

ADMM for Distributed Optimization

Hailiang ZHAO @ ZJU-CS

<http://hliangzhao.me>

October 29, 2022

Outline

I Centralized Distributed Optimization

II Decentralized Distributed Optimization

III Asynchronous Distributed ADMM

IV Nonconvex Distributed ADMM

V ADMM with Generally Linear Constraints

References

Outline

I Centralized Distributed Optimization

II Decentralized Distributed Optimization

III Asynchronous Distributed ADMM

IV Nonconvex Distributed ADMM

V ADMM with Generally Linear Constraints

References

The Distributed Optimization Problem

Consider the following problem in a distributed environment:

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) \equiv \sum_{i=1}^m f_i(\mathbf{x}_i), \quad (1)$$

where m agents form a connected and undirected network and the local function f_i is only accessible by agent i due to storage or privacy reasons. Two kinds of networks:

- ▶ The centralized network with one centralized master agent and m worker agents. Each worker agent is connected to the master agent
- ▶ The decentralized network, which does not have the centralized agent and each agent only communicates with its neighbors

All the agents cooperate to solve (1).

Centralized Optimization

In the centralized network, we reformulate (1) as the following linearly constrained one:

$$\min_{\mathbf{x}_i, \mathbf{z}} \sum_{i=1}^m f_i(\mathbf{x}_i), \quad s.t. \quad \mathbf{x}_i = \mathbf{z}, \forall i \in [m]. \quad (2)$$

Obviously, the vanilla ADMM can solve it. Introduce the augmented Lagrangian function:

$$L_{\beta}(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}) = \sum_{i=1}^m \left(f_i(\mathbf{x}_i) + \langle \boldsymbol{\lambda}_i, \mathbf{x}_i - \mathbf{z} \rangle + \frac{\beta}{2} \|\mathbf{x}_i - \mathbf{z}\|^2 \right). \quad (3)$$

The corresponding iterations are:

$$\begin{aligned}\mathbf{z}^{k+1} &= \underset{\mathbf{z}}{\operatorname{argmin}} \sum_{i=1}^m \left(\langle \boldsymbol{\lambda}_i^k, \mathbf{x}_i^k - \mathbf{z} \rangle + \frac{\beta}{2} \|\mathbf{x}_i^k - \mathbf{z}\|^2 \right) \\ &= \frac{1}{m} \sum_{i=1}^m \left(\mathbf{x}_i^k + \frac{1}{\beta} \boldsymbol{\lambda}_i^k \right),\end{aligned}\tag{4}$$

$$\begin{aligned}\mathbf{x}_i^{k+1} &= \underset{\mathbf{x}_i}{\operatorname{argmin}} \left(f_i(\mathbf{x}_i) + \langle \boldsymbol{\lambda}_i^k, \mathbf{x}_i^k - \mathbf{z}^{k+1} \rangle + \frac{\beta}{2} \|\mathbf{x}_i^k - \mathbf{z}^{k+1}\|^2 \right) \\ &= \operatorname{prox}_{\beta^{-1}f_i}(\mathbf{z}^{k+1} - \frac{1}{\beta} \boldsymbol{\lambda}_i^k), \quad i \in [m],\end{aligned}\tag{5}$$

$$\boldsymbol{\lambda}_i^{k+1} = \boldsymbol{\lambda}_i^k + \beta(\mathbf{x}_i^{k+1} - \mathbf{z}^{k+1}), \quad i \in [m].\tag{6}$$

These steps are executed by different roles.

CADMM

The master should do — In each iteration k :

1. Wait until receiving \mathbf{x}_i^k and $\boldsymbol{\lambda}_i^k$ from all workers
2. Update \mathbf{z}^{k+1} with (4)
3. Send \mathbf{z}^{k+1} to all the workers

Each worker i should do — In each iteration k :

1. Send \mathbf{x}_i^k and $\boldsymbol{\lambda}_i^k$ to the master
2. Wait until receiving \mathbf{z}^{k+1} from the master
3. Update \mathbf{x}_i^{k+1} and $\boldsymbol{\lambda}_i^{k+1}$ with (5) and (6), respectively

We call this method Centralized ADMM for Master and Worker (*CADMM-M and CADMM-W*).

Convergence Rate of CADMM

A KKT point of (2) is denoted by $(\mathbf{x}_1^*, \dots, \mathbf{x}_m^*, \mathbf{z}^*, \boldsymbol{\lambda}_1^*, \dots, \boldsymbol{\lambda}_m^*)$. CADMM has a similar convergence to the vanilla ADMM.

Theorem 1.1 (*similar to Theorem 1.3 of Slide 2*)

Suppose that $i \in [m]$ [$f_i(\mathbf{x}_i)$ is convex]. Then for CADMM, we have

$$\left| \sum_i f_i(\hat{\mathbf{x}}_i^{K+1}) - \sum_i f_i(\mathbf{x}_i^*) \right| \leq \frac{C}{2(K+1)} + \frac{2\sqrt{C}\sqrt{\sum_i \|\boldsymbol{\lambda}_i^*\|^2}}{\sqrt{\beta}(K+1)} \quad (7)$$

$$\sqrt{\sum_i \|\hat{\mathbf{x}}_i^{K+1} - \hat{\mathbf{z}}^{K+1}\|^2} \leq \frac{2\sqrt{C}}{\sqrt{\beta}(K+1)}, \quad (8)$$

where $C = \frac{1}{\beta} \sum_i \|\boldsymbol{\lambda}_i^0 - \boldsymbol{\lambda}_i^*\|^2 + \beta \sum_i \|\mathbf{x}_i^0 - \mathbf{x}_i^*\|^2$.

