

Asynchronous Distributed Optimization using a Randomized Alternating Direction Method of Multipliers

Franck Iutzeler, Pascal Bianchi, Philippe Ciblat and Walid Hachem

Abstract— Consider a set of networked agents endowed with private cost functions and seeking to find a consensus on the minimizer of the aggregate cost. A new class of random asynchronous distributed optimization methods is introduced. The methods generalize the standard Alternating Direction Method of Multipliers (ADMM) to an asynchronous setting where isolated components of the network are activated in an uncoordinated fashion. The algorithms rely on the introduction of *randomized Gauss-Seidel* iterations of Douglas-Rachford splitting leading to an asynchronous algorithm based on the ADMM. Convergence to the sought minimizers is provided under mild connectivity conditions.

I. INTRODUCTION

Consider a network represented by a set V of agents seeking to solve the following optimization problem on a Euclidean space X :

$$\inf_{x \in X} \sum_{v \in V} f_v(x), \quad (1)$$

where f_v is a convex real function known by agent v only. Function f_v can be interpreted as the price paid by an agent v when the global network state is equal to x .

This problem arises for instance in *cloud learning* applications where massive data sets are distributed in a network and processed by distinct virtual machines [1]. We investigate distributed optimization algorithms: agents iteratively update a local estimate using their private objective f_v and, simultaneously, exchange information with their neighbors in order to eventually reach a consensus on the global solution. Standard algorithms are generally *synchronous*: all agents are supposed to complete their local computations synchronously at each tick of an external clock, and then synchronously merge their local results. However, in many situations, one faces variable sizes of the local data sets along with heterogeneous computational abilities of the virtual machines. Synchronism then becomes a burden, as the global convergence rate is expected to depend on the local computation times of the slowest agents. It is crucial to introduce asynchronous methods which allow the estimates to be updated in a non-coordinated fashion, rather than all together or in some frozen order.

The literature contains at least three classes of distributed optimization methods for solving (1). The first one is based on the simultaneous use of a local *first-order* optimization algorithm (subgradient algorithm [2], [3], [4], Nesterov-like method [5], [6]) and a gossip process which drives the

network to a consensus. A second class of methods is formed by distributed Newton-Raphson methods [7]. This paper focuses on a third class of methods derived from proximal splitting methods [8], [9], [10]. Perhaps the most emblematic proximal splitting method is the so-called *Alternating Direction Method of Multipliers* (ADMM) recently popularized to multiagent systems by the monograph [11]. Schizas *et al.* demonstrated the remarkable potential of ADMM to handle distributed optimization problems and introduce a useful framework to encompass graph-constrained communications [13]. We also refer to [14], [15] for recent contributions. However, all of these works share a common perspective: Algorithms are synchronous. They require a significant amount of coordination or scheduling between agents. In [13], [14], agents operate in parallel, whereas [15] proposes a sequential version of ADMM where agents operate one after the other in a predetermined order.

Contributions. This paper introduces a novel class of distributed algorithms to solve (1). The algorithms are asynchronous in the sense that some components of the network are allowed to wake up at random and perform local updates, while the rest of the network stands still. No coordinator or global clock is needed. The frequency of activation of the various network components is likely to vary.

The mathematical framework of the paper is provided by the monotone operator theory (see *e.g.* [18]). The approach relies on the introduction of *randomized Gauss-Seidel* iterations of the so called Douglas-Rachford operator. These iterations stem from a new powerful method for finding the zeros of a sum of two monotone operators. Application of our method to Problem (1) yields a randomized ADMM-like algorithm, which is proved to converge to the sought minimizers.

The paper is organized as follows. The distributed optimization problem is rigorously stated in Section II. The synchronous ADMM algorithm that solves this problem is then described in Section III. Section IV forms the core of the paper. After quickly recalling the monotone operator formalism, the random Gauss-Seidel form of the proximal algorithm is described and its convergence is shown there. These results will eventually lead to an asynchronous version of the well-known Douglas-Rachford splitting algorithm. In Section V, the results of Section IV are applied towards developing an asynchronous version of the ADMM algorithm. An implementation example is finally provided in Section VI along with some simulations in Section VII.

Notations. Consider a non-directed graph $G = (V, E)$ where V is a set of vertices and E a set of edges. We sometimes note $v \sim w$ when $\{v, w\} \in E$. For any $A \subset V$,

This work was partially funded by the French Defense Agency (DGA).
The authors are with Telecom ParisTech, CNRS LTCI, 75013 Paris, France {lastname}@telecom-paristech.fr

we denote by $G(A)$ the subgraph of G induced by A (i.e., $G(A)$ has vertices A and for any $(v, w) \in A^2$, $\{v, w\}$ is an edge of $G(A)$ if and only if it is an edge of G). Let X be a Euclidean space. We denote by X^A the set of functions on $A \rightarrow \mathsf{X}$. It is endowed with the inner product $\langle x, y \rangle_A = \sum_{v \in A} \langle x(v), y(v) \rangle_{\mathsf{X}}$ where $\langle \cdot, \cdot \rangle_{\mathsf{X}}$ is the inner product on X . We will omit subscripts X and A when no confusion occurs. For any finite collection $A_1, \dots, A_L \subset V$, we endow the space $\mathsf{X}^{A_1} \times \dots \times \mathsf{X}^{A_L}$ with the scalar product $\langle x, y \rangle = \sum_{\ell=1}^L \langle x_\ell, y_\ell \rangle_{A_\ell}$ for any $x = (x_1, \dots, x_L)$ and $y = (y_1, \dots, y_L)$.

