# Bregman Alternating Direction Method of Multipliers

Huahua Wang
Dept of Computer Science & Engg
University of Minnesota, Twin Cities
huwang@cs.umn.edu

Arindam Banerjee
Dept of Computer Science & Engg
University of Minnesota, Twin Cities
banerjee@cs.umn.edu

**Abstract**

The mirror descent algorithm (MDA) generalizes gradient descent by using a Bregman divergence to replace squared Euclidean distance. In this paper, we similarly generalize the alternating direction method of multipliers (ADMM) to Bregman ADMM (BADMM), which allows the choice of different Bregman divergences to exploit the structure of problems. BADMM provides a unified framework for ADMM and its variants, including generalized ADMM, inexact ADMM and Bethe ADMM. We establish the global convergence and the $O(1/T)$ iteration complexity for BADMM. In some cases, BADMM can be faster than ADMM by a factor of $O(n/\log(n))$ where $n$ is the dimension of the problem. Experimental results are illustrated on the mass transportation problem, which can be solved in parallel by BADMM. BADMM is faster than ADMM and highly optimized commercial software Gurobi, particularly when implemented on GPU.

## 1 Introduction

In recent years, the alternating direction method of multipliers (ADM or ADMM) [4] has been successfully applied in a broad spectrum of applications, ranging from image processing [11, 14] to applied statistics and machine learning [27, 26, 12]. For further understanding of ADMM, we refer the readers to the comprehensive review by [4] and references therein. In particular, ADMM considers the problem of minimizing composite objective functions subject to an equality constraint:

$$\min_{\mathbf{x}\in\mathcal{X},\mathbf{z}\in\mathcal{Z}} f(\mathbf{x}) + g(\mathbf{z}) \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} = \mathbf{c} , \tag{1}$$

where $f$ and $g$ are convex functions, $\mathbf{A} \in \mathbb{R}^{m\times n_1}, \mathbf{B} \in \mathbb{R}^{m\times n_2}, \mathbf{c} \in \mathbb{R}^{m\times 1}, \mathbf{x} \in \mathcal{X} \in \mathbb{R}^{n_1\times 1}, \mathbf{z} \in \mathcal{Z} \in \mathbb{R}^{n_2\times 1}$, and $\mathcal{X}$ and $\mathcal{Z}$ are convex sets. $f$ and $g$ can be non-smooth functions, including indicator functions of convex sets. Many machine learning problems can be cast into the framework of minimizing a composite objective [23, 10], where $f$ is a loss function such as hinge or logistic loss, and $g$ is a regularizer, e.g., $\ell_1$ norm, $\ell_2$ norm, nuclear norm or total variation. The two functions usually have different structures and constraints because they have different tasks in data mining. Therefore, it is useful and sometimes necessary to split and solve them separately, which is exactly the forte of ADMM.

In each iteration, ADMM updates splitting variables separately and alternatively by solving the augmented Lagrangian of (1), which is defined as follows:

$$L_\rho(\mathbf{x}, \mathbf{z}, \mathbf{y}) = f(\mathbf{x}) + g(\mathbf{z}) + \langle \mathbf{y}, \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} - \mathbf{c}\rangle + \frac{\rho}{2}\|\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} - \mathbf{c}\|_2^2, \tag{2}$$

1

where $\mathbf{y} \in \mathbb{R}^m$ is dual variable, $\rho > 0$ is penalty parameter, and the quadratic penalty term is to penalize the violation of the equality constraint. ADMM consists of the following three updates:

$$\mathbf{x}_{t+1} = \operatorname{argmin}_{\mathbf{x}} \ f(\mathbf{x}) + \langle \mathbf{y}_t, \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z}_t - \mathbf{c} \rangle + \frac{\rho}{2}\|\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z}_t - \mathbf{c}\|_2^2 \,, \tag{3}$$

$$\mathbf{z}_{t+1} = \operatorname{argmin}_{\mathbf{z}} \ g(\mathbf{z}) + \langle \mathbf{y}_t, \mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z} - \mathbf{c} \rangle + \frac{\rho}{2}\|\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z} - \mathbf{c}\|_2^2 \,, \tag{4}$$

$$\mathbf{y}_{t+1} = \mathbf{y}_t + \rho(\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c}) \,. \tag{5}$$

Since the computational complexity of $\mathbf{y}$ update (5) is trivial, the computational complexity of ADMM lies in the $\mathbf{x}$ and $\mathbf{z}$ updates (3)-(4) which amount to solving proximal minimization problems using the quadratic penalty term. Inexact ADMM [27, 4] and generalized ADMM [8] have also been proposed to solve the updates inexactly by linearizing the functions and adding additional quadratic terms. Recently, online ADMM [26] and Bethe-ADMM [12] add an additional Bregman divergence on the $\mathbf{x}$ update by keeping or linearizing the quadratic penalty term $\|\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} - \mathbf{c}\|_2^2$. As far as we know, all existing ADMMs use quadratic penalty terms.

A large amount of literature shows that replacing the quadratic term by Bregman divergence in gradient-type methods could greatly boost their performance in solving the constrained optimization problem. First, the use of Bregman divergence could effectively exploit the structure of problems [6, 2, 10] , e.g., in computerized tomography [3], clustering problems and exponential family distributions [1]. Second, in some cases, the gradient descent method with Kullback-Leibler (KL) divergence can outperform the method with the quadratic term by a factor of $O(\sqrt{n \ln n})$ where $n$ is the dimensionality of the problem [2, 3]. Mirror descent algorithm (MDA) and composite objective mirror descent (COMID) [10] use Bregman divergence to replace the quadratic term in gradient descent or proximal gradient [7]. Proximal point method with D-functions (PMD) [6, 5] and Bregman proximal minimization (BPM) [20] generalize proximal point method by using Bregman divegence to replace the quadratic term.

On the side of ADMM, it is still unknown whether the quadratic penalty term in ADMM can be replaced by Bregman divergence, although the convergence of ADMM is well understood. The proof of global convergence of ADMM can be found in [13, 4]. Recently, it has been shown that ADMM converges at a rate of $O(1/T)$ [26, 17], where $T$ is the number of iterations. For strongly convex functions, the dual objective of an accelerated version of ADMM can converge at a rate of $O(1/T^2)$ [15]. Under suitable assumptions like strongly convex functions or a sufficiently small step size for the dual variable update, ADMM can achieve a linear convergence rate [8, 19]. However, as pointed out by [4], "There is currently no proof of convergence known for ADMM with nonquadratic penalty terms."

In this paper, we propose Bregman ADMM (BADMM) which uses Bregman divergences to replace the quadratic penalty term in ADMM, answering the question raised in [4]. More specifically, the quadratic penalty term in the $\mathbf{x}$ and $\mathbf{z}$ updates (3)-(4) will be replaced by a Bregman divergence in BADMM. We also introduce a generalized version of BADMM where two additional Bregman divergences are added to the $\mathbf{x}$ and $\mathbf{z}$ updates. The generalized BADMM (BADMM for short) provides a unified framework for solving (1), which allows one to choose suitable Bregman divergence so that the $\mathbf{x}$ and $\mathbf{z}$ updates can be solved efficiently. BADMM includes ADMM and its variants as special cases. In particular, BADMM replaces all quadratic terms in generalized ADMM [8] with Bregman divergences. By choosing a proper Bregman divergence, we also show that inexact ADMM [27] and Bethe ADMM [12] can be considered as special cases of BADMM. BADMM generalizes ADMM similar to how MDA generalizes gradient descent and how PMD generalizes proximal methods. In BADMM, the $\mathbf{x}$ and $\mathbf{z}$ updates can take the form of MDA or PMD. We establish the global convergence and the $O(1/T)$ iteration complexity for BADMM. In some cases, we show that BADMM can outperform ADMM by a factor $O(n/\ln n)$. We evaluate the performance

of BADMM in solving the linear program problem of mass transportation [18]. By exploiting the structure of the problem, BADMM leads to massive parallelism and can easily run on GPU. BADMM can even be orders of magnitude faster than highly optimized commercial software Gurobi. While Gurobi breaks down in solving a linear program of hundreds of millions of parameters in a server, BADMM takes hundreds of seconds running in a single GPU.

The rest of the paper is organized as follows. In Section 2, we propose Bregman ADMM and discuss several special cases of BADMM. In Section 3, we establish the convergence of BADMM. In Section 4, we consider illustrative applications of BADMM, and conclude in Section 5.

## 2 Bregman Alternating Direction Method of Multipliers

Let $\phi : \Omega \to \mathbb{R}$ be a continuously differentiable and strictly convex function on the relative interior of a convex set $\Omega$. Denote $\nabla \phi(\mathbf{y})$ as the gradient of $\phi$ at $\mathbf{y}$. We define Bregman divergence[1] $B_\phi : \Omega \times \mathrm{ri}(\Omega) \to \mathbb{R}_+$ induced by $\phi$ as

$$B_\phi(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x}) - \phi(\mathbf{y}) - \langle \nabla \phi(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle .$$

Since $\phi$ is convex, $B_\phi(\mathbf{x}, \mathbf{y}) \geq 0$ where the equality holds if and only if $\mathbf{x} = \mathbf{y}$. More details about Bregman divergence can be found in [6, 1]. Two of the most commonly used examples are squared Euclidean distance $B_\phi(\mathbf{x}, \mathbf{y}) = \frac{1}{2}\|\mathbf{x} - \mathbf{y}\|_2^2$ and KL divergence $B_\phi(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n x_i \log \frac{x_i}{y_i}$.

Assuming $B_\phi(\mathbf{c} - \mathbf{A}\mathbf{x}, \mathbf{B}\mathbf{z})$ is well defined, we replace the quadratic penalty term in the augmented Lagrangian (2) by a Bregman divergence as follows:

$$L_\rho^\phi(\mathbf{x}, \mathbf{z}, \mathbf{y}) = f(\mathbf{x}) + g(\mathbf{z}) + \langle \mathbf{y}, \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} - \mathbf{c} \rangle + \rho B_\phi(\mathbf{c} - \mathbf{A}\mathbf{x}, \mathbf{B}\mathbf{z}). \tag{6}$$

Unfortunately, we can not derive Bregman ADMM (BADMM) updates by simply solving $L_\rho^\phi(\mathbf{x}, \mathbf{z}, \mathbf{y})$ alternatinly as ADMM does because Bregman divergences are not necessarily convex in the second argument. More specifically, given $(\mathbf{z}_t, \mathbf{y}_t)$, $\mathbf{x}_{t+1}$ can be obtained by solving $\min_\mathbf{x} L_\rho^\phi(\mathbf{x}, \mathbf{z}_t, \mathbf{y}_t)$, where the quadratic penalty term $\frac{1}{2}\|\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z}_t - \mathbf{c}\|_2^2$ for ADMM in (3) is replaced with $B_\phi(\mathbf{c} - \mathbf{A}\mathbf{x}, \mathbf{B}\mathbf{z}_t)$ in the $\mathbf{x}$ update of BADMM. However, given $(\mathbf{x}_{t+1}, \mathbf{y}_t)$, we cannot obtain $\mathbf{z}_{t+1}$ by solving $\min_\mathbf{z} L_\rho^\phi(\mathbf{x}_{t+1}, \mathbf{z}, \mathbf{y}_t)$, since the term $B_\phi(\mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}, \mathbf{B}\mathbf{z})$ need not be convex in $\mathbf{z}$. The observation motivates a closer look at the role of the quadratic term in ADMM.