Convergence Rate of CADMM

We also have:

Theorem 1.2 (similar to Theorem 1.4 of Slide 2)

Suppose that *each f_i is μ -strongly convex and L -smooth*. Then for CADMM we have

$$\begin{aligned} & \sum_i \left(\frac{1}{2\beta} \|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*\|^2 + \frac{\beta}{2} \|\mathbf{x}_i^{k+1} - \mathbf{x}_i^*\|^2 \right) \\ & \leq \left(1 + \frac{1}{2} \sqrt{\frac{\mu}{L}} \right)^{-1} \sum_i \left(\frac{1}{2\beta} \|\boldsymbol{\lambda}_i^k - \boldsymbol{\lambda}_i^*\|^2 + \frac{\beta}{2} \|\mathbf{x}_i^k - \mathbf{x}_i^*\|^2 \right). \quad (9) \end{aligned}$$

CLADMM

LADMM can also be used to solve (2). We can linearize f_i in (5) simplify the computation, if the proximal mapping of f_i is not easily computable. The iterations are:

$$\begin{aligned}\mathbf{z}^{k+1} &= \underset{\mathbf{z}}{\operatorname{argmin}} \sum_{i=1}^m \left(\langle \boldsymbol{\lambda}_i^k, \mathbf{x}_i^k - \mathbf{z} \rangle + \frac{\beta}{2} \|\mathbf{x}_i^k - \mathbf{z}\|^2 \right) \\ &= \frac{1}{m} \sum_{i=1}^m \left(\mathbf{x}_i^k + \frac{1}{\beta} \boldsymbol{\lambda}_i^k \right), \quad (\text{unchanged})\end{aligned} \tag{10}$$

CLADMM

The iterations are (cont'd):

$$\begin{aligned}\mathbf{x}_i^{k+1} &= \operatorname{argmin}_{\mathbf{x}_i} \left(f_i(\mathbf{x}_i) + \langle \boldsymbol{\lambda}_i^k, \mathbf{x}_i^k - \mathbf{z}^{k+1} \rangle \right. \\ &\quad \left. + \frac{\beta}{2} \|\mathbf{x}_i^k - \mathbf{z}^{k+1}\|^2 + D_{\Psi_i}(\mathbf{x}_i, \mathbf{x}_i^k) \right) \\ &= \operatorname{argmin}_{\mathbf{x}_i} \left(\langle \nabla f_i(\mathbf{x}_i^k), \mathbf{x}_i - \mathbf{x}_i^k \rangle + \frac{L}{2} \|\mathbf{x}_i - \mathbf{x}_i^k\|^2 \right. \\ &\quad \left. + \langle \boldsymbol{\lambda}_i^k, \mathbf{x}_i - \mathbf{z}^{k+1} \rangle + \frac{\beta}{2} \|\mathbf{x}_i - \mathbf{z}^{k+1}\|^2 \right) \\ &= \frac{1}{L + \beta} \left(L\mathbf{x}_i^k + \beta\mathbf{z}^{k+1} - \nabla f_i(\mathbf{x}_i^k) - \boldsymbol{\lambda}_i^k \right),\end{aligned}\tag{11}$$

$$\boldsymbol{\lambda}_i^{k+1} = \boldsymbol{\lambda}_i^k + \beta(\mathbf{x}_i^{k+1} - \mathbf{z}^{k+1}), \quad i \in [m], \quad (\text{unchanged})\tag{12}$$

by choosing $\Psi_i(\mathbf{x}_i) = \frac{L}{2} \|\mathbf{x}_i\|^2 - f_i(\mathbf{x}_i)$.

CLADMM

The master should do — In each iteration k :

1. Wait until receiving \mathbf{x}_i^k and $\boldsymbol{\lambda}_i^k$ from all workers
2. Update \mathbf{z}^{k+1} with (10)
3. Send \mathbf{z}^{k+1} to all the workers

Each worker i should do — In each iteration k :

1. Send \mathbf{x}_i^k and $\boldsymbol{\lambda}_i^k$ to the master
2. Wait until receiving \mathbf{z}^{k+1} from the master
3. Update \mathbf{x}_i^{k+1} and $\boldsymbol{\lambda}_i^{k+1}$ with (11) and (12), respectively

We call this method Centralized Linearized ADMM for Master and Worker (*CLADMM-M and CLADMM-W*).

Convergence Rate of CLADMM

Of course CLADMM can achieve sublinear convergence rate (Theorem II.1 of *Slide 2*). Further, based on Theorem 3.8 of the ADMM book and using $L_\Psi \leq L - \mu$, where $\Psi(\mathbf{x}) = \sum_i \Psi_i(\mathbf{x}_i)$, CLADMM can achieve linear convergence rate.