We denote by $\Pi_A x$ the restriction of x to A i.e., $\Pi_A : \mathsf{X}^V \rightarrow \mathsf{X}^A$ is the linear operator defined for any $x \in \mathsf{X}^V$ as $\Pi_A x : (v \in A) \mapsto x(v)$. We denote by $\mathbf{1}_A$ the vector of X^A whose entries are all equal to one and by $\text{sp}(\mathbf{1}_A)$ the linear span of $\mathbf{1}_A$ i.e., the set of vectors of X^A whose entries are all equal. Notation $|A|$ represents the cardinal of a set A .

For a closed proper convex function $h : \mathsf{X} \rightarrow (-\infty, +\infty]$ we define the proximal operator $\text{prox}_{h, \rho}(x) = \arg \min_y h(y) + \frac{\rho}{2} \|y - x\|^2$.

II. DISTRIBUTED OPTIMIZATION ON A GRAPH

Consider a network of agents represented by a non-oriented graph $G = (V, E)$ where V is a finite set of vertices (i.e., the agents) and E is a set of edges. Each agent $v \in V$ has a private cost function $f_v : \mathsf{X} \rightarrow (-\infty, +\infty]$ where X is a Euclidean space. We start with the following assumption:

Assumption 1:

- i) For all $v \in V$, f_v is a proper closed convex function.
- ii) The infimum in (1) is finite and is attained at some point $x^* \in \mathsf{X}$.

In order to solve the optimization problem (1) on the graph G , we first provide an equivalent formulation of (1) that will be revealed useful. For some integer $L \geq 1$, consider a finite collection A_1, A_2, \dots, A_L of subsets of V which we shall refer to as *components*. We assume the following condition.

Assumption 2: i) $\bigcup_{\ell=1}^L A_\ell = V$.

ii) $\bigcup_{\ell=1}^L G(A_\ell)$ is connected.

Assumption 2i) implies that any vertex appears in one of the components A_1, \dots, A_L at least. We stress the fact that two distinct components A_ℓ and $A_{\ell'}$ are not necessarily disjoint, though. Assumption 2ii) means that the union of all subgraphs is connected. As the latter union is also a subgraph of G , this implies that G is connected. As will be made clear below, our method will assume that all agents in the same component are able to perform simple operations in a coordinated fashion (i.e., compute a local average over a component). Thus, in practice, it is reasonable to require that each subgraph $G(A_\ell)$ is itself connected.

We introduce some notations. We set for any $x \in \mathsf{X}^V$,

$$f(x) \triangleq \sum_{v \in V} f_v(x(v)).$$

For any $z = (z_1, \dots, z_L) \in \mathsf{Z} \triangleq \mathsf{X}^{A_1} \times \dots \times \mathsf{X}^{A_L}$, we define the closed proper convex function

$$g(z) \triangleq \sum_{\ell=1}^L \iota_{\text{sp}(\mathbf{1}_{A_\ell})}(z_\ell)$$

where ι_H is the indicator function of a set H (equal to zero on H and to $+\infty$ outside). Here $g(z)$ is equal to zero if for every ℓ , z_ℓ is constant. Otherwise, $g(z)$ is infinite. For any $x \in \mathsf{X}^V$, we define $Mx \triangleq (\Pi_{A_1} x, \dots, \Pi_{A_L} x)$. We consider the following optimization problem:

$$\inf_{x \in \mathsf{X}^V} f(x) + g(Mx) \quad (2)$$

Lemma 1: Under Assumption 2, x is a minimizer of (2) if and only if $x = \bar{x} \mathbf{1}_V$ where $\bar{x} \in \mathsf{X}$ is a minimizer of (1).

Proof: Let $x \in \mathsf{X}^V$ such that $g(Mx)$ is finite. Then x is constant on each component. Let v, w be two arbitrary vertices in V . There exists a path in $\bigcup_{\ell=1}^L G(A_\ell)$ connecting v and w . Each edge of this path connects two vertices which belong to a common component. Thus, x is constant on two consecutive vertices of the path. This proves that $x(v) = x(w)$. Thus, x is constant and the result follows. ■

One possible approach for this problem is to use ADMM. Although the choice of the sets A_1, \dots, A_L does not change the minimizers of the initial problem, it has an impact on the particular form of ADMM used to find these minimizers, as we shall see below.

In order to be more explicit, we provide in this section two important examples of possible choices for the components A_1, \dots, A_L .

Example 1: Let $L = 1$ and $A_1 = V$. Problem (2) becomes

$$\inf_{x \in \mathsf{X}^V} f(x) + \iota_{\text{sp}(\mathbf{1}_V)}(x),$$

In this case, the formulation is identical to [11, Chapter 7].

Example 2: Let $L = |E|$ and $\{A_1, \dots, A_L\} = E$. That is, each set A_ℓ is a pair of vertices $\{v, w\}$ such that $\{v, w\}$ is an edge. Here, Problem (2) is written as

$$\inf_{x \in \mathsf{X}^V} f(x) + \sum_{v \sim w} \iota_{\text{sp}(\mathbf{1}_2)} \begin{pmatrix} x(v) \\ x(w) \end{pmatrix}$$

where $\mathbf{1}_2$ stands for the vector $(1, 1)^T$.

