

An exact first-order algorithm for decentralized consensus optimization

Xinwei Sun(1301110047), Jingru Zhang(1301110029)

December 22, 2014

1 Introduction

2 Algorithm

- DGD method
- EXTRA method
- EXTRA as Corrected DGD

3 Convergence analysis

- Convex objective with Lipschitz continuous gradient
- Strongly convex with Lipschitz continuous gradient

4 Experiment

Decentralized Consensus Optimization

$$\min_{x \in \mathbb{R}^p} \bar{f}(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \quad f_i : \mathbb{R}^p \rightarrow \mathbb{R}. \quad (1)$$

EXTRA method

- use a fixed step size independent of the network size.
- convergence rate $O(\frac{1}{k})$ for general convex objection with Lipschitz differential.
- linear convergence rate for (restricted) strongly convex.

Previous Methods

- existing first-order decentralized methods:
(sub)gradient,(sub)gradient-push,fast (sub)gradient,dual averaging.
- more restrictive assumptions and worse convergence rate.
- using a fixed step size,do not converge to a solution x^* ,just in its neighborhood.

Notation

- $x_{(i)} \in \mathbb{R}^p$: local copy of the global variable x .
- $x_{(i)}^k$: its value at iteration k .
- $\mathbf{f}(\mathbf{x}) \triangleq \sum_{i=1}^n f_i(x_{(i)})$: an aggregate objective function of the local variables
-

$$\mathbf{x} \triangleq \begin{pmatrix} -x_{(1)}^T - \\ -x_{(2)}^T - \\ \vdots \\ -x_{(n)}^T - \end{pmatrix} \in \mathbb{R}^{n \times p}$$

Notation



$$\nabla \mathbf{f}(\mathbf{x}) \triangleq \begin{pmatrix} -\nabla^T f_1(x_{(1)}) \\ -\nabla^T f_2(x_{(2)}) \\ \vdots \\ -\nabla^T f_n(x_{(n)}) \end{pmatrix} \in \mathbb{R}^{n \times p}.$$

- \mathbf{x} is consensual if all of its rows are identical, i.e., $x_{(1)} = \cdots = x_{(n)}$.
- $\|A\|_G \triangleq \sqrt{\text{trace}(A^T G A)}$: G-matrix norm.
- $\text{null}\{A\} \triangleq \{x \in \mathbb{R}^n \mid Ax = 0\}$: null space of A .
- $\text{span}\{A\} \triangleq \{y \in \mathbb{R}^m \mid y = Ax, \forall x \in \mathbb{R}^n\}$: linear span of all the columns of A .

Algorithm 1 (DGD)

$$x_{(i)}^{k+1} = \sum_{j=1}^n w_{ij} x_{(j)}^k - \alpha^k \nabla f_i(x_{(i)}^k), \text{ for agent } i = 1, \dots, n. \quad (2)$$

matrix version

$$\mathbf{x}^{k+1} = W\mathbf{x}^k - \alpha^k \nabla \mathbf{f}(\mathbf{x}^k). \quad (3)$$

$x_{(i)}^k \in \mathbb{R}^p$ is the local copy of \mathbf{x} held by agent i at iteration k ,

$W = [w_{ij}] \in \mathbb{R}^{n \times n}$ is a symmetric mixing matrix

satisfying $\text{null}\{I - W\} = \text{span}\{\mathbf{1}\}$. If two agents i and j are neither neighbors nor identical, then $w_{ij} = 0$,

$$\sigma_{\max}(W - \frac{1}{n}\mathbf{1}\mathbf{1}^T) < 1,$$

$\alpha^k > 0$ is a step size for iteration k .

Dilemma of DGD

- converge slowly to an exact solution with a sequence of diminishing step sizes.
- converge faster with a fixed step size but stall at an inaccurate solution ($O(\alpha)$ -neighbourhood of a solution).

The Cause of Inexact Convergence With a Fixed Step Size

let \mathbf{x}^∞ be the limit of \mathbf{x}^k . Take the limit over k and get

$$\mathbf{x}^\infty = \mathbf{W}\mathbf{x}^\infty - \alpha \nabla \mathbf{f}(\mathbf{x}^\infty). \quad (4)$$

When α is fixed and nonzero, assuming the consensus of \mathbf{x}^∞ (namely, it has identical rows $x_{(i)}^\infty$) will mean $\mathbf{x}^\infty = \mathbf{W}\mathbf{x}^\infty$, as a result of $\mathbf{W}\mathbf{1} = \mathbf{1}$, and thus $\nabla \mathbf{f}(\mathbf{x}^\infty) = \mathbf{0}$, which is equivalent to $\nabla f_i(x_{(i)}^\infty) = 0, \forall i$, i.e., the same point $x_{(i)}^\infty$ simultaneously minimizes f_i for all agents i . This is impossible in general and is different from our objective to find a point that minimizes $\sum_{i=1}^n f_i$.