Theorem I.3

Suppose that each $f_i(\mathbf{x}_i)$ is μ -strongly convex and L -smooth. Let $\beta = \sqrt{\mu}(2L - \mu)$. Then for CLADMM, we have

$$\begin{aligned} & \sum_i \left(\frac{1}{2\beta} \|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*\|^2 + \frac{\beta}{2} \|\mathbf{x}_i^{k+1} - \mathbf{x}_i^*\|^2 + D_{\Psi_i}(\mathbf{x}_i^*, \mathbf{x}_i^{k+1}) \right) \\ & \leq \left[1 + \frac{1}{3} \min \left(\sqrt{\frac{\mu}{2L - \mu}}, \frac{\mu}{L - \mu} \right) \right]^{-1} \\ & \quad \times \sum_i \left(\frac{1}{2\beta} \|\boldsymbol{\lambda}_i^k - \boldsymbol{\lambda}_i^*\|^2 + \frac{\beta}{2} \|\mathbf{x}_i^k - \mathbf{x}_i^*\|^2 + D_{\Psi_i}(\mathbf{x}_i^*, \mathbf{x}_i^k) \right). \end{aligned} \tag{13}$$

Acc-CLADMM

Accelerated Linearized ADMM can also be used to solve (2). Based on Acc-LADMM-3, we have the following iterations:

$$\mathbf{w}_i^k = \theta \mathbf{x}_i^k + (1 - \theta) \tilde{\mathbf{x}}_i^k, \quad (14)$$

$$\begin{aligned} \mathbf{z}^{k+1} &= \underset{\mathbf{z}}{\operatorname{argmin}} \sum_i \left(\langle \boldsymbol{\lambda}_i^k, \mathbf{x}_i^k - \mathbf{z} \rangle + \frac{\beta\theta}{2} \|\mathbf{x}_i^k - \mathbf{z}\|^2 \right) \\ &= \frac{1}{m} \sum_i \left(\mathbf{x}_i^k + \frac{1}{\beta\theta} \boldsymbol{\lambda}_i^k \right), \end{aligned} \quad (15)$$

$$\begin{aligned} \mathbf{x}_i^{k+1} &= \frac{1}{\frac{\theta}{\alpha} + \mu} \left\{ \mu \mathbf{w}_i^k + \frac{\theta}{\alpha} \mathbf{x}_i^k - [\nabla f_i(\mathbf{x}_i^k) + \boldsymbol{\lambda}_i^k \right. \\ &\quad \left. + \beta\theta(\mathbf{x}_i^k - \mathbf{z}^{k+1})] \right\}, \end{aligned} \quad (16)$$

Acc-CLADMM

Based on Acc-LADMM-3, we have the following iterations (cont'd):

$$\tilde{\mathbf{z}}^{k+1} = \theta \mathbf{z}^{k+1} + (1 - \theta) \tilde{\mathbf{z}}^k, \quad (17)$$

$$\tilde{\mathbf{x}}_i^{k+1} = \theta \mathbf{x}_i^{k+1} + (1 - \theta) \tilde{\mathbf{x}}_i^k, \quad (18)$$

$$\boldsymbol{\lambda}_i^{k+1} = \boldsymbol{\lambda}_i^k + \beta \theta (\mathbf{x}_i^{k+1} - \mathbf{z}^{k+1}). \quad (19)$$

Acc-CLADMM

The master should do — In each iteration k :

1. Wait until receiving \mathbf{x}_i^k and $\boldsymbol{\lambda}_i^k$ from all workers
2. Update \mathbf{z}^{k+1} and $\tilde{\mathbf{z}}^{k+1}$ with (15) and (17), respectively

Each worker i should do — In each iteration k :

1. Send \mathbf{x}_i^k and $\boldsymbol{\lambda}_i^k$ to the master
2. Wait until receiving \mathbf{z}^{k+1} from the master
3. Update \mathbf{x}_i^{k+1} and $\boldsymbol{\lambda}_i^{k+1}$ with (5) and (6), respectively

We call this method Accelerated Centralized Linearized ADMM for Master and Worker (*Acc-CLADMM-M and Acc-CLADMM-W*).

Convergence Rate of Acc-CLADMM

Similar to Theorem III.1 of *ADMM Slide 2*, we have the following convergence rate for Acc-CLADMM.

Theorem I.4

Suppose that each $f_i(\mathbf{x}_i)$ is μ -strongly convex and L -smooth, $i \in [m]$. Let $\alpha = \frac{1}{4L}$, $\beta = L$, and $\theta = \sqrt{\frac{\mu}{L}}$. Then for Acc-CLADMM, we have

$$l_{k+1} \leq \left(1 - \sqrt{\frac{\mu}{L}}\right) l_k, \quad (20)$$

where

$$\begin{aligned} l_k = & (1 - \theta) \sum_i \left(f_i(\tilde{\mathbf{x}}_i^k) - f_i(\mathbf{x}_i^*) + \langle \boldsymbol{\lambda}_i^*, \tilde{\mathbf{x}}_i^k - \tilde{\mathbf{z}}^k \rangle \right) \\ & + \frac{\theta^2}{2\alpha} \sum_i \|\mathbf{x}_i^k - \mathbf{x}_i^*\|^2 + \frac{1}{2\beta} \|\boldsymbol{\lambda}_i^k - \boldsymbol{\lambda}^*\|^2. \end{aligned} \quad (21)$$

Complexity Comparisons

We can find that Acc-CLADMM is faster than CLAMM with a better dependence on the condition number L/μ . CADMM has the same convergence rate as the Acc-CLADMM. However, CADMM may need to solve subproblem iteratively at each iteration, while Acc-CLADMM only performs a gradient descent update.