III. SYNCHRONOUS ADMM

A. General facts

We now apply the standard ADMM to Problem (2). Perhaps the most direct way to describe ADMM is to reformulate the unconstrained problem (2) into the following constrained problem: Minimize $f(x) + g(z)$ subject to $z = Mx$. For any $x \in \mathsf{X}^V$, $\lambda, z \in \mathsf{Z}$, the augmented Lagrangian is given by

$$\mathcal{L}_\rho(x, z; \lambda) \triangleq f(x) + g(z) + \langle \lambda, Mx - z \rangle + \frac{\rho}{2} \|Mx - z\|^2 \quad (3)$$

where $\rho > 0$ is a constant. ADMM consists of the iterations

$$x^{k+1} = \underset{x \in \mathsf{X}^V}{\text{argmin}} \mathcal{L}_\rho(x, z^k; \lambda^k) \quad (4a)$$

$$z^{k+1} = \underset{z \in \mathsf{Z}}{\text{argmin}} \mathcal{L}_\rho(x^{k+1}, z; \lambda^k) \quad (4b)$$

$$\lambda^{k+1} = \lambda^k + \rho (Mx^{k+1} - z^{k+1}). \quad (4c)$$

From [11, Chap. 3.2], the following result is immediate.

Theorem 1: Under Assumption 1, the sequence $\{x^k\}_{k \in \mathbb{N}}$ defined in (4a) converges to a minimizer of (2).

B. Decentralized Implementation

One should now make (4) more explicit and convince the reader that the iterations are indeed amenable to distributed implementation. Due to the specific form of function g , it is clear from (4b) that all components z_1^k, \dots, z_L^k of z^k are constant. Otherwise stated, $z^k = (\bar{z}_1^k 1_{A_1}, \dots, \bar{z}_L^k 1_{A_L})$ for some constants $\bar{z}_\ell^k \in X$. For any $v \in V$, we set

$$\sigma(v) \triangleq \{\ell : v \in A_\ell\}.$$

Now consider the first update equation (4a). Getting rid of all quantities in \mathcal{L}_ρ which do not depend on the v th component of x , we obtain for any $v \in V$

$$x^{k+1}(v) = \operatorname{argmin}_{y \in X} f_v(y) + \sum_{\ell \in \sigma(v)} \langle \lambda_\ell^k(v), y \rangle + \frac{\rho}{2} \|y - \bar{z}_\ell^k\|^2.$$

After some algebra, the above equation further simplifies to

$$x^{k+1}(v) = \operatorname{prox}_{f_v, \rho|\sigma(v)|} (Z^k(v) - B^k(v)) \quad (5)$$

where we introduced the following parameters:

$$Z^k(v) = \frac{1}{|\sigma(v)|} \sum_{\ell \in \sigma(v)} \bar{z}_\ell^k, \quad B^k(v) = \frac{1}{\rho|\sigma(v)|} \sum_{\ell \in \sigma(v)} \lambda_\ell^k(v). \quad (6)$$

It is straightforward to show that the second update equation (4b) admits as well a simple decomposable form. After some algebra, we obtain that for any $\ell = 1, \dots, L$,

$$\bar{z}_\ell^{k+1} = \frac{1}{|A_\ell|} \sum_{v \in A_\ell} x^{k+1}(v) + \frac{\lambda_\ell^k(v)}{\rho}. \quad (7)$$

Finally, for all $\ell = 1, \dots, L$ and $v \in A_\ell$, equation (4c) reads

$$\lambda_\ell^{k+1}(v) = \lambda_\ell^k(v) + \rho(x^{k+1}(v) - \bar{z}_\ell^{k+1}). \quad (8)$$

Averaging (8) w.r.t. v and using (7) yields $\sum_{v \in A_\ell} \lambda_\ell^k(v) = 0$. Thus, the second term in the RHS of (7) can be deleted. Finally, averaging (8) w.r.t. ℓ leads to

$$B^{k+1}(v) = B^k(v) + x^{k+1}(v) - Z^{k+1}(v). \quad (9)$$

Synchronous ADMM:

At each iteration k ,

For each agent v , compute $x^{k+1}(v)$ using (5).

In each components $\ell = 1, \dots, L$, compute

$$\bar{z}_\ell^{k+1} = \frac{1}{|A_\ell|} \sum_{w \in A_\ell} x^{k+1}(w).$$

For each agent v , compute $Z^{k+1}(v)$ and $B^{k+1}(v)$ using (6) and (9) respectively.

The above algorithm implicitly requires the existence of a routine for computing an average, in each component A_ℓ . This requirement is mild when the components coincide with edges of the graph as in Example 2. In this case, one only needs that the two vertices of an edge share their current estimate and find an agreement on the average. In the general case, the objective can be achieved by selecting a leader in each component whose role is to gather the estimates,

compute the average and send the result to all agents in this component.

It is worth noting that in the case of Example 1, the synchronous ADMM described above coincides with the algorithm of [11] and is very close to the so called PPXA algorithm introduced in [12].

IV. A RANDOMIZED PROXIMAL ALGORITHM

A. Monotone operators

An operator T on a Euclidean space Y is a set valued mapping $T : Y \rightarrow 2^Y$. An operator can be equivalently identified with a subset of $Y \times Y$, which makes us write $(x, y) \in T$ when $y \in T(x)$. Given two operators T_1 and T_2 on Y and two real numbers α_1 and α_2 , the operator $\alpha_1 T_1 + \alpha_2 T_2$ is defined as $\alpha_1 T_1 + \alpha_2 T_2 = \{(x, \alpha_1 y_1 + \alpha_2 y_2) : (x, y_1) \in T_1, (x, y_2) \in T_2\}$. The identity operator is $I = \{(x, x) : x \in Y\}$ and the inverse of the operator T is $T^{-1} = \{(x, y) : (y, x) \in T\}$. The operator T is said *monotone* if

$$\forall (x, y), (x', y') \in T, \langle x - x', y - y' \rangle \geq 0.$$

A monotone operator is said *maximal* if it is not strictly contained in any monotone operator (as a subset of $Y \times Y$). Finally, T is said *firmly non-expansive* if

$$\forall (x, y), (x', y') \in T, \langle x - x', y - y' \rangle \geq \|y - y'\|^2.$$