Development of EXTRA

Consider the DGD update at iterations $k+1$ and k as follows

$$\mathbf{x}^{k+2} = \mathbf{W}\mathbf{x}^{k+1} - \alpha \nabla \mathbf{f}(\mathbf{x}^{k+1}), \quad (5)$$

$$\mathbf{x}^{k+1} = \mathbf{W}\mathbf{x}^k - \alpha \nabla \mathbf{f}(\mathbf{x}^k), \quad (6)$$

$$\tilde{\mathbf{W}} = \frac{\mathbf{I} + \mathbf{W}}{2}, \text{ the choice of } \tilde{\mathbf{W}} \text{ will be generalized later.} \quad (7)$$

Subtracting the above two iterations of DGD, we get the update formula of EXTRA:

$$\mathbf{x}^{k+2} - \mathbf{x}^{k+1} = \mathbf{W}\mathbf{x}^{k+1} - \tilde{\mathbf{W}}\mathbf{x}^k - \alpha \nabla \mathbf{f}(\mathbf{x}^{k+1}) + \alpha \nabla \mathbf{f}(\mathbf{x}^k). \quad (8)$$

Conclusion 1

Provided that $\text{null}\{\mathbf{I} - \mathbf{W}\} = \text{span}\{\mathbf{1}\}$, $\tilde{\mathbf{W}} = \frac{\mathbf{I} + \mathbf{W}}{2}$, $\mathbf{1}^\top(\mathbf{W} - \tilde{\mathbf{W}}) = \mathbf{0}$ and the continuity of ∇f , if a sequence following EXTRA converges to a point \mathbf{x}^ , then \mathbf{x}^* is consensual and any of its identical row vectors solves problem.*

The Algorithm EXTRA

Algorithm 2 (EXTRA)

Choose $\alpha > 0$ and mixing matrices $\mathbf{W} \in \mathbb{R}^{n \times n}$ and $\tilde{\mathbf{W}} \in \mathbb{R}^{n \times n}$

Pick any $\mathbf{x}^0 \in \mathbb{R}^{n \times p}$

1. $\mathbf{x}^1 \leftarrow \mathbf{W}\mathbf{x}^0 - \alpha \nabla \mathbf{f}(\mathbf{x}^0)$;

2. for $k = 0, 1, \dots$ do

$\mathbf{x}^{k+2} \leftarrow (\mathbf{I} + \mathbf{W})\mathbf{x}^{k+1} - \tilde{\mathbf{W}}\mathbf{x}^k - \alpha[\nabla f(\mathbf{x}^{k+1}) - \nabla f(\mathbf{x}^k)]$;

end for

Assumptions on The Mixing Matrices \mathbf{W} and $\tilde{\mathbf{W}}$

Assumption 1

(Mixing matrix). Consider a connected network $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ consisting of a set of agents $\mathcal{V} = \{1, 2, \dots, n\}$ and a set of undirected edges \mathcal{E} . The mixing matrices $W = [w_{ij}] \in R^{n \times n}$ satisfy

1. (Decentralized property) If $i \neq j$ and $(i, j) \notin \mathcal{E}$, then $w_{ij} = \tilde{w}_{ij} = 0$.
2. (Symmetry) $W = W^T$, $\tilde{W} = \tilde{W}^T$.
3. (Null space property) $\text{null}\{W - \tilde{W}\} = \text{span}\{\mathbf{1}\}$, $\text{null}\{I - \tilde{W}\} \supseteq \text{span}\{\mathbf{1}\}$.
4. (Spectral property) $\tilde{W} \succ 0$ and $\frac{I+W}{2} \succcurlyeq \tilde{W} \succcurlyeq W$.

In fact, EXTRA can use the same \mathbf{W} used in DGD and simply take $\tilde{\mathbf{W}} = \frac{I+\mathbf{W}}{2}$

Mixing Matrices

The mixing matrices \mathbf{W} and $\tilde{\mathbf{W}}$ diffuse information throughout the network, which can significantly affect performance. We will verify this point in the numerical experiments!!!