CADMM	CLADMM	Acc-CLADMM
$\mathcal{O}(\sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon})$	$\mathcal{O}(\frac{L}{\mu} \log \frac{1}{\epsilon})$	$\mathcal{O}(\sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon})$

Outline

I Centralized Distributed Optimization

II Decentralized Distributed Optimization

III Asynchronous Distributed ADMM

IV Nonconvex Distributed ADMM

V ADMM with Generally Linear Constraints

References

Decentralized Optimization

Now we consider the decentralized topology. In this case, we cannot introduce the variable \mathbf{z} since there is no central node any more.

Denote \mathcal{E} as the set of edges. Assume that all the nodes are ordered from 1 to m . For any two nodes i and j , if i and j are directly connected in the network **and** $i < j$, we say $(i, j) \in \mathcal{E}$.

To simplify the presentation, we order the edges from 1 to $|\mathcal{E}|$. For each node i , we denote \mathcal{N}_i as its neighborhood:

$$\mathcal{N}_i = \{j \mid (i, j) \in \mathcal{E} \text{ or } (j, i) \in \mathcal{E}\}, \quad (22)$$

and $d_i = |\mathcal{N}_i|$ as its degree.

Decentralized Optimization

Introduce auxiliary variables \mathbf{z}_{ij} if $(i, j) \in \mathcal{E}$. Then we can reformulate (1) as follows:

$$\min_{\mathbf{x}_i, \mathbf{z}_{ij}} \sum_{i=1}^m f_i(\mathbf{x}_i) \quad s.t. \quad \mathbf{x}_i = \mathbf{z}_{ij}, \mathbf{x}_j = \mathbf{z}_{ij}, \forall (i, j) \in \mathcal{E}. \quad (23)$$

That is to say, each variable \mathbf{x}_i corresponds to one node, while each variable \mathbf{z}_{ij} ($i < j$) corresponds to one edge. The augmented Lagrangian function is

$$\begin{aligned} L_{\beta}(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}) = & \sum_i f_i(\mathbf{x}_i) + \sum_{(i,j) \in \mathcal{E}} \left(\langle \boldsymbol{\lambda}_{ij}, \mathbf{x}_i - \mathbf{z}_{ij} \rangle + \langle \boldsymbol{\gamma}_{ij}, \mathbf{x}_j - \mathbf{z}_{ij} \rangle \right. \\ & \left. + \frac{\beta}{2} \|\mathbf{x}_i - \mathbf{z}_{ij}\|^2 + \frac{\beta}{2} \|\mathbf{x}_j - \mathbf{z}_{ij}\|^2 \right). \end{aligned} \quad (24)$$

DADMM

We can use the vanilla ADMM to solve it. The iterations are:

$$\begin{aligned}\mathbf{x}_i^{k+1} = \operatorname{argmin}_{\mathbf{x}_i} & \left[f_i(\mathbf{x}_i) + \sum_{j:(i,j) \in \mathcal{E}} \left(\langle \boldsymbol{\lambda}_{ij}^k, \mathbf{x}_i - \mathbf{z}_{ij}^k \rangle \right. \right. \\ & \left. \left. + \frac{\beta}{2} \|\mathbf{x}_i - \mathbf{z}_{ij}^k\|^2 \right) + \sum_{j:(i,j) \in \mathcal{E}} \left(\langle \boldsymbol{\gamma}_{ji}^k, \mathbf{x}_i - \mathbf{z}_{ji}^k \rangle + \frac{\beta}{2} \|\mathbf{x}_i - \mathbf{z}_{ji}^k\|^2 \right) \right],\end{aligned}\tag{25}$$

$$\begin{aligned}\mathbf{z}_{ij}^{k+1} = \operatorname{argmin}_{\mathbf{z}_{ij}} & \left(- \langle \boldsymbol{\lambda}_{ij}^k + \boldsymbol{\gamma}_{ij}^k, \mathbf{z}_{ij} \rangle + \frac{\beta}{2} \|\mathbf{x}_i^{k+1} - \mathbf{z}_{ij}\|^2 \right. \\ & \left. + \frac{\beta}{2} \|\mathbf{x}_j^{k+1} - \mathbf{z}_{ij}\|^2 \right) \\ & = \frac{1}{2\beta} (\boldsymbol{\lambda}_{ij}^k + \boldsymbol{\gamma}_{ij}^k) + \frac{1}{2} (\mathbf{x}_i^{k+1} + \mathbf{x}_j^{k+1}).\end{aligned}\tag{26}$$

DADMM

We can use the vanilla ADMM to solve it. The iterations are (cont'd):

$$\boldsymbol{\lambda}_{ij}^{k+1} = \boldsymbol{\lambda}_{ij}^k + \beta(\mathbf{x}_i^{k+1} - \mathbf{z}_{ij}^{k+1}), \quad (27)$$

$$\boldsymbol{\gamma}_{ij}^{k+1} = \boldsymbol{\gamma}_{ij}^k + \beta(\mathbf{x}_j^{k+1} - \mathbf{z}_{ij}^{k+1}). \quad (28)$$