The typical example of a monotone operator is the subdifferential ∂f of a convex function $f : Y \rightarrow \mathbb{R}$. Finding a minimum of f amounts to finding a point in $\operatorname{zer}(\partial f)$, where $\operatorname{zer}(T) = \{x : 0 \in T(x)\}$ is the set of zeroes of an operator T . A common technique for finding a zero of a maximal monotone operator T is the so-called *proximal point algorithm* [16] that we now describe. The *resolvent* of T is the operator $J_{\rho T} \triangleq (I + \rho T)^{-1}$ for $\rho > 0$. One key result (see e.g. [9]) says that T is maximal monotone if and only if $J_{\rho T}$ is firmly non expansive and its domain is Y . Observe that a firmly non expansive operator is single valued and denote by $\operatorname{fix}(J_{\rho T})$ the set of fixed points of $J_{\rho T}$. It is clear that $\operatorname{fix}(J_{\rho T}) = \operatorname{zer}(T)$. The firm non expansiveness of $J_{\rho T}$ plays a central role in the proof of the following result:

Lemma 2 (Proximal point algorithm [16]): If T is a maximal monotone operator and $\rho > 0$, then the iterates $\zeta^{k+1} = J_{\rho T}(\zeta^k)$ starting at any point of Y converge to a point of $\operatorname{fix}(J_{\rho T})$ whenever this set is non-empty.

B. Random Gauss-Seidel iterations

Assume now that the Euclidean space Y is a Cartesian product of Euclidean spaces of the form $Y = Y_1 \times \dots \times Y_L$ where L is a given integer, and write any $\zeta \in Y$ as $\zeta = (\zeta_1, \dots, \zeta_L)$ where $\zeta_\ell \in Y_\ell$ for $\ell = 1, \dots, L$. Let S be a firmly non expansive operator on Y and write

$$S(\zeta) = (S_1(\zeta), \dots, S_L(\zeta))$$

where $S_\ell(\zeta) \in Y_\ell$. For $\ell = 1, \dots, L$, define the single valued operator $\hat{S}_\ell : Y \rightarrow Y$ as

$$\hat{S}_\ell(\zeta) = (\zeta_1, \dots, \zeta_{\ell-1}, S_\ell(\zeta), \zeta_{\ell+1}, \dots, \zeta_L). \quad (10)$$

Considering an iterative algorithm of the form $\zeta^{k+1} = S(\zeta^k)$, its *Gauss-Seidel* version would be an algorithm of the form $\zeta^{k+1} = \hat{S}_L \circ \dots \circ \hat{S}_1(\zeta^k)$. We are interested here in a *randomized version* of these iterates. On a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, let $(\xi^k)_{k \in \mathbb{N}}$ be a random process satisfying the following assumption:

Assumption 3: The random variables ξ^k are independent and identically distributed. They are valued in the set $\{1, \dots, L\}$ with $\mathbb{P}[\xi^1 = \ell] = p_\ell > 0$ for all $\ell = 1, \dots, L$. We are interested here in the convergence of the random iterates $\zeta^{k+1} = \hat{S}_{\xi^{k+1}}(\zeta^k)$ towards a (generally random) point of $\text{fix}(S)$, provided this set is non empty:

Theorem 2 (Main result): Let S be a firmly non-expansive operator on Y with domain Y . Let $\{\xi^k\}_{k \in \mathbb{N}}$ be a sequence of random variables satisfying Assumption 3. Assume that $\text{fix}(S) \neq \emptyset$. Then for any initial value ζ^0 , the sequence of iterates $\zeta^{k+1} = \hat{S}_{\xi^{k+1}}(\zeta^k)$ converges almost surely to a random variable supported by $\text{fix}(S)$.

Proof: Denote by $\langle \zeta, \eta \rangle = \sum_{\ell=1}^L \langle \zeta_\ell, \eta_\ell \rangle_{Y_\ell}$ the inner product of Y , and by $\|\zeta\| = \langle \zeta, \zeta \rangle^{1/2}$ its associated norm. Define a new inner product $\zeta \bullet \eta = \sum_{\ell=1}^L p_\ell^{-1} \langle \zeta_\ell, \eta_\ell \rangle_{Y_\ell}$ on Y , and let $\|\zeta\| = (\zeta \bullet \zeta)^{1/2}$ be its associated norm. Fix ζ^* in $\text{fix}(S)$. Conditionally to the sigma-field $\mathcal{F}_k = \sigma(\xi^1, \dots, \xi^k)$ we have

$$\begin{aligned} \mathbb{E}[\|\zeta^{k+1} - \zeta^*\|^2 | \mathcal{F}_k] &= \sum_{\ell=1}^L p_\ell \left\| \hat{S}_\ell(\zeta^k) - \zeta^* \right\|^2 \\ &= \sum_{\ell=1}^L p_\ell \left(\frac{1}{p_\ell} \|S_\ell(\zeta^k) - \zeta_\ell^*\|_{Y_\ell}^2 + \sum_{\substack{i=1 \\ i \neq \ell}}^L \frac{1}{p_i} \|\zeta_i^k - \zeta_i^*\|_{Y_i}^2 \right) \\ &= \|S(\zeta^k) - \zeta^*\|^2 + \sum_{\ell=1}^L \frac{1-p_\ell}{p_\ell} \|\zeta_\ell^k - \zeta_\ell^*\|_{Y_\ell}^2 \\ &= \|\zeta^k - \zeta^*\|^2 + \|S(\zeta^k) - \zeta^*\|^2 - \|\zeta^k - \zeta^*\|^2 \end{aligned}$$