In standard ADMM, the quadratic augmentation term added to the Lagrangian is just a penalty term to ensure the new updates do not violate the constraint significantly. Staying with these goals, we propose the $\mathbf{z}$ update augmentation term of BADMM to be: $B_\phi(\mathbf{B}\mathbf{z}, \mathbf{c} - \mathbf{A}\mathbf{x}_{t+1})$, instead of the quadratic penalty term $\frac{1}{2}\|\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z} - \mathbf{c}\|_2^2$ in (3). Then, we get the following updates for BADMM:

$$\mathbf{x}_{t+1} = \underset{\mathbf{x} \in \mathcal{X}}{\mathrm{argmin}} \; f(\mathbf{x}) + \langle \mathbf{y}_t, \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z}_t - \mathbf{c} \rangle + \rho B_\phi(\mathbf{c} - \mathbf{A}\mathbf{x}, \mathbf{B}\mathbf{z}_t) , \tag{7}$$

$$\mathbf{z}_{t+1} = \underset{\mathbf{z} \in \mathcal{Z}}{\mathrm{argmin}} \; g(\mathbf{z}) + \langle \mathbf{y}_t, \mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z} - \mathbf{c} \rangle + \rho B_\phi(\mathbf{B}\mathbf{z}, \mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}) , \tag{8}$$

$$\mathbf{y}_{t+1} = \mathbf{y}_t + \rho(\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c}) . \tag{9}$$

Compared to ADMM (3)-(5), BADMM simply uses a Bregman divergence to replace the quadratic penalty term in the $\mathbf{x}$ and $\mathbf{z}$ updates. It is worth noting that the same Bregman divergence $B_\phi$ is used in the $\mathbf{x}$ and $\mathbf{z}$ updates.

---

[1]The definition of Bregman divergence has been generalized to nondifferentiable functions [20, 24].

We consider a special case when $\mathbf{A} = -\mathbf{I}, \mathbf{B} = \mathbf{I}, \mathbf{c} = \mathbf{0}$. (7) is reduced to

$$\mathbf{x}_{t+1} = \operatorname*{argmin}_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) + \langle \mathbf{y}_t, -\mathbf{x} + \mathbf{z}_t \rangle + \rho B_\phi(\mathbf{x}, \mathbf{z}_t) . \tag{10}$$

If $\phi$ is a quadratic function, the constrained problem (10) requires the projection onto the constraint set $\mathcal{X}$. However, in some cases, if choosing a proper Bregman divergence, (10) can be solved efficiently or has a closed-form solution. For example, if $f$ is a linear function and $\mathcal{X}$ is the unit simplex, $B_\phi$ should be KL divergence, leading to the exponentiated gradient [2, 3, 22]. Interestingly, if the $\mathbf{z}$ update is also the exponentiated gradient, we have alternating exponentiated gradients. In Section 4, we will show the mass transportation problem can be cast into this scenario.

While the updates (7)-(8) use the same Bregman divergences, efficiently solving the $\mathbf{x}$ and $\mathbf{z}$ updates may not be feasible, especially when the structure of the original functions $f, g$, the function $\phi$ used for augmentation, and the constraint sets $\mathcal{X}, \mathcal{Z}$ are rather different. For example, if $f(\mathbf{x})$ is a logistic function in (10), it will not have a closed-form solution even $B_\phi$ is the KL divergence and $\mathcal{X}$ is the unit simplex. To address such concerns, we propose a generalized version of BADMM in Section 2.1.

## 2.1 Generalized BADMM

To allow the use of different Bregman divergences in the $\mathbf{x}$ and $\mathbf{z}$ updates (7)-(9) of BADMM, the generalized BADMM simply introduces an additional Bregman divergence for each update. The generalized BADMM has the following updates:

$$\mathbf{x}_{t+1} = \operatorname*{argmin}_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) + \langle \mathbf{y}_t, \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z}_t - \mathbf{c} \rangle + \rho B_\phi(\mathbf{c} - \mathbf{A}\mathbf{x}, \mathbf{B}\mathbf{z}_t) + \rho_\mathbf{x} B_{\varphi_\mathbf{x}}(\mathbf{x}, \mathbf{x}_t) , \tag{11}$$

$$\mathbf{z}_{t+1} = \operatorname*{argmin}_{\mathbf{z} \in \mathcal{Z}} g(\mathbf{z}) + \langle \mathbf{y}_t, \mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z} - \mathbf{c} \rangle + \rho B_\phi(\mathbf{B}\mathbf{z}, \mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}) + \rho_\mathbf{z} B_{\varphi_\mathbf{z}}(\mathbf{z}, \mathbf{z}_t) , \tag{12}$$

$$\mathbf{y}_{t+1} = \mathbf{y}_t + \tau(\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c}) . \tag{13}$$

where $\rho > 0, \tau > 0, \rho_\mathbf{x} \geq 0, \rho_\mathbf{z} \geq 0$. Note that we allow the use of a different step size $\tau$ in the dual variable update [8, 19]. There are three Bregman divergences in the generalized BADMM. While the Bregman divergence $B_\phi$ is shared by the $\mathbf{x}$ and $\mathbf{z}$ updates, the $\mathbf{x}$ update has its own Bregman divergence $B_{\varphi_\mathbf{x}}$ and the $\mathbf{z}$ update has its own Bregman divergence $B_{\varphi_\mathbf{z}}$. The two additional Bregman divergences in generalized BADMM are variable specific, and can be chosen to make sure that the $\mathbf{x}_{t+1}, \mathbf{z}_{t+1}$ updates are efficient. If all three Bregman divergences are quadratic functions, the generalized BADMM reduces to the generalized ADMM [8]. We prove convergence of generalized BADMM in Section 3, which yields the convergence of BADMM with $\rho_x = \rho_z = 0$.

In the following, we illustrate how to choose a proper Bregman divergence $B_{\varphi_\mathbf{x}}$ so that the $\mathbf{x}$ update can be solved efficiently, e.g., a closed-form solution, noting that the same arguments apply to the $\mathbf{z}$-updates. Consider the first three terms in (11) as $s(\mathbf{x}) + h(\mathbf{x})$, where $s(\mathbf{x})$ denotes an easy term and $h(\mathbf{x})$ is the problematic term which needs to be linearized for an efficient $\mathbf{x}$-update. We illustrate the idea with several examples later in the section. Now, we have

$$\mathbf{x}_{t+1} = \min_{\mathbf{x} \in \mathcal{X}} s(\mathbf{x}) + h(\mathbf{x}) + \rho_\mathbf{x} B_{\varphi_\mathbf{x}}(\mathbf{x}, \mathbf{x}_t) . \tag{14}$$

where efficient updates are difficult due to the mismatch in structure between $h$ and $\mathcal{X}$. The goal is to 'linearize' the function $h$ by using the fact that the Bregman divergence $B_h(\mathbf{x}, \mathbf{x}_t)$ captures all the higher-order (beyond linear) terms in $h(\mathbf{x})$ so that:

$$h(\mathbf{x}) - B_h(\mathbf{x}, \mathbf{x}_t) = h(\mathbf{x}_t) + \langle \mathbf{x} - \mathbf{x}_t, \nabla h(\mathbf{x}_t) \rangle \tag{15}$$

is a linear function of $\mathbf{x}$. Let $\psi$ be another convex function such that one can efficiently solve $\min_{x \in \mathcal{X}} s(\mathbf{x}) + \psi(x) + \langle \mathbf{x}, \mathbf{b} \rangle$ for any constant $\mathbf{b}$. Assuming $\varphi_{\mathbf{x}}(\mathbf{x}) = \psi(\mathbf{x}) - \frac{1}{\eta}h(\mathbf{x})$ is convex, we construct a Bregman divergence based proximal term to the original problem so that:

$$\operatorname*{argmin}_{\mathbf{x} \in \mathcal{X}} s(\mathbf{x}) + h(\mathbf{x}) + \rho_{\mathbf{x}}B_{\varphi_{\mathbf{x}}}(\mathbf{x}, \mathbf{x}_t) = \operatorname*{argmin}_{\mathbf{x} \in \mathcal{X}} s(\mathbf{x}) + \psi(\mathbf{x}) + \langle \mathbf{x}, \frac{1}{\rho_{\mathbf{x}}}\nabla h(\mathbf{x}_t) - \nabla \psi(\mathbf{x}_t) \rangle , \qquad (16)$$

where the latter problem can be solved efficiently, by our assumption. To ensure $\varphi_{\mathbf{x}}$ is convex, we need the following condition:

**Proposition 1** *If $h$ is smooth and has Lipschitz continuous gradients with constant $\nu$ under a p-norm, then $\varphi_{\mathbf{x}}$ is $\nu/\rho_{\mathbf{x}}$-strongly convex w.r.t. the p-norm.*

This condition has been widely used in gradient-type methods, including MDA and COMID. Note that the convergence analysis of generalized ADMM in Section 4 holds for any additional Bregman divergence based proximal terms, and does not rely on such specific choices. Using the above idea, one can 'linearize' different parts of the $\mathbf{x}$ update to yield an efficient update.

We consider three special cases, respectively focusing on linearizing the function $f(\mathbf{x})$, linearizing the Bregman divergence based augmentation term $B_\phi(\mathbf{c} - \mathbf{Ax}, \mathbf{Bz}_t)$, and linearizing both terms, along with examples for each case.

**Case 1: Linearization of smooth function $f$:** Let $h(\mathbf{x}) = f(\mathbf{x})$ in (16), we have

$$\mathbf{x}_{t+1} = \operatorname*{argmin}_{\mathbf{x} \in \mathcal{X}} \langle \nabla f(\mathbf{x}_t), \mathbf{x} - \mathbf{x}_t \rangle + \langle \mathbf{y}_t, \mathbf{Ax} \rangle + \rho B_\phi(\mathbf{c} - \mathbf{Ax}, \mathbf{Bz}_t) + \rho_{\mathbf{x}}B_{\psi_{\mathbf{x}}}(\mathbf{x}, \mathbf{x}_t) . \qquad (17)$$

where $\nabla f(\mathbf{x}_t)$ is the gradient of $f(\mathbf{x})$ at $\mathbf{x}_t$.

