Dasgupta and David McAllester, editors, *Proceedings of the*
396 *30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine*
397 *Learning Research*, pages 80–88, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. URL
398 <http://proceedings.mlr.press/v28/ouyang13.html>.
- 399 [43] Yuyuan Ouyang, Yunmei Chen, Guanghui Lan, and Eduardo Pasiliao Jr. An accelerated
400 linearized alternating direction method of multipliers. *SIAM Journal on Imaging Sciences*, 8(1):
401 644–681, 2015.
- 402 [44] D. H. Peaceman and H. H. Rachford. The numerical solution of parabolic elliptic differential
403 equations. *SIAM Journal on Applied Mathematics*, 3:28–41, 1955.
- 404 [45] Clarice Poon and Jingwei Liang. Trajectory of alternating direction method of multipliers
405 and adaptive acceleration. In *Advances in Neural Information Processing Systems 32*,
406 pages 7357–7365. Curran Associates, Inc., 2019. URL [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/8955-trajectory-of-alternating-direction-method-of-multipliers-and-adaptive-acceleration.pdf)
407 [8955-trajectory-of-alternating-direction-method-of-multipliers-and-adaptive-acceleration.](http://papers.nips.cc/paper/8955-trajectory-of-alternating-direction-method-of-multipliers-and-adaptive-acceleration.pdf)
408 [pdf](http://papers.nips.cc/paper/8955-trajectory-of-alternating-direction-method-of-multipliers-and-adaptive-acceleration.pdf)
- 409 [46] Ohad Shamir and Tong Zhang. Stochastic gradient descent for non-smooth optimization:
410 Convergence results and optimal averaging schemes. In *Proceedings of the 30th International*
411 *Conference on Machine Learning*, pages 71–79, 2013.
- 412 [47] Bin Shi, Simon S Du, Michael I Jordan, and Weijie J Su. Understanding the acceleration
413 phenomenon via high-resolution differential equations. *arXiv preprint arXiv:1810.08907*, 2018.

- 414 [48] Bin Shi, Simon S Du, Weijie Su, and Michael I Jordan. Acceleration via symplectic discretiza-
415 tion of high-resolution differential equations. In *Advances in Neural Information Processing*
416 *Systems*, pages 5745–5753, 2019.
- 417 [49] Weijie Su, Stephen Boyd, and Emmanuel J Candes. A differential equation for modeling
418 Nesterov’s accelerated gradient method: theory and insights. *Journal of Machine Learning*
419 *Research*, 17(153):1–43, 2016.
- 420 [50] Qingyun Sun and Stephen Boyd. Distributional robust kelly gambling. *arXiv preprint*
421 *arXiv:1812.10371*, 2018.
- 422 [51] Qingyun Sun, Mengyuan Yan, David Donoho, and Stephen Boyd. Convolutional imputation of
423 matrix networks. In *International Conference on Machine Learning*, pages 4818–4827, 2018.
- 424 [52] Taiji Suzuki. Dual averaging and proximal gradient descent for online alternating direction
425 multiplier method. In *International Conference on Machine Learning*, pages 392–400, 2013.
- 426 [53] Apidopoulos Vassilis, Aujol Jean-François, and Dossal Charles. The differential inclusion
427 modeling FISTA algorithm and optimality of convergence rate in the case $b \leq 3$. *SIAM Journal*
428 *on Optimization*, 28(1):551–574, 2018.
- 429 [54] B. C. Vu. A splitting algorithm for dual monotone inclusions involving cocoercive operators.
430 *Advances in Computational Mathematics*, 38(3):667–681, 2013.
- 431 [55] Huahua Wang and Arindam Banerjee. Online alternating direction method. In *Proceedings of*
432 *the 29th International Conference on Machine Learning, ICML 2012*, Proceedings of the 29th
433 International Conference on Machine Learning, ICML 2012, page 1699–1706, 10 2012. ISBN
434 9781450312851.
- 435 [56] Y. Wang, J. Yang, W. Yin, and Y. Zhang. A new alternating minimization algorithm for total
436 variation image reconstruction. *SIAM Journal on Imaging Sciences*, 1(3):248–272, 2008.
- 437 [57] Roscoe B White. *Asymptotic Analysis of Differential Equations*. Imperial College Press, 2010.
- 438 [58] Andre Wibisono, Ashia C Wilson, and Michael I Jordan. A variational perspective on accelerated
439 methods in optimization. *Proceedings of the National Academy of Sciences*, page 201614734,
440 2016.
- 441 [59] Ashia C Wilson, Benjamin Recht, and Michael I Jordan. A lyapunov analysis of momentum
442 methods in optimization. *arXiv preprint arXiv:1611.02635*, 2016.
- 443 [60] Y. Xu. Alternating proximal gradient method for sparse nonnegative Tucker decomposition.
444 *Mathematical Programming Computation*, 7(1):39–70, 2015.
- 445 [61] J. Yang and X. Yuan. Linearized augmented lagrangian and alternating direction methods for
446 nuclear norm minimization. *Mathematics of Computation*, 82(281):301–329, 2013.
- 447 [62] Junfeng Yang and Yin Zhang. Alternating direction algorithms for ℓ_1 -problems in compressive
448 sensing. *SIAM journal on scientific computing*, 33(1):250–278, 2011.
- 449 [63] Huizhuo Yuan, Yuren Zhou, Chris Junchi Li, and Qingyun Sun. Differential inclusions for
450 modeling nonsmooth ADMM variants: A continuous limit theory. In Kamalika Chaudhuri and
451 Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine*
452 *Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 7232–7241, Long
453 Beach, California, USA, 09–15 Jun 2019. PMLR. URL [http://proceedings.mlr.press/
454 v97/yuan19c.html](http://proceedings.mlr.press/v97/yuan19c.html).
- 455 [64] Xiao-Tong Yuan and Tong Zhang. Truncated power method for sparse eigenvalue problems.
456 *Journal of Machine Learning Research*, 14:899–925, 2013.
- 457 [65] X. Zhang, M. Burger, X. Bresson, and S. Osher. Bregmanized nonlocal regularization for
458 deconvolution and sparse reconstruction. *SIAM Journal on Imaging Science*, 3:253–276, 2010.
- 459 [66] Yuchen Zhang, Percy Liang, and Moses Charikar. A hitting time analysis of stochastic gradient
460 langevin dynamics. *arXiv preprint arXiv:1702.05575*, 2017.

- 461 [67] Wenliang Zhong and James Kwok. Fast stochastic alternating direction method of multipliers.
462 In Eric P. Xing and Tony Jebara, editors, *Proceedings of the 31st International Conference*
463 *on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 46–
464 54, Beijing, China, 22–24 Jun 2014. PMLR. URL [http://proceedings.mlr.press/v32/
465 zhong14.html](http://proceedings.mlr.press/v32/zhong14.html).
- 466 [68] Zhengyuan Zhou, Panayotis Mertikopoulos, Nicholas Bambos, Stephen Boyd, and Peter W
467 Glynn. Stochastic mirror descent in variationally coherent optimization problems. In *Advances*
468 *in Neural Information Processing Systems*, pages 7040–7049, 2017.