Actually (25) \sim (28) can be simplified. Summing (27) and (28) and using (26), we have

$$\boldsymbol{\lambda}_{ij}^{k+1} + \boldsymbol{\gamma}_{ij}^{k+1} = \mathbf{0}, \forall k \geq 0. \quad (29)$$

Initialize $\boldsymbol{\lambda}_{ij}^0 = \boldsymbol{\gamma}_{ij}^0 = \mathbf{0}$ we have

$$\boldsymbol{\lambda}_{ij}^k + \boldsymbol{\gamma}_{ij}^k = \mathbf{0}, \forall k \geq 0. \quad (30)$$

DADMM

Plugging (30) into (26) we have

$$\mathbf{z}_{ij}^{k+1} = \frac{1}{2}(\mathbf{x}_i^{k+1} + \mathbf{x}_j^{k+1}), \forall k \geq 0. \quad (31)$$

Similarly, we may initialize $\mathbf{z}_{ij}^0 = \frac{1}{2}(\mathbf{x}_i^0 + \mathbf{x}_j^0)$. Combing (31) and (27) we have

$$\boldsymbol{\lambda}_{ij}^{k+1} = \boldsymbol{\lambda}_{ij}^k + \frac{\beta}{2}(\mathbf{x}_i^{k+1} - \mathbf{x}_j^{k+1}) \quad (32)$$

Thus we have

$$\boldsymbol{\lambda}_{ij}^{k+1} = \beta \sum_{t=1}^{k+1} \frac{1}{2}(\mathbf{x}_i^t - \mathbf{x}_j^t). \quad (33)$$

Similarly we have

$$\boldsymbol{\gamma}_{ij}^{k+1} = \beta \sum_{t=1}^{k+1} \frac{1}{2}(\mathbf{x}_j^t - \mathbf{x}_i^t). \quad (34)$$

Note that we only define λ_{ij} , γ_{ij} , and z_{ij} for $i < j$. Now we define

$$\lambda_{ij} \equiv \gamma_{ji} \text{ and } z_{ij} \equiv z_{ji} \text{ for } i > j. \quad (35)$$

Then (31), (32) and (33) hold for both $i < j$ and $i > j$.

DADMM

Thus, we can simplify (25) to

$$\begin{aligned}
\mathbf{x}_i^{k+1} &= \underset{\mathbf{x}_i}{\operatorname{argmin}} \left[f_i(\mathbf{x}_i) + \sum_{j:(i,j) \in \mathcal{E}} \left(\langle \boldsymbol{\lambda}_{ij}^k - \beta \mathbf{z}_{ij}^k, \mathbf{x}_i \rangle + \frac{\beta}{2} \|\mathbf{x}_i\|^2 \right) \right. \\
&\quad \left. + \sum_{j:(j,i) \in \mathcal{E}} \left(\langle \boldsymbol{\gamma}_{ji}^k - \beta \mathbf{z}_{ji}^k, \mathbf{x}_i \rangle + \frac{\beta}{2} \|\mathbf{x}_i\|^2 \right) \right] \\
&= \underset{\mathbf{x}_i}{\operatorname{argmin}} \left[f_i(\mathbf{x}_i) + \sum_{j \in \mathcal{N}_i} \left(\langle \boldsymbol{\lambda}_{ij}^k - \beta \mathbf{z}_{ij}^k, \mathbf{x}_i \rangle + \frac{\beta}{2} \|\mathbf{x}_i\|^2 \right) \right] \\
&= \underset{\mathbf{x}_i}{\operatorname{argmin}} \left[f_i(\mathbf{x}_i) + \sum_{j \in \mathcal{N}_i} \left(\langle \boldsymbol{\lambda}_{ij}^k - \beta \mathbf{z}_{ij}^k + \beta \mathbf{x}_i^k, \mathbf{x}_i \rangle + \frac{\beta}{2} \|\mathbf{x}_i - \mathbf{x}_i^k\|^2 \right) \right] \\
&= \underset{\mathbf{x}_i}{\operatorname{argmin}} \left[f_i(\mathbf{x}_i) + \sum_{j \in \mathcal{N}_i} \left(\langle \boldsymbol{\lambda}_{ij}^k + \frac{\beta}{2} (\mathbf{x}_i^k - \mathbf{x}_j^k), \mathbf{x}_i \rangle + \frac{\beta}{2} \|\mathbf{x}_i - \mathbf{x}_i^k\|^2 \right) \right].
\end{aligned} \tag{36}$$

DADMM

Denote $\mathbf{L} \in \mathbf{R}^{m \times m}$ as the Laplacian matrix and \mathbf{D} as the diagonal degree matrix with $\mathbf{D}_{ii} = d_i$. Note that \mathbf{L} is symmetric and satisfies $\mathbf{0} \preceq \mathbf{L} \preceq 2\mathbf{D}$. Define

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_m^T \end{pmatrix} \in \mathbf{R}^{m \times d}, f(\mathbf{X}) = \sum_{i=1}^m f_i(\mathbf{x}_i), \quad (37)$$

$$\mathbf{v}_i = \sum_{j \in \mathcal{N}_i} \lambda_{ij}, \text{ and } \mathbf{\Theta} = \begin{pmatrix} \mathbf{v}_1^T \\ \vdots \\ \mathbf{v}_m^T \end{pmatrix} \in \mathbf{R}^{m \times d}. \quad (38)$$

Then we have

$$\mathbf{L}_i^T \mathbf{X} = d_i \mathbf{x}_i^T - \sum_{j \in \mathcal{N}_i} \mathbf{x}_j^T, \quad (39)$$

where \mathbf{L}_i is the i -th column of \mathbf{L} .