Since $(I - S)(\zeta^*) = 0$, we have

$$\begin{aligned} &\|S(\zeta^k) - \zeta^*\|^2 - \|\zeta^k - \zeta^*\|^2 \\ &= \|S(\zeta^k) - \zeta^k + \zeta^k - \zeta^*\|^2 - \|\zeta^k - \zeta^*\|^2 \\ &= \|S(\zeta^k) - \zeta^k\|^2 + 2\langle S(\zeta^k) - \zeta^k, \zeta^k - \zeta^* \rangle \\ &= \|S(\zeta^k) - \zeta^k\|^2 - 2\langle (I - S)(\zeta^k) - (I - S)(\zeta^*), \zeta^k - \zeta^* \rangle \\ &\leq -\|S(\zeta^k) - \zeta^k\|^2 \end{aligned}$$

where the inequality comes from the easily verifiable fact that $(I - S)$ is firmly non-expansive when S is. This leads to the inequality

$$\mathbb{E}[\|\zeta^{k+1} - \zeta^*\|^2 | \mathcal{F}_k] \leq \|\zeta^k - \zeta^*\|^2 - \|S(\zeta^k) - \zeta^k\|^2 \quad (11)$$

which shows that $\|\zeta^k - \zeta^*\|^2$ is a nonnegative supermartingale with respect to the filtration (\mathcal{F}_k) . As such, it converges with probability one towards a random variable X_{ζ^*} satisfying $0 \leq X_{\zeta^*} < \infty$ almost everywhere. Given a countable dense subset H of $\text{fix}(S)$, there is a probability one set on which $\|\zeta^k - \zeta\| \rightarrow X_\zeta \in [0, \infty)$ for all $\zeta \in H$. Let $\zeta^* \in \text{fix}(S)$, let $\varepsilon > 0$, and choose $\zeta \in H$ such that

$\|\zeta^* - \zeta\| \leq \varepsilon$. With probability one, we have

$$\|\zeta^k - \zeta^*\| \leq \|\zeta^k - \zeta\| + \|\zeta - \zeta^*\| \leq X_\zeta + 2\varepsilon$$

for k large enough. Similarly, $\|\zeta^k - \zeta^*\| \geq X_\zeta - 2\varepsilon$ for k large enough. We therefore obtain:

C1 : There is a probability one set on which $\|\zeta^k - \zeta^*\|$ converges for every $\zeta^* \in \text{fix}(S)$.

Getting back to Inequality (11), taking the expectations on both sides of this inequality and iterating over k , we obtain

$$\sum_{k=0}^{\infty} \mathbb{E}[\|S(\zeta^k) - \zeta^k\|^2] \leq (\zeta^0 - \zeta^*)^2.$$

By Markov's inequality and Borel Cantelli's lemma, we therefore obtain:

C2 : $S(\zeta^k) - \zeta^k \rightarrow 0$ almost surely.

We now consider an elementary event in the probability one set where **C1** and **C2** hold. On this event, since $\|\zeta^k - \zeta^*\|$ converges for $\zeta^* \in \text{fix}(S)$, the sequence $\{\zeta^k\}_{k \in \mathbb{N}}$ is bounded. Since S is firmly non expansive, it is continuous, and **C2** shows that all the accumulation points of $\{\zeta^k\}_{k \in \mathbb{N}}$ are in $\text{fix}(S)$. It remains to show that these accumulation points reduce to one point. Assume that ζ_1^* is an accumulation point. By **C1**, $\|\zeta^k - \zeta_1^*\|$ converges. Therefore, $\lim \|\zeta^k - \zeta_1^*\| = \liminf \|\zeta^k - \zeta_1^*\| = 0$, which shows that ζ_1^* is unique. ■

V. RANDOM ADMM

We now return to the optimization problem (2). It is a well known fact that the standard ADMM can be seen as a special case of the so-called Douglas-Rachford algorithm [9]. The Douglas-Rachford algorithm can itself be seen as a special case of a proximal point algorithm. By the results of the previous section, this suggests that random Gauss-Seidel iterations applied to the Douglas-Rachford operator produce a sequence which eventually converges to the sought solutions. It turns out that the latter random iterations can be written under the form of practical asynchronous ADMM-like algorithm.

A. Douglas-Rachford operator

Consider the following dual problem associated with (2)

$$\min_{\lambda \in Z} f^*(-M^*\lambda) + g^*(\lambda), \quad (12)$$

where f^*, g^* are the Fenchel conjugates of f and g and M^* is the adjoint of M . By Assumption 1 along with [17, Th.3.3.5], the minimum in (12) is attained and its opposite coincides with the minimum of (2). Note that λ is a minimizer of (12) iff zero belongs to the subdifferential of the objective function in (12). By [17, Th.3.3.5] again, this reads $0 \in -M \cdot \partial f^*(-M^*\lambda) + \partial g^*(\lambda)$. Otherwise stated, finding minimizers of the dual problem (12) boils down to searching zeros of the sum of two maximal monotone operators $T + U$ defined by $T = -M \cdot \partial f^* \circ (-M^*)$ and $U = \partial g^*$. For a fixed $\rho > 0$, the Douglas-Rachford (or Lions-Mercier) operator R is defined as

$$\{(\nu + \rho b, \mu - \nu) : (\mu, b) \in U, (\nu, a) \in T, \nu + \rho a = \mu - \rho b\}.$$

The following Lemma is an immediate consequence of [9].

Lemma 3: Under Assumption 1, R is maximal monotone, and $\text{zer}(R) \neq \emptyset$. Moreover, $J_{\rho U}(\zeta) \in \text{zer}(T + U)$ for any $\zeta \in \text{zer}(R)$.