**Example 1** Consider the following ADMM form for sparse logistic regression problem [16, 4]:

$$\min_{\mathbf{x}} h(\mathbf{x}) + \lambda \|\mathbf{z}\|_1 , \text{ s.t. } \mathbf{x} = \mathbf{z} , \qquad (18)$$

where $h(\mathbf{x})$ is the logistic function. If we use ADMM to solve (18), the $\mathbf{x}$ update is as follows [4]:

$$\mathbf{x}_{t+1} = \operatorname*{argmin}_{\mathbf{x}} h(\mathbf{x}) + \langle \mathbf{y}_t, \mathbf{x} - \mathbf{z}_t \rangle + \frac{\rho}{2}\|\mathbf{x} - \mathbf{z}_t\|_2^2 , \qquad (19)$$

which is a ridge-regularized logistic regression problem and one needs an iterative algorithm like L-BFGS to solve it. Instead, if we linearize $h(\mathbf{x})$ at $\mathbf{x}_t$ and set $B_\psi$ to be a quadratic function, then

$$\mathbf{x}_{t+1} = \operatorname*{argmin}_{\mathbf{x}} \langle \nabla h(\mathbf{x}_t), \mathbf{x} - \mathbf{x}_t \rangle + \langle \mathbf{y}_t, \mathbf{x} - \mathbf{z}_t \rangle + \frac{\rho}{2}\|\mathbf{x} - \mathbf{z}_t\|_2^2 + \frac{\rho_{\mathbf{x}}}{2}\|\mathbf{x} - \mathbf{x}_t\|_2^2 , \qquad (20)$$

the $\mathbf{x}$ update has a simple closed-form solution.

**Case 2: Linearization of the quadratic penalty term:** In ADMM, $B_\phi(\mathbf{c} - \mathbf{Ax}, \mathbf{Bz}_t) = \frac{1}{2}\|\mathbf{Ax} + \mathbf{Bz}_t - \mathbf{c}\|_2^2$. Let $h(\mathbf{x}) = \frac{1}{2}\|\mathbf{Ax} + \mathbf{Bz}_t - \mathbf{c}\|_2^2$. Then $\nabla h(\mathbf{x}_t) = \mathbf{A}^T(\mathbf{Ax}_t + \mathbf{Bz}_t - \mathbf{c})$, we have

$$\mathbf{x}_{t+1} = \operatorname*{argmin}_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) + \langle \mathbf{y}_t + \rho(\mathbf{Ax}_t + \mathbf{Bz}_t - \mathbf{c}), \mathbf{Ax} \rangle + \rho_{\mathbf{x}}B_\psi(\mathbf{x}, \mathbf{x}_t) . \qquad (21)$$

The case mainly solves the problem due to the $\mathbf{Ax}$ term which makes $\mathbf{x}$ updates nonseparable, whereas the linearized version can be solved with separable (parallel) updates. Several problems have been benefited

from the linearization of quadratic term [8], e.g., when $f$ is $\ell_1$ loss function [16], and projection onto the unit simplex or $\ell_1$ ball [9].

**Case 3: Mirror Descent:** In some settings, we want to linearize both the function $f$ and the quadratic augmentation term $B_\phi(\mathbf{c} - \mathbf{Ax}, \mathbf{Bz}_t) = \frac{1}{2}\|\mathbf{Ax} + \mathbf{Bz}_t - \mathbf{c}\|_2^2$. Let $h(\mathbf{x}) = f(\mathbf{x}) + \langle \mathbf{y}^t, \mathbf{Ax} \rangle + \frac{\rho}{2}\|\mathbf{Ax} + \mathbf{Bz}_t - \mathbf{c}\|_2^2$, we have

$$\mathbf{x}_{t+1} = \underset{\mathbf{x} \in \mathcal{X}}{\operatorname{argmin}} \langle \nabla h(\mathbf{x}_t), \mathbf{x} \rangle + \rho_{\mathbf{x}} B_\psi(\mathbf{x}, \mathbf{x}_t) . \tag{22}$$

Note that (22) is a MDA-type update. Further, one can do a similar exercise with a general Bregman divergence based augmentation term $B_\phi(\mathbf{c} - \mathbf{Ax}, \mathbf{Bz}_t)$, although there has to be a good motivation for going to this route.

**Example 2 [Bethe-ADMM [12]]** Given an undirected graph $G = (V, E)$, where $V$ is the vertex set and $E$ is the edge set. Assume a random discrete variable $X_i$ associated with node $i \in V$ can take $K$ values. In a pairwise MRF, the joint distribution of a set of discrete random variables $X = \{X_1, \cdots, X_n\}$ ($n$ is the number of nodes in the graph) is defined in terms of nodes and cliques [25]. Consider solving the following graph-structured problem :

$$\min \ l(\boldsymbol{\mu}) \ \text{s.t.} \ \boldsymbol{\mu} \in \mathbb{L}(G) , \tag{23}$$

where $l(\boldsymbol{\mu})$ is a decomposable function of $\boldsymbol{\mu}$ and $\mathbb{L}(G)$ is the so-called local polytope [25] determined by the marginalization and normalization (MN) constraints for each node and edge in the graph $G$:

$$\mathbb{L}(G) = \{\boldsymbol{\mu} \geq 0 , \sum_{x_i} \mu_i(x_i) = 1 , \sum_{x_j} \mu_{ij}(x_i, x_j) = \mu_i(x_i)\} , \tag{24}$$

where $\mu_i, \mu_{ij}$ are pseudo-marginal distributions of node $i$ and edge $ij$ respectively. In particular, (23) serves as a LP relaxation of MAP inference probem in a pairwise MRF if $l(\boldsymbol{\mu})$ is defined as follows:

$$l(\boldsymbol{\mu}) = \sum_i \sum_{x_i} \theta_i(x_i)\mu_i(x_i) + \sum_{ij \in E} \sum_{x_{ij}} \theta_{ij}(x_i, x_j)\mu_{ij}(x_i, x_j), \tag{25}$$

where $\theta_i, \theta_{ij}$ are the potential functions of node $i$ and edge $ij$ respectively.

The complexity of polytope $\mathbb{L}(G)$ makes (23) difficult to solve. One possible way is to decompose the graph into trees such that

$$\min \ \sum_\tau c_\tau l_\tau(\boldsymbol{\mu}_\tau) \ \text{s.t.} \ \boldsymbol{\mu}_\tau \in \mathbb{T}_\tau, \boldsymbol{\mu}_\tau = \mathbf{m}_\tau , \tag{26}$$

where $\mathbb{T}_\tau$ denotes the MN constraints (24) in the tree $\tau$. $\boldsymbol{\mu}_\tau$ is a vector of pseudo-marginals of nodes and edges in the tree $\tau$. $\mathbf{m}$ is a global variable which contains all trees and $\mathbf{m}_\tau$ corresponds to the tree $\tau$ in the global variable. $c_\tau$ is the weight for sharing variables. The augmented Lagrangian is

$$L_\rho(\boldsymbol{\mu}_\tau, \mathbf{m}, \boldsymbol{\lambda}_\tau) = \sum_\tau c_\tau l_\tau(\boldsymbol{\mu}_\tau) + \langle \boldsymbol{\lambda}_\tau, \boldsymbol{\mu}_\tau - \mathbf{m}_\tau \rangle + \frac{\rho}{2}\|\boldsymbol{\mu}_\tau - \mathbf{m}_\tau\|_2^2 . \tag{27}$$

which leads to the following update for $\boldsymbol{\mu}_\tau^{t+1}$ in ADMM:

$$\boldsymbol{\mu}_\tau^{t+1} = \underset{\boldsymbol{\mu}_\tau \in \mathbb{T}_\tau}{\operatorname{argmin}} \ c_\tau l_\tau(\boldsymbol{\mu}_\tau) + \langle \boldsymbol{\lambda}_\tau^t, \boldsymbol{\mu}_\tau \rangle + \frac{\rho}{2}\|\boldsymbol{\mu}_\tau - \mathbf{m}_\tau^t\|_2^2 \tag{28}$$

6

(28) is difficult to solve due to the MN constraints in the tree. Let $h(\boldsymbol{\mu}_\tau)$ be the objective of (28). If linearizing $h(\boldsymbol{\mu}_\tau)$ and adding a Bregman divergence in (28), we have:

$$\boldsymbol{\mu}_\tau^{t+1} = \underset{\boldsymbol{\mu}_\tau \in \mathbb{T}_\tau}{\operatorname{argmin}} \ \langle \nabla h(\boldsymbol{\mu}_\tau^t), \boldsymbol{\mu}_\tau \rangle + \rho_{\mathbf{x}} B_\psi(\boldsymbol{\mu}_\tau, \boldsymbol{\mu}_\tau^t)$$

$$= \underset{\boldsymbol{\mu}_\tau \in \mathbb{T}_\tau}{\operatorname{argmin}} \ \langle \nabla h(\boldsymbol{\mu}_\tau^t) - \rho_{\mathbf{x}} \nabla \psi(\boldsymbol{\mu}_\tau^t), \boldsymbol{\mu}_\tau \rangle + \rho_{\mathbf{x}} \psi(\boldsymbol{\mu}_\tau) \,,$$

If $\psi(\boldsymbol{\mu}_\tau)$ is the negative Bethe entropy of $\boldsymbol{\mu}_\tau$, the update of $\boldsymbol{\mu}_\tau^{t+1}$ becomes the Bethe entropy problem [25] and can be solved exactly by the sum-product algorithm in a linear time in the tree.

## 3  Convergence Analysis of BADMM

We need the following assumption in establishing the convergence of BADMM:

**Assumption 1**
*(a) $f : \mathbb{R}^{n_1} \to \mathbb{R} \cup \{+\infty\}$ and $g : \mathbb{R}^{n_2} \to \mathbb{R} \cup \{+\infty\}$ are closed, proper and convex.*
*(b) An optimal solution exists.*
*(c) The Bregman divergence $B_\phi$ is defined on an $\alpha$-strongly convex function $\phi$ with respect to a $p$-norm $\| \cdot \|_p^2$, i.e., $B_\phi(\mathbf{u}, \mathbf{v}) \geq \frac{\alpha}{2} \|\mathbf{u} - \mathbf{v}\|_p^2$, where $\alpha > 0$.*

We start wth the Lagrangian, which is defined as follows:

$$L(\mathbf{x}, \mathbf{y}, \mathbf{z}) = f(\mathbf{x}) + g(\mathbf{z}) + \langle \mathbf{y}, \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} - \mathbf{c} \rangle. \tag{29}$$

Assume that $\{\mathbf{x}^*, \mathbf{z}^*, \mathbf{y}^*\}$ satisfies the KKT conditions of (29), i.e.,

$$-\mathbf{A}^T \mathbf{y}^* \in \partial f(\mathbf{x}^*) \,, \tag{30}$$

$$-\mathbf{B}^T \mathbf{y}^* \in \partial g(\mathbf{z}^*) \,, \tag{31}$$

$$\mathbf{A}\mathbf{x}^* + \mathbf{B}\mathbf{z}^* - \mathbf{c} = 0 \,. \tag{32}$$