\mathbf{W} can be chosen through following methods:

- symmetric doubly stochastic matrix: $\mathbf{W} = \mathbf{W}^T$, $\mathbf{W}\mathbf{1} = \mathbf{1}$ and $w_{ij} \geq 0$.
- Laplacian-based constant edge weight matrix
- Metropolis constant edge weight matrix
- symmetric fastest distributed linear averaging (FDLA) matrix

$\tilde{\mathbf{W}} = \frac{\mathbf{I} + \mathbf{W}}{2}$ is found to be very efficient.

Due to the update formula

$$\mathbf{x}^{k+2} - \mathbf{x}^{k+1} = \mathbf{W}\mathbf{x}^{k+1} - \tilde{\mathbf{W}}\mathbf{x}^k - \alpha\nabla f(\mathbf{x}^{k+1}) + \alpha\nabla f(\mathbf{x}^k),$$

We add the iterations k together and obtain

$$\mathbf{x}^{k+1} = \mathbf{W}\mathbf{x}^k - \alpha\nabla f(\mathbf{x}^k) + \sum_{t=0}^{k-1} (\mathbf{W} - \tilde{\mathbf{W}})\mathbf{x}^t, \quad k = 0, 1, \dots \quad (9)$$

An EXTRA update is a DGD update with a cumulative correction term. The role of the cumulative term $\sum_{t=0}^{k-1} (\mathbf{W} - \tilde{\mathbf{W}})\mathbf{x}^t$ is to neutralize $-\alpha\nabla f(\mathbf{x}^k)$ in $(\text{span}\{\mathbf{1}\})^\perp$

EXTRA

Preliminaries

Lemma. 1

Assume $\text{null}\{I - W\} = \text{span}\{\mathbf{1}\}$. If

$$\mathbf{x}^* = \begin{bmatrix} \text{---} & x_{(1)}^{*\top} & \text{---} \\ \text{---} & x_{(2)}^{*\top} & \text{---} \\ & \cdot & \\ & \cdot & \\ & \cdot & \\ \text{---} & x_{(n)}^{*\top} & \text{---} \end{bmatrix} \quad (10)$$

satisfies conditions:

1. $\mathbf{x}^* = W\mathbf{x}^*$ (consensus),
2. $\mathbf{1}^\top \nabla f(\mathbf{x}^*) = 0$ (optimality),

then $x^* = x_{(i)}^*$, for any i , is a solution to consensus optimization problem.

EXTRA

Lemma. 2

Given mixing matrices W and \tilde{W} , define $U = (\tilde{W} - W)^{\frac{1}{2}}$ by letting $U \triangleq VS^{\frac{1}{2}}V^{\top} \in \mathbb{R}^{n \times n}$ where $VSV^{\top} = \tilde{W} - W$ is the economical-form singular value decomposition. Then, under some assumptions, x^* is consensual and $x_{(1)}^* = x_{(2)}^* = \dots = x_{(n)}^*$ is optimal to optimization problem if and only if there exists $q^* = Up$ for some $p \in \mathbb{R}^{n \times p}$ such that

$$Uq^* + \alpha \nabla f(x^*) = 0 \quad (11)$$

$$Ux^* = 0 \quad (12)$$

EXTRA

Theorem 1

Under assumptions 1 – 3, If α satisfies $0 < \alpha < \frac{2\lambda_{\min}(\tilde{W})}{L_f}$, then

$$\|z^k - z^*\|_G^2 - \|z^{k+1} - z^*\|_G^2 \geq \zeta \|z^k - z^{k+1}\|_G^2, k = 0, 1, \dots, \quad (13)$$

where $\zeta = 1 - \frac{\alpha L_f}{2\lambda_{\min}(\tilde{W})}$.

$$z^k = \begin{bmatrix} q^k \\ x^k \end{bmatrix}, z^k = \begin{bmatrix} q^k \\ x^k \end{bmatrix}, G = \begin{bmatrix} I & 0 \\ 0 & \tilde{W} \end{bmatrix}.$$

EXTRA

The key steps in the proof of above theorem.