DADMM

With \mathbf{L} and \mathbf{v}_i , (36) can be written as

$$\begin{aligned}\mathbf{x}_i^{k+1} &= \underset{\mathbf{x}_i}{\operatorname{argmin}} \left[f_i(\mathbf{x}_i) + \langle \mathbf{v}_i^k, \mathbf{x}_i \rangle \right. \\ &\quad \left. + \frac{\beta}{2} \langle \sum_{j \in \mathcal{N}_i} \mathbf{L}_{ij} \mathbf{x}_j^k, \mathbf{x}_i \rangle + \frac{\beta d_i}{2} \|\mathbf{x}_i - \mathbf{x}_i^k\|^2 \right] \\ &= \operatorname{prox}_{(\beta d_i)^{-1} f_i} \left(\mathbf{x}_i^k - \frac{1}{\beta d_i} \left(\mathbf{v}_i^k + \frac{\beta}{2} \sum_{j \in \mathcal{N}_i} \mathbf{L}_{ij} \mathbf{x}_j^k \right) \right), \quad \forall i \in [m].\end{aligned}\tag{40}$$

Summing (32) over $j \in \mathcal{N}_i$, we have

$$\mathbf{v}_i^{k+1} = \mathbf{v}_i^k + \frac{\beta}{2} \sum_{j \in \mathcal{N}_i} \mathbf{L}_{ij} \mathbf{x}_j^{k+1}, \forall i \in [m].\tag{41}$$

DADMM

Each node i should do — In each iteration k :

1. Update \mathbf{x}_i^{k+1} by (40)
2. Send \mathbf{x}_i^{k+1} to its neighbors
3. Wait until receiving \mathbf{x}_j^{k+1} from all its neighbors, $j \in \mathcal{N}_i$
4. Update \mathbf{v}_i^{k+1} by (41)

We call this method Decentralized ADMM (*DADMM*).

(40) and (41) can be written in a compact form:

$$\begin{aligned} \mathbf{X}^{k+1} = \operatorname{argmin}_{\mathbf{X}} & \left(f(\mathbf{X}) + \langle \boldsymbol{\Theta}^k + \frac{\beta}{2} \mathbf{L} \mathbf{X}^k, \mathbf{X} \rangle \right. \\ & \left. + \frac{\beta}{2} \|\sqrt{\mathbf{D}}(\mathbf{X} - \mathbf{X}^k)\|^2 \right), \end{aligned} \quad (42)$$

$$\boldsymbol{\Theta}^{k+1} = \boldsymbol{\Theta}^k + \frac{\beta}{2} \mathbf{L} \mathbf{X}^{k+1} = \boldsymbol{\Theta}^k + \beta \mathbf{W}^2 \mathbf{X}^{k+1}, \quad (43)$$

where $\mathbf{W} = \sqrt{\mathbf{L}/2}$. Let $\boldsymbol{\Theta}^0 \in \operatorname{span}(\mathbf{W}^2)$, we know that

$$\boldsymbol{\Theta}^k \in \operatorname{span}(\mathbf{W}^2), \forall k \geq 0, \quad (44)$$

and *there exists* $\boldsymbol{\Omega}^k$ *such that* $\boldsymbol{\Theta}^k = \mathbf{W} \boldsymbol{\Omega}^k$.

(42) and (43) can be further written as

$$\begin{aligned}
 \mathbf{X}^{k+1} &= \underset{\mathbf{X}}{\operatorname{argmin}} \left(f(\mathbf{X}) + \langle \boldsymbol{\Omega}^k, \mathbf{W}\mathbf{X} \rangle \right. \\
 &\quad \left. + \beta \langle \mathbf{W}^2 \mathbf{X}^k, \mathbf{X} \rangle + \frac{\beta}{2} \|\sqrt{\mathbf{D}}(\mathbf{X} - \mathbf{X}^k)\|^2 \right. \\
 &= \underset{\mathbf{X}}{\operatorname{argmin}} \left(f(\mathbf{X}) + \langle \boldsymbol{\Omega}^k, \mathbf{W}\mathbf{X} \rangle + \frac{\beta}{2} \|\mathbf{W}\mathbf{X}\|^2 + D_{\Psi}(\mathbf{X}, \mathbf{X}^k) \right),
 \end{aligned} \tag{45}$$

$$\boldsymbol{\Omega}^{k+1} = \boldsymbol{\Omega}^k + \beta \mathbf{W}\mathbf{X}^{k+1}, \tag{46}$$

where

$$\Psi(\mathbf{X}) = \frac{\beta}{2} \|\sqrt{\mathbf{D}}\mathbf{X}\|^2 - \frac{\beta}{2} \|\mathbf{W}\mathbf{X}\|^2. \tag{47}$$

DADMM'

We can find that (45) and (46) is equal to using the LADMM to solve

$$\min_{\mathbf{X}} f(\mathbf{X}), \quad s.t. \quad \mathbf{W}\mathbf{X} = \mathbf{0}. \quad (48)$$

We call (45) and (46) with general Ψ settings *DADMM'*.
DADMM' is not implementable in the distributed manner because $\mathbf{W} = \sqrt{\mathbf{L}/2}$, which requires the overall graph structure information.