Lemma 3 implies that the search for a zero of $T + U$ boils down to the search of a zero of R up to a resolvent step $J_{\rho U}$. To that end, a standard approach is to use a proximal point algorithm of the form $\zeta^{k+1} = J_R(\zeta^k)$. By [9], it can be shown that this approach is equivalent to the ADMM derived in Section II. Here, our aim is different. We shall consider random Gauss-Seidel iterations in order to derive an asynchronous version of the ADMM.

B. Random Gauss-Seidel Iterations

Define $S \triangleq J_R$ as the resolvent associated with the Douglas-Rachford operator R . On the space $Z = X^{A_1} \times \cdots \times X^{A_L}$, define the operator \hat{S}_ℓ as in (10) for any $\ell = 1, \dots, L$. Let $\{\zeta^k\}_{k \in \mathbb{N}}$ be a random process satisfying Assumption 3. The following result is a consequence of Theorem 2 combined with Lemma 3.

Theorem 3: Let Assumptions 1, 2 and 3 hold true. Consider the sequence $\{\zeta^k\}_{k \in \mathbb{N}}$ defined by $\zeta^{k+1} = \hat{S}_{\xi^{k+1}}(\zeta^k)$. Then for any initial value ζ^0 , the sequence $\lambda^k \triangleq J_{\rho U}(\zeta^k)$ converges almost surely to a minimizer of (12).

In order to complete the above result, we still must justify the fact that, as claimed, the above iterations can be seen as an asynchronous distributed algorithm.

C. Distributed Algorithm

We make the above random Gauss-Seidel iterations more explicit. In the sequel we shall always denote by ζ_ℓ the ℓ th component of a function $\zeta \in Z$ i.e., $\zeta = (\zeta_1, \dots, \zeta_L)$. For any ℓ , we introduce the average $\bar{\zeta}_\ell = \sum_{v \in A_\ell} \zeta_\ell(v) / |A_\ell|$. Lemma 4 below states that any $\zeta \in Z$ is uniquely represented by a couple $(\lambda, z) \in U$ whose expression is provided. Moreover, it provides the explicit form of the ℓ th block S_ℓ of the resolvent S . This will be the basis of our asynchronous distributed algorithm.

Lemma 4: For any $\zeta \in Z$, the following holds true:

- i) There exists a unique $(\lambda, z) \in U$ such that $\lambda + \rho z = \zeta$.
- ii) $J_{\rho U}(\zeta) = \lambda$.
- iii) For any $\ell = 1, \dots, L$, $\lambda_\ell = \zeta_\ell - \bar{\zeta}_\ell 1_{A_\ell}$ and $z_\ell = \frac{\bar{\zeta}_\ell}{\rho} 1_{A_\ell}$.
- iv) For any $\ell = 1, \dots, L$, and any $v \in A_\ell$

$$S_\ell(\zeta) : v \mapsto \lambda_\ell(v) + \rho x(v) \quad (13)$$

where $x(v)$ is defined by

$$x(v) = \text{prox}_{f_v, \rho|\sigma(v)|} \left(\frac{1}{|\sigma(v)|} \sum_{\ell \in \sigma(v)} \bar{z}_\ell - \frac{\lambda_\ell(v)}{\rho} \right). \quad (14)$$

Proof: i)-ii) Since $J_{\rho U}$ is the resolvent of a maximal monotone operator, $J_{\rho U}(\zeta)$ exists and is unique for every $\zeta \in Z$. Let us define $\lambda = J_{\rho U}(\zeta)$ and $z = (\zeta - \lambda) / \rho$. Trivially, $\lambda + \rho z = \zeta$. As $\zeta \in \lambda + \rho U(\lambda)$, we get that $(\lambda, z) \in U$. This decomposition is unique since for any $(\lambda, z) \in U$ satisfying $\lambda + \rho z = \zeta$, one has $\zeta \in (I + \rho U)(\lambda)$ and thus $\lambda = J_{\rho U}(\zeta)$.

iii) We use $\lambda = J_{\rho U}(\zeta) = \text{prox}_{g^*, \rho}(\zeta) = \zeta - 1/\rho \text{prox}_{g, 1/\rho}(\rho\zeta)$ (see [18, Th. 14.3]). As g is the indicator function of the set $\text{sp}(1_{A_1}) \times \cdots \times \text{sp}(1_{A_L})$, $\text{prox}_{g, 1/\rho}$ coincides with the projection operator onto that set. Thus, for any ℓ , $\lambda_\ell = \zeta_\ell - \bar{\zeta}_\ell 1_{A_\ell}$. The expression of z follows from $z = (\zeta - \lambda) / \rho$.

iv) Operator $S = J_R$ can be written as

$$\{(\mu + \rho b, \nu + \rho b) : (\mu, b) \in U, (\nu, a) \in T, \nu + \rho a = \mu - \rho b\}.$$

Moreover, as R is monotone, $S(\zeta)$ is a singleton. Representing $\zeta = \lambda + \rho z$ with $(\lambda, z) \in U$, it follows from the above expression of S that $S(\zeta) = \nu + \rho z$ where ν is such that $\nu + \rho a = \lambda - \rho z$ for some $a \in T(\nu)$. Using $T = -M \cdot \partial f^* \circ (-M^*)$, condition $a \in T(\nu)$ translates to: there exists $x \in \partial f^*(M^*\nu)$ s.t. $a = -Mx$. The output-resolvent is obtained by $\nu + \rho z = \lambda + \rho Mx$. For a given component ℓ , this boils down to equation (13). The remaining task is to provide the expression of x . By the Fenchel-Young equality $\partial f^* = \partial f^{-1}$ [17, Prop.3.3.4], condition $x \in \partial f^*(M^*\nu)$ is equivalent to $-M^*\nu \in \partial f(x)$. Using that $\nu = \lambda - \rho(z - Mx)$, we obtain $0 \in \partial f(x) + M^*\lambda + \rho M^*(Mx - z)$. Otherwise stated, $x = \arg \min_{y \in X^V} \mathcal{L}_\rho(y, z; \lambda)$ where \mathcal{L}_ρ is the augmented Lagrangian defined in (3). Using the results of Section III, $x(v)$ is given by (14) for any v . ■