$\{\mathbf{x}^*, \mathbf{z}^*, \mathbf{y}^*\}$ is an optimal solution. The optimality conditions of (11) and (12) are

$$-\mathbf{A}^T \{\mathbf{y}_t + \rho(-\nabla\phi(\mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}) + \nabla\phi(\mathbf{B}\mathbf{z}_t))\} - \rho_{\mathbf{x}}(\nabla\varphi_{\mathbf{x}}(\mathbf{x}_{t+1}) - \nabla\varphi_{\mathbf{x}}(\mathbf{x}_t)) \in \partial f(\mathbf{x}_{t+1}) \,, \tag{33}$$

$$-\mathbf{B}^T \{\mathbf{y}_t + \rho(\nabla\phi(\mathbf{B}\mathbf{z}_{t+1}) - \nabla\phi(\mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}))\} - \rho_{\mathbf{z}}(\nabla\varphi_{\mathbf{z}}(\mathbf{z}_{t+1}) - \nabla\varphi_{\mathbf{z}}(\mathbf{z}_t)) \in \partial g(\mathbf{z}_{t+1}) \,. \tag{34}$$

If $\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} = \mathbf{c}$, then $\mathbf{y}_{t+1} = \mathbf{y}_t$. Therefore, (30) is satisfied if $\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_t = \mathbf{c}$ , $\mathbf{x}_{t+1} = \mathbf{x}_t$ in (33). Similarly, (31) is satisfied if $\mathbf{z}_{t+1} = \mathbf{z}_t$ in (34). Overall, the KKT conditions (30)-(32) are satisfied if the following optimality conditions are satisfied:

$$B_{\varphi_{\mathbf{x}}}(\mathbf{x}_{t+1}, \mathbf{x}_t) = 0 \,, \ B_{\varphi_{\mathbf{z}}}(\mathbf{z}_{t+1}, \mathbf{z}_t) = 0 \,, \tag{35a}$$

$$\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_t - \mathbf{c} = 0 \,, \ \mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c} = 0 \,. \tag{35b}$$

For the exact BADMM, $\rho_{\mathbf{x}} = \rho_{\mathbf{z}} = 0$ in (11) and (12), the optimality conditions are (35b), which is equivalent to the optimality conditions used in the proof of ADMM in [4], i.e.,

$$\mathbf{B}\mathbf{z}_{t+1} - \mathbf{B}\mathbf{z}_t = 0 \,, \ \mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c} = 0 \,. \tag{36}$$

Define the residuals of optimality conditions (35) at $(t+1)$ as:

$$R(t+1) = \frac{\rho_{\mathbf{x}}}{\rho} B_{\varphi_{\mathbf{x}}}(\mathbf{x}_{t+1}, \mathbf{x}_t) + \frac{\rho_{\mathbf{z}}}{\rho} B_{\varphi_{\mathbf{z}}}(\mathbf{z}_{t+1}, \mathbf{z}_t) + B_\phi(\mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}, \mathbf{B}\mathbf{z}_t) + \gamma \|\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c}\|_2^2 \,, \tag{37}$$

where $\gamma > 0$. If $R(t+1) = 0$, the optimality conditions (35) and (35b) are satisfied. It is sufficient to show the convergence of BADMM by showing $R(t+1)$ converges to zero. We need the following lemma.

**Lemma 1** *Let the sequence $\{\mathbf{x}_t, \mathbf{z}_t, \mathbf{y}_t\}$ be generated by Bregman ADMM (11)-(13). For any $\mathbf{x}^*, \mathbf{z}^*$ satisfying $\mathbf{A}\mathbf{x}^* + \mathbf{B}\mathbf{z}^* = \mathbf{c}$, we have*

$$\begin{aligned}
&f(\mathbf{x}_{t+1}) + g(\mathbf{z}_{t+1}) - (f(\mathbf{x}^*) + g(\mathbf{z}^*)) \\
&\leq -\langle \mathbf{y}_t, \mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c}\rangle - \rho(B_\phi(\mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}, \mathbf{B}\mathbf{z}_t) + B_\phi(\mathbf{B}\mathbf{z}_{t+1}, \mathbf{c} - \mathbf{A}\mathbf{x}_{t+1})) \\
&\quad + \rho(B_\phi(\mathbf{B}\mathbf{z}^*, \mathbf{B}\mathbf{z}_t) - B_\phi(\mathbf{B}\mathbf{z}^*, \mathbf{B}\mathbf{z}_{t+1})) + \rho_{\mathbf{x}}(B_{\varphi_{\mathbf{x}}}(\mathbf{x}^*, \mathbf{x}_t) - B_{\varphi_{\mathbf{x}}}(\mathbf{x}^*, \mathbf{x}_{t+1}) - B_{\varphi_{\mathbf{x}}}(\mathbf{x}_{t+1}, \mathbf{x}_t)) \\
&\quad + \rho_{\mathbf{z}}(B_{\varphi_{\mathbf{z}}}(\mathbf{z}^*, \mathbf{z}_t) - B_{\varphi_{\mathbf{z}}}(\mathbf{z}^*, \mathbf{z}_{t+1}) - B_{\varphi_{\mathbf{z}}}(\mathbf{z}_{t+1}, \mathbf{z}_t)) \,.
\end{aligned} \tag{38}$$

*Proof:* Using the convexity of $f$ and its subgradient given in (33), we have

$$\begin{aligned}
&f(\mathbf{x}_{t+1}) - f(\mathbf{x}) \\
&\leq \langle -\mathbf{A}^T\{\mathbf{y}_t + \rho(-\nabla\phi(\mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}) + \nabla\phi(\mathbf{B}\mathbf{z}_t)\} - \rho_{\mathbf{x}}(\nabla\varphi_{\mathbf{x}}(\mathbf{x}_{t+1}) - \nabla\varphi_{\mathbf{x}}(\mathbf{x}_t)), \mathbf{x}_{t+1} - \mathbf{x}\rangle \\
&= -\langle \mathbf{y}_t, \mathbf{A}(\mathbf{x}_{t+1} - \mathbf{x})\rangle + \rho\langle \nabla\phi(\mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}) - \nabla\phi(\mathbf{B}\mathbf{z}_t), \mathbf{A}(\mathbf{x}_{t+1} - \mathbf{x})\rangle \\
&\quad - \rho_{\mathbf{x}}\langle \nabla\varphi_{\mathbf{x}}(\mathbf{x}_{t+1}) - \nabla\varphi_{\mathbf{x}}(\mathbf{x}_t), \mathbf{x}_{t+1} - \mathbf{x}\rangle \,.
\end{aligned} \tag{39}$$

Setting $\mathbf{x} = \mathbf{x}^*$ and using $\mathbf{A}\mathbf{x}^* + \mathbf{B}\mathbf{z}^* = \mathbf{c}$, we have

$$\begin{aligned}
&f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*) \\
&\leq -\langle \mathbf{y}_t, \mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}^* - \mathbf{c}\rangle + \rho\langle \nabla\phi(\mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}) - \nabla\phi(\mathbf{B}\mathbf{z}_t), \mathbf{B}\mathbf{z}^* - (\mathbf{c} - \mathbf{A}\mathbf{x}_{t+1})\rangle \\
&\quad - \rho_{\mathbf{x}}\langle \nabla\varphi_{\mathbf{x}}(\mathbf{x}_{t+1}) - \nabla\varphi_{\mathbf{x}}(\mathbf{x}_t), \mathbf{x}_{t+1} - \mathbf{x}\rangle \\
&= -\langle \mathbf{y}_t, \mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}^* - \mathbf{c}\rangle + \rho(B_\phi(\mathbf{B}\mathbf{z}^*, \mathbf{B}\mathbf{z}_t) - B_\phi(\mathbf{B}\mathbf{z}^*, \mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}) - B_\phi(\mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}, \mathbf{B}\mathbf{z}_t)) \\
&\quad + \rho_{\mathbf{x}}(B_{\varphi_{\mathbf{x}}}(\mathbf{x}^*, \mathbf{x}_t) - B_{\varphi_{\mathbf{x}}}(\mathbf{x}^*, \mathbf{x}_{t+1}) - B_{\varphi_{\mathbf{x}}}(\mathbf{x}_{t+1}, \mathbf{x}_t)) \,.
\end{aligned} \tag{40}$$

where the last equality uses the three point property of Bregman divergence, i.e.,

$$\langle \nabla\phi(\mathbf{u}) - \nabla\phi(\mathbf{v}), \mathbf{w} - \mathbf{u}\rangle = B_\phi(\mathbf{w}, \mathbf{v}) - B_\phi(\mathbf{w}, \mathbf{u}) - B_\phi(\mathbf{u}, \mathbf{v}) \,. \tag{41}$$

Similarly, using the convexity of $g$ and its subgradient given in (34), for any $\mathbf{z}$,

$$\begin{aligned}
&g(\mathbf{z}_{t+1}) - g(\mathbf{z}) \\
&\leq \langle -\mathbf{B}^T\{\mathbf{y}_t + \rho(\nabla\phi(\mathbf{B}\mathbf{z}_{t+1}) - \nabla\phi(\mathbf{c} - \mathbf{A}\mathbf{x}_{t+1})\} - \rho_{\mathbf{z}}(\nabla\varphi_{\mathbf{z}}(\mathbf{z}_{t+1}) - \nabla\varphi_{\mathbf{z}}(\mathbf{z}_t)), \mathbf{z}_{t+1} - \mathbf{z}\rangle \\
&= -\langle \mathbf{y}_t, \mathbf{B}(\mathbf{z}_{t+1} - \mathbf{z})\rangle + \rho\langle \nabla\phi(\mathbf{B}\mathbf{z}_{t+1}) - \nabla\phi(\mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}), \mathbf{B}\mathbf{z} - \mathbf{B}\mathbf{z}_{t+1}\rangle \\
&\quad - \rho_{\mathbf{z}}\langle \nabla\varphi_{\mathbf{z}}(\mathbf{z}_{t+1}) - \nabla\varphi_{\mathbf{z}}(\mathbf{z}_t), \mathbf{z}_{t+1} - \mathbf{z}\rangle \\
&= -\langle \mathbf{y}_t, \mathbf{B}(\mathbf{z}_{t+1} - \mathbf{z})\rangle + \rho\{B_\phi(\mathbf{B}\mathbf{z}, \mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}) - B_\phi(\mathbf{B}\mathbf{z}, \mathbf{B}\mathbf{z}_{t+1}) - B_\phi(\mathbf{B}\mathbf{z}_{t+1}, \mathbf{c} - \mathbf{A}\mathbf{x}_{t+1})\} \\
&\quad + \rho_{\mathbf{z}}(B_{\varphi_{\mathbf{z}}}(\mathbf{z}, \mathbf{z}_t) - B_{\varphi_{\mathbf{z}}}(\mathbf{z}, \mathbf{z}_{t+1}) - B_{\varphi_{\mathbf{z}}}(\mathbf{z}_{t+1}, \mathbf{z}_t)) \,.
\end{aligned} \tag{42}$$

where the last equality uses the three point property of Bregman divergence (41). Set $\mathbf{z} = \mathbf{z}^*$ in (42). Adding (40) and (42) completes the proof. ∎