Convergence Rate of DADMM'

Theorem II.1

Assume that *each* f_i is μ -strongly convex and L -smooth, $i \in [m]$, and $\Psi(\mathbf{y})$ is convex and L_Ψ -smooth. Initialize $\mathbf{\Omega}^0 = \mathbf{0}$. The for DADMM', we have

$$\begin{aligned} & \frac{1}{2\beta} \|\mathbf{\Omega}^{k+1} - \mathbf{\Omega}^*\|^2 + \frac{\beta}{2} \|\mathbf{W}\mathbf{X}^{k+1} - \mathbf{W}\mathbf{X}^*\|^2 + D_\Psi(\mathbf{X}^*, \mathbf{X}^{k+1}) \\ & \leq \left(1 + \frac{1}{3} \min \left\{ \frac{\beta\sigma_{\mathbf{L}}}{2(L + L_\Psi)}, \frac{\mu}{\beta\|\mathbf{W}\|^2}, \frac{\mu}{L_\Psi} \right\} \right)^{-1} \\ & \times \left(\frac{1}{2\beta} \|\mathbf{\Omega}^k - \mathbf{\Omega}^*\|^2 + \frac{\beta}{2} \|\mathbf{W}\mathbf{X}^k - \mathbf{W}\mathbf{X}^*\|^2 + D_\Psi(\mathbf{X}^*, \mathbf{X}^k) \right), \end{aligned} \tag{49}$$

where $\sigma_{\mathbf{L}}$ is the smallest positive eigenvalue of \mathbf{L} .

Convergence Rate of DADMM

For DLADMM', when Ψ is set as (47) and $L_\Psi = \beta d_{\max}$, where $d_{\max} = \max_i \{d_i\}$. DLADMM' reduces to DLADMM, and

$$\|\mathbf{W}\|^2 \leq \frac{1}{2}\|\mathbf{L}\| \leq \|\mathbf{D}\| \leq d_{\max}. \quad (50)$$

We thus have the following convergence result.

Theorem II.2

Assume that *each f_i is μ -strongly convex and L -smooth, $i \in [m]$.*

Initialize $\mathbf{\Omega}^0 = \mathbf{0}$ and let $\beta = \mathcal{O}(\sqrt{\frac{\mu L}{\sigma_L d_{\max}}})$. Then DADMM

needs $\mathcal{O}((\sqrt{\frac{L d_{\max}}{\mu \sigma_L}} + \frac{d_{\max}}{\sigma_L}) \log \frac{1}{\epsilon})$ iterations to find an ϵ -approximate solution $(\mathbf{X}, \mathbf{\Omega})$, i.e.,

$$\frac{1}{2\beta}\|\mathbf{\Omega} - \mathbf{\Omega}^*\|^2 + \frac{\beta}{2}\|\mathbf{W}\mathbf{X} - \mathbf{W}\mathbf{X}^*\|^2 + D_\Psi(\mathbf{X}^*, \mathbf{X}) \leq \epsilon. \quad (51)$$

DLADMM

LADMM can also be used to solve (23). We may linearize f_i :

$$\begin{aligned}
\mathbf{x}_i^{k+1} &= \operatorname{argmin}_{\mathbf{x}_i} \left[\langle \nabla f_i(\mathbf{x}_i^k), \mathbf{x}_i - \mathbf{x}_i^k \rangle + \frac{L}{2} \|\mathbf{x}_i - \mathbf{x}_i^k\|^2 \right. \\
&\quad + \sum_{j:(i,j) \in \mathcal{E}} \left(\langle \boldsymbol{\lambda}_{ij}^k, \mathbf{x}_i - \mathbf{z}_{ij}^k \rangle + \frac{\beta}{2} \|\mathbf{x}_i - \mathbf{z}_{ij}^k\|^2 \right) \\
&\quad \left. + \sum_{j:(j,i) \in \mathcal{E}} \left(\langle \boldsymbol{\gamma}_{ji}^k, \mathbf{x}_i - \mathbf{z}_{ji}^k \rangle + \frac{\beta}{2} \|\mathbf{x}_i - \mathbf{z}_{ji}^k\|^2 \right) \right] \\
&= \operatorname{argmin}_{\mathbf{x}_i} \left[\langle \nabla f_i(\mathbf{x}_i^k), \mathbf{x}_i - \mathbf{x}_i^k \rangle + \frac{L}{2} \|\mathbf{x}_i - \mathbf{x}_i^k\|^2 \right. \\
&\quad \left. + \sum_{j \in \mathcal{N}_i} \left(\langle \boldsymbol{\lambda}_{ij}^k + \frac{\beta}{2} (\mathbf{x}_i^k - \mathbf{x}_j^k), \mathbf{x}_i \rangle + \frac{\beta}{2} \|\mathbf{x}_i - \mathbf{x}_i^k\|^2 \right) \right] \\
&= \mathbf{x}_i^k - \frac{1}{L + \beta d_i} \left\{ \nabla f_i(\mathbf{x}_i^k) + \sum_{j \in \mathcal{N}_i} \left[\boldsymbol{\lambda}_{ij}^k + \frac{\beta}{2} (\mathbf{x}_i^k - \mathbf{x}_j^k) \right] \right\}. \quad (52)
\end{aligned}$$