$$(I + W - 2\tilde{W})x^* = 0 \quad (14)$$

$$\frac{2\alpha}{L_f} \|\nabla f(x^k) - \nabla f(x^*)\|_F^2 \leq 2\alpha \langle x^k - x^*, \nabla f(x^k) - \nabla f(x^*) \rangle \quad (15)$$

$$\langle x^{k+1} - x^*, U(q^* - q^{k+1}) \rangle = \langle U(x^{k+1} - x^*), q^* - q^{k+1} \rangle = \langle Ux^{k+1}, q^* - q^{k+1} \rangle \quad (16)$$

$$\langle x^{k+1} - x^*, \tilde{W}(q^* - q^{k+1}) \rangle = \langle \tilde{W}(x^{k+1} - x^*), q^* - q^{k+1} \rangle \quad (17)$$

EXTRA

The key steps in the proof of above theorem

$$\|z^k - z^*\|_G^2 - \|z^{k+1} - z^*\|_G^2 - \|z^k - z^{k+1}\|_G^2 - 2\|x^{k+1} - x^*\|_{I+W-2\tilde{W}}^2 + \frac{\alpha L_f}{2} \|x^k - x^{k+1}\|_F^2 \geq 0 \quad (18)$$

$$\|z^k - z^*\|_G^2 - \|z^{k+1} - z^*\|_G^2 - \|z^k - z^{k+1}\|_G^2 + \frac{\alpha L_f}{2} \|x^k - x^{k+1}\|_F^2 \geq 0 \quad (19)$$

$$\|z^k - z^*\|_G^2 - \|z^{k+1} - z^*\|_G^2 \geq \|z^k - z^{k+1}\|_G^2, \quad (20)$$

$$\|z^k - z^{k+1}\|_G^2 \geq \zeta \|z^k - z^{k+1}\|_G^2 \quad (21)$$

EXTRA

From the analysis above, we can give these three necessary assumptions:

Assumption 2

(Mixing matrix). Consider a connected network $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ consisting of a set of agents $\mathcal{V} = \{1, 2, \dots, n\}$ and a set of undirected edges \mathcal{E} . The mixing matrices $W = [w_{ij}] \in \mathbb{R}^{n \times n}$ satisfy

1. (Decentralized property) If $i \neq j$ and $(i, j) \notin \mathcal{E}$, then $w_{ij} = \tilde{w}_{ij} = 0$.
2. (Symmetry) $W = W^\top$, $\tilde{W} = \tilde{W}^\top$.
3. (Null space property) $\text{null}\{W - \tilde{W}\} = \text{span}\{\mathbf{1}\}$, $\text{null}\{I - \tilde{W}\} \supseteq \text{span}\{\mathbf{1}\}$.
4. (Spectral property) $\tilde{W} \succ 0$ and $\frac{I+W}{2} \succcurlyeq \tilde{W} \succcurlyeq W$.

EXTRA

Assumption 3

(Convex objective with Lipschitz continuous gradient) Objective functions f_i are proper closed convex and Lipschitz differentiable:

$$\|\nabla f_i(x_a) - \nabla f_i(x_b)\|_2 \leq L_{f_i} \|x_a - x_b\|_2, \forall x_a, x_b \in R^P,$$

where $L_{f_i} \geq 0$ are constant.

Assumption 4

(Solution existence) The optimization problem has a nonempty set of optimal solutions: $\mathcal{X}^ \neq \emptyset$.*

EXTRA

From the theorem above, we can get:

Theorem 2

In the same setting of Theorem 3, the following rates hold:

(1) *Running-average progress:*

$$\frac{1}{k} \sum_{t=1}^k \|z^t - z^{t+1}\|_G^2 = O\left(\frac{1}{k}\right);$$

(2) *Running-best progress:*

$$\min_{t \leq k} \{\|z^t - z^{t+1}\|_G^2\} = o\left(\frac{1}{k}\right);$$

(3) *Running-average optimality residuals:*

$$\frac{1}{k} \sum_{t=1}^k \|Uq^t + \alpha \nabla f^{x^t}\|_{\tilde{W}}^2 = O\left(\frac{1}{k}\right) \text{ and } \frac{1}{k} \sum_{t=1}^k \|Ux^t\|_F^2 = O\left(\frac{1}{k}\right);$$

(4) *Running-best optimality residuals:*

$$\min_{t \leq k} \{\|Uq^t + \alpha \nabla f^{x^t}\|_{\tilde{W}}^2\} = o\left(\frac{1}{k}\right) \text{ and } \min_{t \leq k} \|Ux^t\|_F^2 = o\left(\frac{1}{k}\right);$$