Under Assumption 1(c), the following lemma shows that (37) is bounded by a telescoping series of $D(\mathbf{w}^*, \mathbf{w}_t) - D(\mathbf{w}^*, \mathbf{w}_{t+1})$, where $D(\mathbf{w}^*, \mathbf{w}_t)$ defines the distance from the current iterate $\mathbf{w}_t = (\mathbf{x}_t, \mathbf{z}_t, \mathbf{y}_t)$ to a KKT point $\mathbf{w}^* = (\mathbf{x}^*, \mathbf{z}^*, \mathbf{y}^*)$ as follows:

$$D(\mathbf{w}^*, \mathbf{w}_t) = \frac{1}{2\tau\rho}\|\mathbf{y}^* - \mathbf{y}_t\|_2^2 + B_\phi(\mathbf{Bz}^*, \mathbf{Bz}_t) + \frac{\rho_{\mathbf{x}}}{\rho}B_{\varphi_{\mathbf{x}}}(\mathbf{x}^*, \mathbf{x}_t) + \frac{\rho_{\mathbf{z}}}{\rho}B_{\varphi_{\mathbf{z}}}(\mathbf{z}^*, \mathbf{z}_t). \tag{43}$$

**Lemma 2** *Let the sequence* $\{\mathbf{x}_t, \mathbf{z}_t, \mathbf{y}_t\}$ *be generated by Bregman ADMM (11)-(13) and* $\{\mathbf{x}^*, \mathbf{z}^*, \mathbf{y}^*\}$ *satisfying (30)-(32). Let the Assumption 1 hold.* $R(t+1)$ *and* $D(\mathbf{w}^*, \mathbf{w}_t)$ *are defined in (37) and (43) respectively. Set* $\tau \le (\alpha\sigma - 2\gamma)\rho$, *where* $\sigma = \min\{1, m^{\frac{2}{p}-1}\}$ *and* $0 < \gamma < \frac{\alpha\sigma}{2}$. *Then*

$$R(t+1) \le D(\mathbf{w}^*, \mathbf{w}_t) - D(\mathbf{w}^*, \mathbf{w}_{t+1}). \tag{44}$$

*Proof:* Assume $\{\mathbf{x}^*, \mathbf{y}^*\}$ satisfies (30). Since $f$ is convex, then

$$f(\mathbf{x}^*) - f(\mathbf{x}_{t+1}) \le -\langle \mathbf{A}^T\mathbf{y}^*, \mathbf{x}^* - \mathbf{x}_{t+1}\rangle = -\langle \mathbf{y}^*, \mathbf{Ax}^* - \mathbf{Ax}_{t+1}\rangle. \tag{45}$$

Similarly, for convex function $g$ and $\{\mathbf{z}^*, \mathbf{y}^*\}$ satisfying (31), we have

$$g(\mathbf{z}^*) - g(\mathbf{z}_{t+1}) \le -\langle \mathbf{B}^T\mathbf{y}^*, \mathbf{z}^* - \mathbf{z}_{t+1}\rangle = -\langle \mathbf{y}^*, \mathbf{Bz}^* - \mathbf{Bz}_{t+1}\rangle. \tag{46}$$

Adding them together and using the fact that $\mathbf{Ax}^* + \mathbf{Bz}^* = \mathbf{c}$, we have

$$f(\mathbf{x}^*) + g(\mathbf{z}^*) - (f(\mathbf{x}_{t+1}) + g(\mathbf{z}_{t+1})) \le \langle \mathbf{y}^*, \mathbf{Ax}_{t+1} + \mathbf{Bz}_{t+1} - \mathbf{c}\rangle. \tag{47}$$

Adding (47) and (38) together yields

$$
\begin{aligned}
0 \le {}& \langle \mathbf{y}^* - \mathbf{y}_t, \mathbf{Ax}_{t+1} + \mathbf{Bz}_{t+1} - \mathbf{c}\rangle - \rho(B_\phi(\mathbf{c} - \mathbf{Ax}_{t+1}, \mathbf{Bz}_t) + B_\phi(\mathbf{Bz}_{t+1}, \mathbf{c} - \mathbf{Ax}_{t+1})) \\
& + \rho(B_\phi(\mathbf{Bz}^*, \mathbf{Bz}_t) - B_\phi(\mathbf{Bz}^*, \mathbf{Bz}_{t+1})) + \rho_{\mathbf{x}}(B_{\varphi_{\mathbf{x}}}(\mathbf{x}^*, \mathbf{x}_t) - B_{\varphi_{\mathbf{x}}}(\mathbf{x}^*, \mathbf{x}_{t+1}) - B_{\varphi_{\mathbf{x}}}(\mathbf{x}_{t+1}, \mathbf{x}_t)) \\
& + \rho_{\mathbf{z}}(B_{\varphi_{\mathbf{z}}}(\mathbf{z}^*, \mathbf{z}_t) - B_{\varphi_{\mathbf{z}}}(\mathbf{z}^*, \mathbf{z}_{t+1}) - B_{\varphi_{\mathbf{z}}}(\mathbf{z}_{t+1}, \mathbf{z}_t)).
\end{aligned}
\tag{48}
$$

Using $\mathbf{Ax}_{t+1} + \mathbf{Bz}_{t+1} - \mathbf{c} = \frac{1}{\tau}(\mathbf{y}_{t+1} - \mathbf{y}_t)$, the first term can be rewritten as

$$
\begin{aligned}
\langle \mathbf{y}^* - \mathbf{y}_t, \mathbf{Ax}_{t+1} + \mathbf{Bz}_{t+1} - \mathbf{c}\rangle &= \frac{1}{\tau}\langle \mathbf{y}^* - \mathbf{y}_t, \mathbf{y}_{t+1} - \mathbf{y}_t\rangle \\
&= \frac{1}{2\tau}\left(\|\mathbf{y}^* - \mathbf{y}_t\|_2^2 - \|\mathbf{y}^* - \mathbf{y}_{t+1}\|_2^2 + \|\mathbf{y}_{t+1} - \mathbf{y}_t\|_2^2\right) \\
&= \frac{1}{2\tau}\left(\|\mathbf{y}^* - \mathbf{y}_t\|_2^2 - \|\mathbf{y}^* - \mathbf{y}_{t+1}\|_2^2\right) + \frac{\tau}{2}\|\mathbf{Ax}_{t+1} + \mathbf{Bz}_{t+1} - \mathbf{c}\|_2^2.
\end{aligned}
\tag{49}
$$

Plugging into (48) and rearranging the terms, we have

$$
\begin{aligned}
& \frac{1}{2\tau}\left(\|\mathbf{y}^* - \mathbf{y}_t\|_2^2 - \|\mathbf{y}^* - \mathbf{y}_{t+1}\|_2^2\right) + \rho(B_\phi(\mathbf{Bz}^*, \mathbf{Bz}_t) - B_\phi(\mathbf{Bz}^*, \mathbf{Bz}_{t+1})) \\
& \rho_{\mathbf{x}}(B_{\varphi_{\mathbf{x}}}(\mathbf{x}^*, \mathbf{x}_t) - B_{\varphi_{\mathbf{x}}}(\mathbf{x}^*, \mathbf{x}_{t+1})) + \rho_{\mathbf{z}}(B_{\varphi_{\mathbf{z}}}(\mathbf{z}^*, \mathbf{z}_t) - B_{\varphi_{\mathbf{z}}}(\mathbf{z}^*, \mathbf{z}_{t+1})) \\
& \ge \rho_{\mathbf{x}}B_{\varphi_{\mathbf{x}}}(\mathbf{x}_{t+1}, \mathbf{x}_t) + \rho_{\mathbf{z}}B_{\varphi_{\mathbf{z}}}(\mathbf{z}_{t+1}, \mathbf{z}_t) + \rho B_\phi(\mathbf{c} - \mathbf{Ax}_{t+1}, \mathbf{Bz}_t) \\
& + \rho B_\phi(\mathbf{Bz}_{t+1}, \mathbf{c} - \mathbf{Ax}_{t+1}) - \frac{\tau}{2}\|\mathbf{Ax}_{t+1} + \mathbf{Bz}_{t+1} - \mathbf{c}\|_2^2.
\end{aligned}
\tag{50}
$$

9

Dividing both sides by $\rho$ and letting $R(t+1)$ and $D(\mathbf{w}^*, \mathbf{w}_t)$ be defined in (37) and (43) respectively, we have

$$D(\mathbf{w}^*, \mathbf{w}_t) - D(\mathbf{w}^*, \mathbf{w}_{t+1}) \geq R(t+1) + B_\phi(\mathbf{Bz}_{t+1}, \mathbf{c} - \mathbf{Ax}_{t+1}) - (\frac{\tau}{2\rho} + \gamma)\|\mathbf{Ax}_{t+1} + \mathbf{Bz}_{t+1} - \mathbf{c}\|_2^2$$

$$\geq R(t+1) + \frac{\alpha}{2}\|\mathbf{Ax}_{t+1} + \mathbf{Bz}_{t+1} - \mathbf{c}\|_p^2 - (\frac{\tau}{2\rho} + \gamma)\|\mathbf{Ax}_{t+1} + \mathbf{Bz}_{t+1} - \mathbf{c}\|_2^2 \,, \qquad (51)$$

where the last inequality uses the Assumption 1(c).

If $0 < p \leq 2$, $\|\mathbf{u}\|_p \geq \|\mathbf{u}\|_2$. Set $\frac{\alpha}{2} \geq \frac{\tau}{2\rho} + \gamma$ in (51), i.e., $\tau \leq (\alpha - 2\gamma)\rho$. We can always find a $\gamma < \frac{\alpha}{2}$, thus (44) follows.

If $p > 2$, $\|\mathbf{u}\|_2 \leq m^{\frac{1}{2}-\frac{1}{p}}\|\mathbf{u}\|_p$ for any $\mathbf{u} \in \mathbb{R}^{m \times 1}$, so $\|\mathbf{u}\|_p^2 \geq m^{\frac{2}{p}-1}\|\mathbf{u}\|_2^2$. In (51), set $\frac{\alpha}{2}m^{\frac{2}{p}-1} \geq \frac{\tau}{2\rho} + \gamma$, i.e., $\tau \leq (\alpha m^{\frac{2}{p}-1} - 2\gamma)\rho$. As long as $\gamma < \frac{\alpha}{2}m^{\frac{2}{p}-1}$, we have (44). ∎

**Remark 1** *(a) If $0 < p \leq 2$, then $\sigma = 1$ and $\tau \leq (\alpha - 2\gamma)\rho$. The case that $0 < p \leq 2$ includes two widely used Bregman divergences, i.e., Euclidean distance and KL divergence. For KL divergence in the unit simplex, we have $\alpha = 1, p = 1$ in the Assumption 1 (c), i.e., $KL(\mathbf{u}, \mathbf{v}) \geq \frac{1}{2}\|\mathbf{u} - \mathbf{v}\|_1^2$ [2].*

*(b) Since we often set $B_\phi$ to be a quadratic function ($p = 2$), the three special cases in Section 2.1 could choose step size $\tau = (\alpha - 2\gamma)\rho$.*

*(c) If $p > 2$, the proof requires a sufficiently small step size $\tau$, which may not be needed in practice. It would be interesting to see whether we can use a same $\tau = O(\rho)$ for any $p > 0$ using other proof techniques.*

The following theorem establishes the global convergence for BADMM.