We are now in position to state the main algorithm. It simply consists in an explicit writing of the random Gauss-Seidel iterations $\zeta^{k+1} = \hat{S}_{\xi^{k+1}}(\zeta^k)$ using Lemma 4iv). Note that, by Lemma 4i), the definition of a sequence $(\zeta_k)_k$ on Z is equivalent to the definition of two sequences $\{\lambda^k\}_{k \in \mathbb{N}}$ and $\{z^k\}_{k \in \mathbb{N}}$ such that for all k , $\zeta^k = \lambda^k + \rho z^k$ with $(\lambda^k, z^k) \in U$. Moreover, by Lemma 4iii), each component z_ℓ^k of z^k is a constant. The definition of z^k thus reduces to the definition of L constants $\bar{z}_1^k, \dots, \bar{z}_L^k$ in X .

Asynchronous ADMM:

At each iteration k , draw r.v. ξ^{k+1} .

For $\ell = \xi^{k+1}$, set for any $v \in A_\ell$:

$$\begin{aligned} x^{k+1}(v) &= \text{prox}_{f_v, \rho|\sigma(v)|} \left(\frac{1}{|\sigma(v)|} \sum_{\ell \in \sigma(v)} \bar{z}_\ell^k - \frac{\lambda_\ell^k(v)}{\rho} \right) \\ \bar{z}_\ell^{k+1} &= \frac{1}{|A_\ell|} \sum_{w \in A_\ell} x^{k+1}(w) \\ \lambda_\ell^{k+1}(v) &= \lambda_\ell^k(v) + \rho (x^{k+1}(v) - \bar{z}_\ell^{k+1}). \end{aligned}$$

For any $\ell \neq \xi^{k+1}$, set $\lambda_\ell^{k+1} = \lambda_\ell^k$.

For any $w \notin A_{\xi^{k+1}}$, set $x^{k+1}(w) = x^k(w)$.

VI. IMPLEMENTATION EXAMPLE

In order to illustrate our results, we consider herein an asynchronous version of the ADMM algorithm in the context of Section II-Example 2. The scenario is the following: first, Agent $v \in \{1, \dots, |V|\}$ wakes up at time $k + 1$ with the probability q_v . Denoting by \mathcal{N}_v the neighborhood of Agent v in the Graph G , this agent then chooses one of its neighbors, say w , with the probability $1/|\mathcal{N}_v|$ and sends an activation message to w . In this setting, the edge $\{v, w\}$ coincides with

one of the A_ℓ of Example 2 in Section II. It is easily seen that the samples of the activation process $\{\xi^k\}_{k \in \mathbb{N}}$ which is of course valued in $\{1, \dots, |E|\}$ are governed by the probability law

$$\mathbb{P}[\xi^1 = \{v, w\}] = \frac{q_v}{|\mathcal{N}_v|} + \frac{q_w}{|\mathcal{N}_w|} > 0.$$

When the edge $\{v, w\}$ is activated, the following two prox(\cdot) operations are performed by the agents:

$$x^{k+1}(v) = \text{prox}_{f_v, \rho|\mathcal{N}(v)|} \left(\frac{1}{|\mathcal{N}(v)|} \sum_{\ell \in \mathcal{N}(v)} \bar{z}_\ell^k - \frac{\lambda_\ell^k(v)}{\rho} \right)$$

$$x^{k+1}(w) = \text{prox}_{f_w, \rho|\mathcal{N}(w)|} \left(\frac{1}{|\mathcal{N}(w)|} \sum_{\ell \in \mathcal{N}(w)} \bar{z}_\ell^k - \frac{\lambda_\ell^k(w)}{\rho} \right).$$

The two agents exchange then the values $x^{k+1}(v)$ and $x^{k+1}(w)$ and perform the following operations:

$$\bar{z}_\ell^{k+1} = \frac{x^{k+1}(v) + x^{k+1}(w)}{2}$$

$$\lambda_\ell^{k+1}(v) = \lambda_\ell^k(v) + \rho \frac{x^{k+1}(v) - x^{k+1}(w)}{2}$$

$$\lambda_\ell^{k+1}(w) = \lambda_\ell^k(w) + \rho \frac{x^{k+1}(w) - x^{k+1}(v)}{2}.$$

We remark that this communication scheme is reminiscent of the so-called *Random Gossip* algorithm introduced in [19] in the context of distributed averaging.

VII. NUMERICAL RESULTS

We consider a network with $V = \{1, \dots, 5\}$ and with $E = \{\{1, 2\}, \{2, 3\}, \{3, 4\}, \{4, 5\}, \{5, 3\}\}$. We evaluate the behavior of: i) the synchronous *Distributed Gradient* with $1/k$ step size [20] using Metropolis exchange matrix; ii) the *Asynchronous Distributed Gradient* with *Random Gossip* as a communication algorithm [19]; iii) the *Synchronous ADMM*; and iv) the *Asynchronous ADMM*. Each agent maintains a different quadratic convex function and their goal is to reach consensus over the minimizer of problem (1).

In Figure 1, we plot the squared error versus the number of iterations for the considered algorithms. We observe that our algorithm clearly outperforms the distributed gradient descents.