EXTRA

Theorem 3

If $g(x) \triangleq f(x) + \frac{1}{4\alpha} \|x\|_{\tilde{W}-W}^2$ is restricted strongly convex with respect to x^* with constant $\mu_g > 0$, then with proper step size $\alpha < \frac{2\mu_g \lambda_{\min}(\tilde{W})}{L_f^2}$, there exists $\delta > 0$ such that the sequence $\{z^k\}$ generated by EXTRA satisfies

$$\|z^k - z^*\|_G^2 \geq (1 + \delta) \|z^{k+1} - z^*\|_G^2. \quad (22)$$

EXTRA

Denote μ_f is the constant of the restrictedly convex function $f(x)$, then in the proof of the theorem above,

$$\mu_g = \min\left\{\mu_f - 2L_f\gamma, \frac{\tilde{\lambda}_{\min}(\tilde{W} - W)}{2\alpha(1 + \frac{1}{\gamma^2})}\right\} \quad (23)$$

So we set

$$\alpha = \frac{\tilde{\lambda}_{\min}(\tilde{W} - W)}{2(1 + \frac{1}{\gamma^2})(\mu_f - 2L_f\gamma)} = O\left(\frac{\mu_g}{L_f^2}\right) \quad (24)$$

The paper said that in this case if we set $\alpha = O(\frac{1}{L_f})$, the algorithm still converges and become faster, but it remains an open question to prove linear convergence under this step.

EXTRA

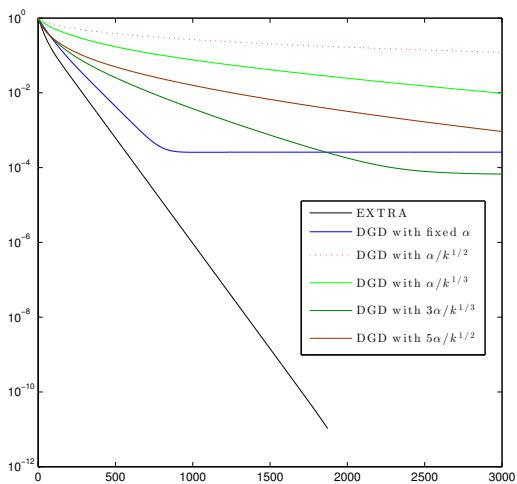
Experiment set:

We consider solving problem

$$\underset{x}{\text{minimize}} \bar{f}(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|M_{(i)}x - y_{(i)}\|_2^2 \quad (25)$$

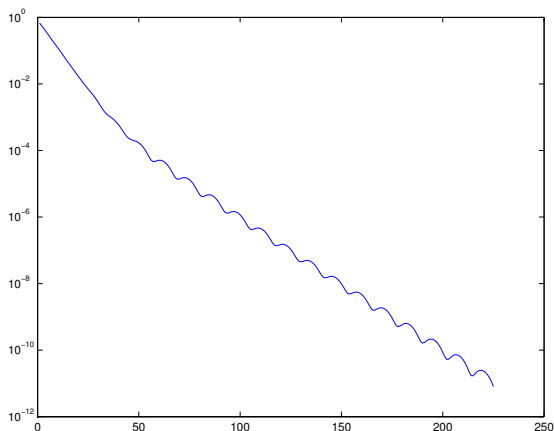
Here $y_{(i)} = M_{(i)}x + e_{(i)}$, where $y_{(i)} \in R^{m_i}$ and $M_{(i)} \in R^{m_i \times p}$ are measured data, $x \in R^p$ is unknown signal, and $e_{(i)} \in R^{m_i}$ is unknown noise. In this set, $n = 10$, $m_i = 1, \forall i$, $p = 5$. Data $M_{(i)}$ and $e_{(i)}, \forall i$, are generated following the standard normal equation. The algorithm starts from $x_{(i)}^0 = 0, \forall i$, and we set $\|x - x_i^0\| = 300$.

EXTRA



EXTRA

We adjust $M_{(i)}$, $\forall i$, to 20 to make $f_{(i)}(x)$ strongly convex. All other parameters are the same. And we set $\alpha = \frac{1}{2L_f}$.



Other Mixing Matrices

We have tried other methods to choose mixing matrices \mathbf{W} . The result of using metropolis method with different connectivity ratio $r = 0.2, 0.5, 0.7$ as follows. We can see the significant effect from mixing matrices.