**Theorem 1** *Let the sequence $\{\mathbf{x}_t, \mathbf{z}_t, \mathbf{y}_t\}$ be generated by Bregman ADMM (11)-(13) and $\{\mathbf{x}^*, \mathbf{z}^*, \mathbf{y}^*\}$ satisfying (30)-(32). Let the Assumption 1 hold and $\tau, \gamma$ satisfy the conditions in Lemma 2. Then $R(t+1)$ converges to zero and $\{\mathbf{x}_t, \mathbf{z}_t, \mathbf{y}_t\}$ converges to a KKT point $\{\mathbf{x}^*, \mathbf{z}^*, \mathbf{y}^*\}$ of (29).*

*Proof:* Since $R(t+1) \geq 0$, (44) implies $D(\mathbf{w}^*, \mathbf{w}_{t+1}) \leq D(\mathbf{w}^*, \mathbf{w}_t)$. Therefore, $D(\mathbf{w}^*, \mathbf{w}_t)$ is monotonically nonincreasing and $\mathbf{w}_t$ converges to a KKT point $\mathbf{w}^*$. Summing (44) over $t$ from 0 to $\infty$ yields

$$\sum_{t=0}^{\infty} R(t+1) \leq D(\mathbf{w}^*, \mathbf{w}_0) \,. \qquad (52)$$

Since $R(t+1) \geq 0$, $R(t+1) \to 0$ as $t \to \infty$, which completes the proof. ∎

The following theorem establishs a $O(1/T)$ convergence rate for the objective and residual of constraints in an ergodic sense.

**Theorem 2** *Let the sequences $\{\mathbf{x}_t, \mathbf{z}_t, \mathbf{y}_t\}$ be generated by Bregman ADMM (11),(12),(13) and $\mathbf{y}_0 = \mathbf{0}$. Let $\bar{\mathbf{x}}_T = \frac{1}{T}\sum_{t=1}^T \mathbf{x}_t, \bar{\mathbf{z}}_T = \frac{1}{T}\sum_{t=1}^T \mathbf{z}_t$. Set $\tau \leq (\alpha\sigma - 2\gamma)\rho$, where $\sigma = \min\{1, m^{\frac{2}{p}-1}\}$ and $0 < \gamma < \frac{\alpha\sigma}{2}$. For any $(\mathbf{x}^*, \mathbf{z}^*, \mathbf{y}^*)$ satisfying KKT conditions (30)-(32), we have*

$$f(\bar{\mathbf{x}}_T) + g(\bar{\mathbf{z}}_T) - (f(\mathbf{x}^*) + g(\mathbf{z}^*)) \leq \frac{D_1}{T} \,, \qquad (53)$$

$$\|\mathbf{A}\bar{\mathbf{x}}_T + \mathbf{B}\bar{\mathbf{z}}_T - \mathbf{c}\|_2^2 \leq \frac{D(\mathbf{w}^*, \mathbf{w}_0)}{\gamma T} \,, \qquad (54)$$

*where $D_1 = \rho B_\phi(\mathbf{Bz}^*, \mathbf{Bz}_0) + \rho_{\mathbf{x}} B_{\varphi_{\mathbf{x}}}(\mathbf{x}^*, \mathbf{x}_0) + \rho_{\mathbf{z}} B_{\varphi_{\mathbf{z}}}(\mathbf{z}^*, \mathbf{z}_0)$.*

10

*Proof:* Using (13), we have

$$-\langle \mathbf{y}_t, \mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c} \rangle = -\frac{1}{\tau}\langle \mathbf{y}_t, \mathbf{y}_{t+1} - \mathbf{y}_t \rangle$$

$$= -\frac{1}{2\tau}(\|\mathbf{y}_{t+1}\|_2^2 - \|\mathbf{y}_t\|_2^2 - \|\mathbf{y}_{t+1} - \mathbf{y}_t\|_2^2)$$

$$= \frac{1}{2\tau}(\|\mathbf{y}_t\|_2^2 - \|\mathbf{y}_{t+1}\|_2^2) + \frac{\tau}{2}\|\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c}\|_2^2. \qquad (55)$$

Plugging into (38) and ignoring some negative terms yield

$$f(\mathbf{x}_{t+1}) + g(\mathbf{z}_{t+1}) - (f(\mathbf{x}^*) + g(\mathbf{z}^*))$$

$$\leq \frac{1}{2\tau}(\|\mathbf{y}_t\|_2^2 - \|\mathbf{y}_{t+1}\|_2^2) + \rho(B_\phi(\mathbf{B}\mathbf{z}^*, \mathbf{B}\mathbf{z}_t) - B_\phi(\mathbf{B}\mathbf{z}^*, \mathbf{B}\mathbf{z}_{t+1})) + \rho_\mathbf{x}(B_{\varphi_\mathbf{x}}(\mathbf{x}^*, \mathbf{x}_t) - B_{\varphi_\mathbf{x}}(\mathbf{x}^*, \mathbf{x}_{t+1}))$$

$$+ \rho_\mathbf{z}(B_{\varphi_\mathbf{z}}(\mathbf{z}^*, \mathbf{z}_t) - B_{\varphi_\mathbf{z}}(\mathbf{z}^*, \mathbf{z}_{t+1})) - \rho B_\phi(\mathbf{B}\mathbf{z}_{t+1}, \mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}) + \frac{\tau}{2}\|\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c}\|_2^2. \qquad (56)$$

Assume $B_\phi(\mathbf{B}\mathbf{z}_{t+1}, \mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}) \geq \frac{\alpha}{2}\|\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c}\|_p^2$. If $0 < p \leq 2$, using $\|\mathbf{u}\|_p \leq \|\mathbf{u}\|_2$,

$$-\rho B_\phi(\mathbf{B}\mathbf{z}_{t+1}, \mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}) + \frac{\tau}{2}\|\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c}\|_2^2 \leq -\frac{\alpha\rho - \tau}{2}\|\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c}\|_2^2.$$

Setting $\tau \leq (\alpha - 2\gamma)\rho$, the last two terms on the right hand side of (56) can be removed.

If $p > 2$, $\|\mathbf{u}\|_2 \leq m^{\frac{1}{2}-\frac{1}{p}}\|\mathbf{u}\|_p$ for any $\mathbf{u} \in \mathbb{R}^{m \times 1}$, so $\|\mathbf{u}\|_p^2 \geq m^{\frac{2}{p}-1}\|\mathbf{u}\|_2^2$. Then

$$-\rho B_\phi(\mathbf{B}\mathbf{z}_{t+1}, \mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}) + \frac{\tau}{2}\|\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c}\|_2^2 \leq -\frac{\alpha\rho m^{\frac{2}{p}-1} - \tau}{2}\|\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c}\|_2^2.$$

Setting $\tau \leq (\alpha m^{\frac{2}{p}-1} - 2\gamma)\rho$, the last two terms on the right hand side of (56) can be removed. Summing over $t$ from 0 to $T-1$, we have the following telescoping sum

$$\sum_{t=0}^{T-1}[f(\mathbf{x}_{t+1}) + g(\mathbf{z}_{t+1}) - (f(\mathbf{x}^*) + g(\mathbf{z}^*))]$$

$$\leq \frac{1}{2\tau}\|\mathbf{y}_0\|_2^2 + \rho B_\phi(\mathbf{B}\mathbf{z}^*, \mathbf{B}\mathbf{z}_0) + \rho_\mathbf{x} B_{\varphi_\mathbf{x}}(\mathbf{x}^*, \mathbf{x}_0) + \rho_\mathbf{z}(B_{\varphi_\mathbf{z}}(\mathbf{z}^*, \mathbf{z}_0)$$

$$= \rho B_\phi(\mathbf{B}\mathbf{z}^*, \mathbf{B}\mathbf{z}_0) + \rho_\mathbf{x} B_{\varphi_\mathbf{x}}(\mathbf{x}^*, \mathbf{x}_0) + \rho_\mathbf{z}(B_{\varphi_\mathbf{z}}(\mathbf{z}^*, \mathbf{z}_0). \qquad (57)$$

Dividing both sides by $T$ and applying the Jensen's inequality gives (59).

Dividing both sides of (52) by $T$ and applying the Jensen's inequality yield (60). ∎

We consider one special case of BADMM which could outperform ADMM. Assume $\mathbf{B} = \mathbf{I}$ and $\mathcal{X}, \mathcal{Z}$ are the unit simplex. Let $B_\phi$ be the KL divergence. For $\mathbf{z} \in \mathbb{R}^{n_2 \times 1}$, we have

$$B_\phi(\mathbf{z}^*, \mathbf{z}_0) = \sum_{i=1}^{n_2} z_i^* \ln\frac{z_i^*}{z_{i,0}} = \sum_{i=1}^{n_2} z_i^* \ln z_i^* + \ln n_2 \leq \ln n_2. \qquad (58)$$

Similarly, if $\rho_\mathbf{x} > 0$, by choosing $\mathbf{x}_0 = \mathbf{e}/n_2$, $B_{\varphi_\mathbf{x}}(\mathbf{x}^*, \mathbf{x}_0) \leq \ln n_1$. Setting $\alpha = 1, \sigma = 1$ and $\gamma = \frac{1}{4}$ in Theorem 2 yields the following result:

**Corollary 1** *Let the sequences $\{\mathbf{x}_t, \mathbf{z}_t, \mathbf{y}_t\}$ be generated by Bregman ADMM (11),(12),(13) and $\mathbf{y}_0 = \mathbf{0}$. Assume $\mathbf{B} = \mathbf{I}$, and $\mathcal{X}$ and $\mathcal{Z}$ is the unit simplex. Let $B_\phi, B_{\varphi_\mathbf{x}}, B_{\varphi_\mathbf{z}}$ be KL divergence. Let $\bar{\mathbf{x}}_T = \frac{1}{T}\sum_{t=1}^T \mathbf{x}_t$, $\bar{\mathbf{z}}_T = \frac{1}{T}\sum_{t=1}^T \mathbf{z}_t$. Set $\tau = \frac{3\rho}{4}$. For any $(\mathbf{x}^*, \mathbf{z}^*, \mathbf{y}^*)$ satisfying KKT conditions (??), we have*

$$f(\bar{\mathbf{x}}_T) + g(\bar{\mathbf{z}}_T) - (f(\mathbf{x}^*) + g(\mathbf{z}^*)) \leq \frac{\rho \ln n_2 + \rho_\mathbf{x} \ln n_1 + \rho_\mathbf{z} \ln n_2}{T} , \tag{59}$$

$$\|\mathbf{A}\bar{\mathbf{x}}_T + \mathbf{B}\bar{\mathbf{z}}_T - \mathbf{c}\|_2^2 \leq \frac{\frac{2}{\tau\rho}\|\mathbf{y}^* - \mathbf{y}_0\|_2^2 + 4\ln n_2 + \frac{4\rho_\mathbf{x}}{\rho}\ln n_1 + \frac{4\rho_\mathbf{z}}{\rho}\ln n_2}{T} , \tag{60}$$