ACKNOWLEDGEMENT

The authors would like to thank Jean-Christophe Pesquet from Université Paris-Est for his useful advice.

REFERENCES

- [1] P. A. Forero, A. Cano, and G. B. Giannakis, "Distributed Clustering Using Wireless Sensor Networks," *IEEE Journal of Selected Topics in Signal Processing*, vol. 5, no. 4, pp. 707–724, 2011.
- [2] D. P. Bertsekas and J. N. Tsitsiklis, *Parallel and distributed computation*, Prentice Hall Inc., Old Tappan, 1989.
- [3] S. Ram, A. Nedic, and V. Veeravalli, "Distributed Stochastic Sub-gradient Projection Algorithms for Convex Optimization," *Journal of Optimization Theory and Applications*, vol. 147, no. 3, pp. 516–545, 2010.
- [4] P. Bianchi and J. Jakubowicz, "Convergence of a multi-agent projected stochastic gradient algorithm for non-convex optimization," *IEEE Transactions on Automatic Control*, vol. 58, no. 2, 2013.
- [5] D. Jakovetić, J. M. F. Moura, and X. Joao, "Distributed Nesterov-like gradient algorithms," in *Proc. 51st IEEE Conference on Decision and Control (CDC)*, 2012, pp. 5459–5464.

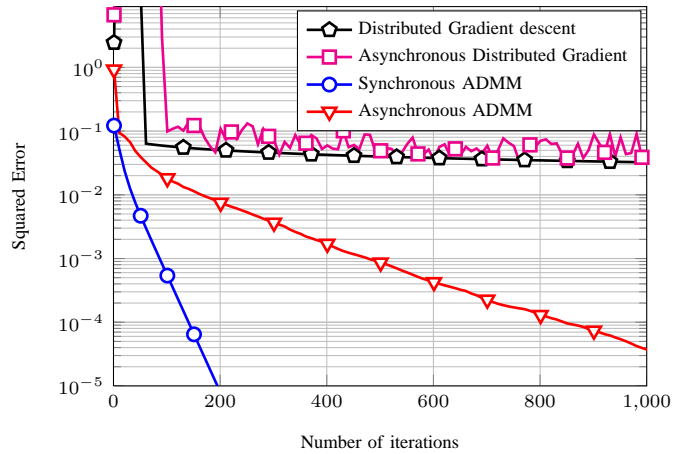


Fig. 1. Squared error of distributed optimization algorithms versus the number of primal updates

- [6] J. C. Duchi, A. Agarwal, and M. J. Wainwright, "Dual Averaging for Distributed Optimization: Convergence Analysis and Network Scaling," *IEEE Transactions on Automatic Control*, vol. 57, no. 3, pp. 592–606, 2012.
- [7] A. Jadbabaie, A. Ozdaglar, and M. Zargham, "A distributed newton method for network optimization," in *Proc. 48th IEEE Conference on Decision and Control (CDC)*, 2009, pp. 2736–2741.
- [8] P. L. Lions and B. Mercier, "Splitting algorithms for the sum of two nonlinear operators," *SIAM Journal on Numerical Analysis*, vol. 16, no. 6, pp. 964–979, 1979.
- [9] J. Eckstein and D. P. Bertsekas, "On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators," *Mathematical Programming*, vol. 55, pp. 293–318, 1992.
- [10] P. L. Combettes and J.-C. Pesquet, "Proximal Splitting Methods in Signal Processing," in *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*, H. H. Bauschke, R. Burachik, P. L. Combettes, V. Elser, D. R. Luke, and H. Wolkowicz editors. Springer-Verlag, New York, pp. 185–212, 2011.
- [11] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," vol. 3 of *Foundations and Trends in Machine Learning*, Now Publishers Inc., 2011.
- [12] P. L. Combettes and J.-C. Pesquet, "A proximal decomposition method for solving convex variational inverse problems," *Inverse Problems*, vol. 24, no. 6, Dec. 2008.
- [13] I. D. Schizas, A. Ribeiro, and G. B. Giannakis, "Consensus in Ad Hoc WSNs With Noisy Links – Part I: Distributed Estimation of Deterministic Signals," *IEEE Transactions on Signal Processing*, vol. 56, no. 1, pp. 350–364, jan. 2008.
- [14] J. F. C. Mota, J. M. F. Xavier, P. M. Q. Aguiar, and M. Puschel, "Distributed ADMM for model predictive control and congestion control," in *Proc. 51st IEEE Conference on Decision and Control (CDC)*, 2012, pp. 5110–5115.
- [15] E. Wei and A. Ozdaglar, "Distributed Alternating Direction Method of Multipliers," in *Proc. 51st IEEE Conference on Decision and Control (CDC)*, 2012, pp. 5445–5450.
- [16] R. T. Rockafellar, "Monotone operators and the proximal point algorithm," *SIAM Journal on Control and Optimization*, vol. 14, no. 5, pp. 877–898, 1976.
- [17] J. M. Borwein and A. S. Lewis, *Convex Analysis and Nonlinear Optimization: Theory and Examples*, Springer Verlag, 2006.
- [18] H. H. Bauschke and P. L. Combettes, *Convex analysis and monotone operator theory in Hilbert spaces*, Springer, 2011.
- [19] S. Boyd, A. Ghosh, B. Prabhakar, and D. Shah, "Randomized gossip algorithms," *IEEE Transactions on Information Theory*, vol. 52, no. 6, pp. 2508–2530, 2006.
- [20] J. N. Tsitsiklis, D. P. Bertsekas, and M. Athans, "Distributed asynchronous deterministic and stochastic gradient optimization algorithms," *IEEE Transactions on Automatic Control*, vol. 31, no. 9, pp. 803–812, 1986.