**Remark 2** *(a) In [2], it shows that MDA yields a smilar $O(\ln n)$ bound where $n$ is dimensionality of the problem. If the diminishing step size of MDA is propotional to $\sqrt{\ln n}$, the bound is $O(\sqrt{\ln n})$. Therefore, MDA can outperform the gradient method by a factor $O((n/\ln n)^{1/2})$.*

*(b) With constant step size, BADMM outperforms ADMM by a factor $O(n/\ln n)$ in an ergodic sense.*

## 4 Experimental Results

In this section, we use BADMM to solve the mass transportation problem [18]:

$$\min \quad \langle \mathbf{C}, \mathbf{X} \rangle \quad \text{s.t.} \quad \mathbf{X}\mathbf{e} = \mathbf{a}, \mathbf{X}^T\mathbf{e} = \mathbf{b}, \mathbf{X} \geq 0 . \tag{61}$$

where $\langle \mathbf{C}, \mathbf{X} \rangle$ denotes $\mathrm{Tr}(\mathbf{C}^T\mathbf{X})$, $\mathbf{C} \in \mathbb{R}^{m \times n}$ is a cost matrix, $\mathbf{e}$ is a column vector of ones. (61) is called the assignment problem and can be solved exactly by the Hungarian method [21]. The mass transportation problem (61) is a linear program and thus can be solved by the simplex method.

We now show that (61) can be solved by ADMM and BADMM. We first introduce a variable $\mathbf{Z}$ to split the constraints into two simplex such that $\Delta_\mathbf{x} = \{\mathbf{X}|\mathbf{X} \geq 0, \mathbf{X}\mathbf{e} = \mathbf{a}\}$ and $\Delta_\mathbf{z} = \{\mathbf{Z}|\mathbf{Z} \geq 0, \mathbf{Z}^T\mathbf{e} = \mathbf{b}\}$. (61) can be rewritten in the following ADMM form:

$$\min \quad \langle \mathbf{C}, \mathbf{X} \rangle \quad \text{s.t.} \quad \mathbf{X} \in \Delta_\mathbf{x}, \mathbf{Z} \in \Delta_\mathbf{z}, \mathbf{X} = \mathbf{Z} . \tag{62}$$

(62) can be solved by ADMM which requires the Euclidean projection onto the simplex $\Delta_\mathbf{x}$ and $\Delta_\mathbf{z}$, although the projection can be done efficiently [9]. We use BADMM to solve (62):

$$\mathbf{X}^{t+1} = \operatorname*{argmin}_{\mathbf{X} \in \Delta_\mathbf{x}} \langle \mathbf{C}, \mathbf{X} \rangle + \langle \mathbf{Y}^t, \mathbf{X} \rangle + \rho \mathrm{KL}(\mathbf{X}, \mathbf{Z}^t) , \tag{63}$$

$$\mathbf{Z}^{t+1} = \operatorname*{argmin}_{\mathbf{Z} \in \Delta_\mathbf{z}} \langle \mathbf{Y}^t, -\mathbf{Z} \rangle + \rho \mathrm{KL}(\mathbf{Z}, \mathbf{X}^{t+1}) , \tag{64}$$

$$\mathbf{Y}^{t+1} = \mathbf{Y}^t + \rho(\mathbf{X}^{t+1} - \mathbf{Z}^{t+1}) . \tag{65}$$

Both (63) and (64) have closed-form solutions, i.e.,

$$X_{ij}^{t+1} = \frac{Z_{ij}^t \exp(-\frac{C_{ij}+Y_{ij}^t}{\rho})}{\sum_{j=1}^n Z_{ij}^t \exp(-\frac{C_{ij}+Y_{ij}^t}{\rho})}a_i , \quad Z_{ij}^{t+1} = \frac{X_{ij}^{t+1} \exp(\frac{Y_{ij}^t}{\rho})}{\sum_{i=1}^m X_{ij}^{t+1} \exp(\frac{Y_{ij}^t}{\rho})}b_j \tag{66}$$

which are exponentiated graident updates and can be done in $O(mn)$. Besides the sum operation which can be done in $O(\log(n))$, (66) amounts to elementwise operation and thus can be done in parallel. According to Corollary 1, BADMM can be faster than ADMM by a factor of $O(n/\log(n))$.

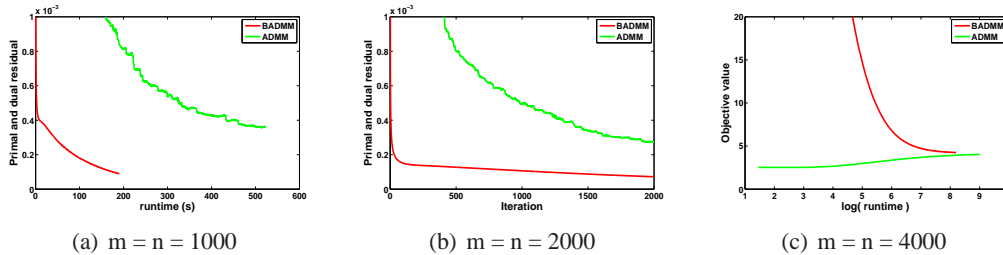(a) m = n = 1000          (b) m = n = 2000          (c) m = n = 4000

Figure 1: Comparison BADMM and ADMM. BADMM converges faster than ADMM.

We compare BADMM with ADMM and a highly optimized commercial linear programming solvers on the mass transportation problem (61) when $m = n$ and $\mathbf{a} = \mathbf{b} = \mathbf{e}$. $\mathbf{C}$ is randomly generated from the uniform distribution. They run 5 times and the average is reported. We choose the 'best' parameter for BADMM ($\rho = 0.001$) and ADMM ($\rho = 0.001$). The stopping condition is either when the number of iterations exceeds 2000 or when the primal-dual residual is less than $10^{-4}$.

**BADMM vs ADMM:** Figure 1 compares BADMM and ADMM with different dimensions $n = \{1000, 2000, 4000\}$ running on a single CPU. Figure 1(a) plots the primal and dual residual against the runtime when the dimension is 1000, and Figure 1(b) plots the convergence of primal and dual residual over iteration when the dimension is 2000. BADMM converges faster than ADMM. Figure 1(c) plots the convergence of objective value against the log of runtime. BADMM converges faster than ADMM even when the initial point is further from the optimum.

**BADMM vs Gurobi:** Gurobi[2] is a highly optimized commercial software where linear programming solvers have been efficiently implemented. We run Gurobi on two settings: a Mac laptop with 6G memory and a server with 86G memory, respectively. For comparison, BADMM is run in parallel on a Tesla M2070 GPU with 5G memory and 448 cores[3]. We experiment with large scale problems and use $m = n = \{1, 5, 10, 15\} \times 2^{10}$. Table 1 shows the runtime and the objective values of BADMM and Gurobi, where a '-' indicates the algorithm did not terminate. In spite of Gurobi being one of the most optimized LP solvers, BADMM running in parallel is several times faster than Gurobi. In fact, for larger values of $n$, Gurobi did not terminate even on the 86G server, whereas BADMM was efficient even with just 5G memory! The complexity of most LP solvers in Gurobi is $O(n^3)$ and can become slow as $n$ increases, especially at the scales we consider. Moreover, the memory consumption of Gurobi increases rapidly with the increase of $n$. When $n = 5 \times 2^{10}$, the memory required by Gurobi surpassed the memory in the laptop, leading to the rapid increase of time. A similar situation was also observed in the server with 86G when $n = 10 \times 2^{10}$. In contrast, the memory required by BADMM is $O(n^2)$—even when $n = 15 \times 2^{10}$ (more than 0.2 billion parameters), BADMM can still run on a single GPU with only 5G memory.

The results clearly illustrate the promise of BADMM. With more careful implementation and code optimization, BADMM has the potential to solve large scale problems efficiently in parallel with small memory foot-print.

## 5 Conclusions

In this paper, we generalized the alternating direction method of multipliers (ADMM) to Bregman ADMM, similar to how mirror descent generalizes gradient descent. BADMM defines a unified framework for ADMM, generalized ADMM, inexact ADMM and Bethe ADMM. The global convergence and the $O(1/T)$

---

[2]http://www.gurobi.com/

[3]GPU code is available on http://www-users.cs.umn.edu/~huwang/badmm_mt.zip

Table 1: Comparison of BADMM (GPU) with Gurobi

| m=n | Gurobi (Laptop) | | Gurobi (Server) | | BADMM (GPU) | |
|---|---|---|---|---|---|---|
| | time | objective | time | objective | time | objective |
| $2^{10}$ | 4.22 | 1.69 | 2.66 | 1.69 | 0.54 | 1.69 |
| $5 \times 2^{10}$ | 377.14 | 1.61 | 92.89 | 1.61 | 22.15 | 1.61 |
| $10 \times 2^{10}$ | - | - | 1235.34 | 1.65 | 117.75 | 1.65 |
| $15 \times 2^{10}$ | - | - | - | - | 303.54 | 1.63 |

iteration complexity of BADMM are also established. In some cases, BADMM is faster than ADMM by a factor of $O(n/\log(n))$. BADMM can also be faster than highly optimized commercial software in solving linear program of mass transportation problem.

## Acknowledgment

## References

[1] A. Banerjee, S. Merugu, I. Dhillon, and J. Ghosh. Clustering with Bregman divergences. *Journal of Machine Learning Research (JMLR)*, 6:1705–1749, 2005.

[2] A. Beck and M. Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31:167–175, 2003.

[3] A. Ben-Tal, T. Margalit, and A. Nemirovski. The ordered subsets mirror descent optimization method with applications to tomography. *SIAM Journal on Optimization*, 12:79–108, 2001.

[4] S. Boyd, E. Chu N. Parikh, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundation and Trends Machine Learning*, 3(1):1–122, 2011.

[5] Y. Censor and S. Zenios. *Parallel Optimization: Theory, Algorithms, and Applications*. Oxford University Press, 1998.

[6] G. Chen and M. Teboulle. Convergence analysis of a proximal-like minimization algorithm using bremgan functions. *SIAM Journal on Optimization*, 3:538–543, 1993.

[7] P. Combettes and J. Pesquet. Proximal splitting methods in signal processsing. *Fixed-Point Algorithms for Inverse Problems in Science and Engineering Springer (Ed.)*, pages 185–212, 2011.

[8] W. Deng and W. Yin. On the global and linear convergence of the generalized alternating direction method of multipliers. *ArXiv*, 2012.

[9] J. Duchi, S. Shalev-Shwartz, Y. Singer, and T. Chandra. Efficient projections onto the $l_1$-ball for learning in high dimensions. In *International Conference on Machine Learning (ICML)*, pages 272–279, 2008.

[10] J. Duchi, S. Shalev-Shwartz, Y. Singer, and A. Tewari. Composite objective mirror descent. In *Conference on Learning Theory (COLT)*, 2010.

[11] M. A. T. Figueiredo and J. M. Bioucas-Dias. Restoration of poissonian images using alternating direction optimization. *IEEE Transactions on Image Processing*, 19:3133–3145, 2010.

[12] Q. Fu, H. Wang, and A. Banerjee. Bethe-ADMM for tree decomposition based parallel MAP inference. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2013.

[13] D. Gabay. Applications of the method of multipliers to variational inequalities. In *Augmented Lagrangian Methods: Applications to the Solution of Boundary-Value Problems*. M. Fortin and R. Glowinski, eds., North-Holland: Amsterdam, 1983.

[14] T. Goldstein, X. Bresson, and S. Osher. Geometric applications of the split Bregman method: segmentation and surface reconstruction. *Journal of Scientific Computing*, 45(1):272–293, 2010.

[15] T. Goldstein, B. Donoghue, and S. Setzer. Fast alternating direction optimization methods. *CAM report 12-35, UCLA*, 2012.

[16] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2009.

[17] B. He and X. Yuan. On the $O(1/n)$ convergence rate of the Douglas-Rachford alternating direction method. *SIAM Journal on Numerical Analysis*, 50:700–709, 2012.

[18] F. L. Hitchcock. The distribution of a product from several sources to numerous localities. *Journal of Mathematical Physics*, 20:224–230, 1941.

[19] M. Hong and Z. Luo. On the linear convergence of the alternating direction method of multipliers. *ArXiv*, 2012.

[20] K. C. Kiwiel. Proximal minimization methods with generalized Bregman functions. *SIAM Journal on Control and Optimization*, 35:1142–1168, 1995.

[21] H. W. Kuhn. The hungarian method for the assignment problem. *Naval Research Logistics Quarterly*, 2:83–97, 1955.

[22] A. Nemirovski and D. Yudin. *Problem Complexity and Method Efficiency in Optimization*. Wiley, 1983.

[23] Y. Nesterov. Gradient methods for minimizing composite objective function. *Technical Report 76, Center for Operation Research and Economics (CORE), Catholic University of Louvain (UCL)*, 2007.

[24] M. Telgarsky and S. Dasgupta. Agglomerative Bregman clustering. In *International Conference on Machine Learning (ICML)*, 2012.

[25] M. J. Wainwright and M. I. Jordan. Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 1:1–305, 2008.

[26] H. Wang and A. Banerjee. Online alternating direction method. In *International Conference on Machine Learning (ICML)*, 2012.

[27] J. Yang and Y. Zhang. Alternating direction algorithms for L1-problems in compressive sensing. *ArXiv*, 2009.

# A  Convergence of BADMM with Time Varying Step Size

Under the assumption that $\mathbf{y}_t$ is bounded, the following theorem requires a large step size to establish the convergence of BADMM.

**Theorem 3** *Let the sequences* $\{\mathbf{x}_t, \mathbf{z}_t, \mathbf{y}_t\}$ *be generated by Bregman ADMM (11)-(13) and* $\{\mathbf{x}^*, \mathbf{z}^*, \mathbf{y}^*\}$ *satisfying (30)-(32). Let the Assumption 1 hold and* $\|\mathbf{y}_t\|_2 \le D_\mathbf{y}$. *Setting* $\rho_\mathbf{x} = \rho_\mathbf{z} = c_1\sqrt{T}, \tau = c_2\sqrt{T}$ *and* $\rho = \sqrt{T}$ *for some positive constant* $c_1, c_2$, *then* $R(t+1)$ *converges to zero.*

*Proof:*  Assuming $\|\mathbf{y}_t\|_2 \le D_\mathbf{y}$ and using (13), we have

$$\|\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c}\|_2^2 = \frac{1}{\tau^2}\|\mathbf{y}_{t+1} - \mathbf{y}_t\|_2^2 \le \frac{2}{\tau^2}(\|\mathbf{y}_{t+1}\|_2^2 + \|\mathbf{y}_t\|_2^2) \le \frac{4D_\mathbf{y}^2}{\tau^2} . \tag{67}$$

Plugging into (51) and rearranging the terms yields

$$R(t+1) \le D(\mathbf{w}^*, \mathbf{w}_t) - D(\mathbf{w}^*, \mathbf{w}_{t+1}) + (\frac{\tau}{2\rho} + \gamma)\frac{4D_\mathbf{y}^2}{\tau^2} . \tag{68}$$

Setting $\rho_\mathbf{x} = \rho_\mathbf{z} = c_1\sqrt{T}, \tau = c_2\sqrt{T}$ and $\rho = \sqrt{T}$ for some positive constant $c_1, c_2$, we have

$$R(t+1) = c_1 B_{\varphi_\mathbf{x}}(\mathbf{x}_{t+1}, \mathbf{x}_t) + c_1 B_{\varphi_\mathbf{z}}(\mathbf{z}_{t+1}, \mathbf{z}_t) + B_\phi(\mathbf{c} - \mathbf{A}\mathbf{x}_{t+1}, \mathbf{B}\mathbf{z}_t) + \gamma\|\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c}\|_2^2 , \tag{69}$$

Summing (68) over $t$ from 0 to $T-1$, we have the following telescoping sum

$$\sum_{t=0}^{T-1} R(t+1) \le D(\mathbf{w}^*, \mathbf{w}_0) + \sum_{t=0}^{T-1}(\frac{\tau}{2\rho} + \gamma)\frac{4D_\mathbf{y}^2}{\tau^2} = D(\mathbf{w}^*, \mathbf{w}_0) + \frac{4(c_2/2 + \gamma)D_\mathbf{y}^2}{c_2^2} . \tag{70}$$

Therefore, $R(t+1) \to 0$ as $t \to \infty$.  ∎

The following theorem establishs the convergence rate for the objective and residual of constraints in an ergodic sense.

**Theorem 4** *Let the sequences* $\{\mathbf{x}_t, \mathbf{z}_t, \mathbf{y}_t\}$ *be generated by Bregman ADMM (11)-(13). Let* $\bar{\mathbf{x}}_T = \frac{1}{T}\sum_{t=1}^T \mathbf{x}_t, \bar{\mathbf{z}}_T = \frac{1}{T}\sum_{t=1}^T \mathbf{z}_t$. *Let the Assumption 1 hold and* $\|\mathbf{y}_t\|_2 \le D_\mathbf{y}$. *Set* $\rho_\mathbf{x} = \rho_\mathbf{z} = c_1\sqrt{T}, \tau = c_2\sqrt{T}, \rho = \sqrt{T}$ *for some positive constants* $c_1, c_2$. *For any* $(\mathbf{x}^*, \mathbf{z}^*, \mathbf{y}^*)$ *satisfying KKT conditions (30)-(32), we have*

$$f(\bar{\mathbf{x}}_T) + g(\bar{\mathbf{z}}_T) - (f(\mathbf{x}^*) + g(\mathbf{z}^*)) \le \frac{2D_\mathbf{y}^2}{c_2\sqrt{T}} + \frac{\|\mathbf{y}_0\|_2^2}{2c_2T\sqrt{T}} + \frac{D_2}{\sqrt{T}} , \tag{71}$$

$$\|\mathbf{A}\bar{\mathbf{x}}_T + \mathbf{B}\bar{\mathbf{z}}_T - \mathbf{c}\|_2^2 \le \frac{D(\mathbf{w}^*, \mathbf{w}_0)}{\gamma T} + \frac{4(c_2/2 + \gamma)D_\mathbf{y}^2}{\gamma c_2^2 T} , \tag{72}$$

*where* $D_2 = B_\phi(\mathbf{B}\mathbf{z}^*, \mathbf{B}\mathbf{z}_0) + c_1(B_{\varphi_\mathbf{x}}(\mathbf{x}^*, \mathbf{x}_0) + B_{\varphi_\mathbf{z}}(\mathbf{z}^*, \mathbf{z}_0))$.

*Proof:* Assuming $\|\mathbf{y}_t\|_2 \leq D_{\mathbf{y}}^2$ and using (13), we have

$$-\langle \mathbf{y}_t, \mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c} \rangle = -\frac{1}{\tau}\langle \mathbf{y}_t, \mathbf{y}_{t+1} - \mathbf{y}_t \rangle \leq \frac{1}{\tau}(\|\mathbf{y}_t\|_2^2 + \|\mathbf{y}_t\|_2 * \|\mathbf{y}_{t+1}\|_2) \leq \frac{2D_{\mathbf{y}}^2}{\tau}. \qquad (73)$$

Plugging into (38) and ignoring some negative terms yield

$$f(\mathbf{x}_{t+1}) + g(\mathbf{z}_{t+1}) - (f(\mathbf{x}^*) + g(\mathbf{z}^*))$$
$$\leq \frac{2D_{\mathbf{y}}^2}{\tau} + \rho(B_\phi(\mathbf{B}\mathbf{z}^*, \mathbf{B}\mathbf{z}_t) - B_\phi(\mathbf{B}\mathbf{z}^*, \mathbf{B}\mathbf{z}_{t+1})) + \rho_{\mathbf{x}}(B_{\varphi_{\mathbf{x}}}(\mathbf{x}^*, \mathbf{x}_t) - B_{\varphi_{\mathbf{x}}}(\mathbf{x}^*, \mathbf{x}_{t+1}))$$
$$+ \rho_{\mathbf{z}}(B_{\varphi_{\mathbf{z}}}(\mathbf{z}^*, \mathbf{z}_t) - B_{\varphi_{\mathbf{z}}}(\mathbf{z}^*, \mathbf{z}_{t+1})). \qquad (74)$$

Summing over $t$ from 0 to $T-1$, we have the following telescoping sum

$$\sum_{t=0}^{T-1} [f(\mathbf{x}_{t+1}) + g(\mathbf{z}_{t+1}) - (f(\mathbf{x}^*) + g(\mathbf{z}^*))]$$
$$\leq \sum_{t=0}^{T-1} \frac{2D_{\mathbf{y}}^2}{\tau} + \frac{1}{2\tau}\|\mathbf{y}_0\|_2^2 + \rho B_\phi(\mathbf{B}\mathbf{z}^*, \mathbf{B}\mathbf{z}_0) + \rho_{\mathbf{x}} B_{\varphi_{\mathbf{x}}}(\mathbf{x}^*, \mathbf{x}_0) + \rho_{\mathbf{z}} B_{\varphi_{\mathbf{z}}}(\mathbf{z}^*, \mathbf{z}_0).$$

Setting $\rho_{\mathbf{x}} = \rho_{\mathbf{z}} = c_1\sqrt{T}, \tau = c_2\sqrt{T}, \rho = \sqrt{T}$, dividing both sides by $T$ and applying the Jensen's inequality yield (71).

Dividing both sides of (70) by $T$ and applying the Jesen's inequality yield (72). ∎