DLADMM

Similar to (45) and (46), we can write these iterations in a compact form:

$$\begin{aligned}\mathbf{X}^{k+1} &= \underset{\mathbf{X}}{\operatorname{argmin}} \left(\langle \nabla f(\mathbf{X}^k), \mathbf{X} \rangle + \frac{L}{2} \|\mathbf{X} - \mathbf{X}^k\|^2 \right. \\ &\quad \left. + \langle \boldsymbol{\Omega}^k, \mathbf{W}\mathbf{X} \rangle + \beta \langle \mathbf{W}^2 \mathbf{W}^k, \mathbf{X} \rangle + \frac{\beta}{2} \|\sqrt{\mathbf{D}}(\mathbf{X} - \mathbf{X}^k)\|^2 \right) \\ &= \mathbf{X}^k - (L\mathbf{I} + \beta\mathbf{D})^{-1} (\beta\mathbf{W}^2 \mathbf{X}^k + \nabla f(\mathbf{X}^k) + \mathbf{W}\boldsymbol{\Omega}^k), \quad (53)\end{aligned}$$

$$\boldsymbol{\Omega}^{k+1} = \boldsymbol{\Omega}^k + \beta\mathbf{W}\mathbf{X}^{k+1}, \quad (54)$$

which is also a special case of DADMM' with

$$\Phi(\mathbf{X}) = \frac{L}{2} \|\mathbf{X}\|^2 - f(\mathbf{X}) + \frac{\beta}{2} \|\sqrt{\mathbf{D}}\mathbf{X}\|^2 - \frac{\beta}{2} \|\mathbf{W}\mathbf{X}\|^2 \quad (55)$$

and

$$L_{\Phi} = L + \beta d_{\max}. \quad (56)$$

DLADMM

Each node i should do — In each iteration k :

1. Update \mathbf{x}_i^{k+1} by (52)
2. Send \mathbf{x}_i^{k+1} to its neighbors
3. Wait until receiving \mathbf{x}_j^{k+1} from all its neighbors, $j \in \mathcal{N}_i$
4. Update \mathbf{v}_i^{k+1} by (41)

We call (52), (26) \sim (28) *DLADMM*.

Convergence Rate of DLADMM

We have the following convergence result for DLADMM.

Theorem II.3

Assume that *each f_i is μ -strongly convex and L -smooth, $i \in [m]$.*

Initialize $\mathbf{\Omega}^0 = \mathbf{0}$ and let $\beta = \mathcal{O}(\sqrt{\frac{\mu L}{\sigma_L d_{\max}}})$. Then DLADMM needs $\mathcal{O}((\frac{L}{\mu} + \frac{d_{\max}}{\sigma_L}) \log \frac{1}{\epsilon})$ iterations to find an ϵ -approximate solution $(\mathbf{X}, \mathbf{\Omega})$, i.e.,

$$\frac{1}{2\beta} \|\mathbf{\Omega} - \mathbf{\Omega}^*\|^2 + \frac{\beta}{2} \|\mathbf{W}\mathbf{X} - \mathbf{W}\mathbf{X}^*\|^2 + D_{\Psi}(\mathbf{X}^*, \mathbf{X}) \leq \epsilon. \quad (57)$$

Acc-DLADMM

Accelerated Linearized ADMM can also be used to solve (23). Based on Acc-LADMM-3, we have the following iterations:

$$\mathbf{Y}^k = \theta \mathbf{X}^k + (1 - \theta) \tilde{\mathbf{X}}^k, \quad (58)$$

$$\mathbf{X}^{k+1} = \frac{1}{\frac{\theta}{\alpha} + \mu} \left[\mu \mathbf{Y}^k + \frac{\theta}{\alpha} \mathbf{X}^k - (\nabla f(\mathbf{Y}^k) + \mathbf{W} \boldsymbol{\Omega}^k + \beta \theta \mathbf{W}^2 \mathbf{X}^k) \right], \quad (59)$$

$$\tilde{\mathbf{X}}^{k+1} = \theta \mathbf{X}^{k+1} + (1 - \theta) \tilde{\mathbf{X}}^k, \quad (60)$$

$$\boldsymbol{\Omega}^{k+1} = \boldsymbol{\Omega}^k + \beta \theta \mathbf{W} \mathbf{X}^{k+1}. \quad (61)$$

The above formulas are written in compact form.

Acc-DLADMM

We call it *Acc-DLADMM* — in the distributed form:

$$\mathbf{y}_i^k = \theta \mathbf{x}_i^k + (1 - \theta) \tilde{\mathbf{x}}_i^k, \quad (62)$$

$$\mathbf{x}_i^{k+1} = \frac{1}{\frac{\theta}{\alpha} + \mu} \left[\mu \mathbf{y}_i^k + \frac{\theta}{\alpha} \mathbf{x}_i^k - (\nabla f_i(\mathbf{y}_i^k) + \mathbf{v}_i^k + \frac{\beta\theta}{2} \sum_{j \in \mathcal{N}_i} \mathbf{L}_{ij} \mathbf{x}_j^k) \right], \quad (63)$$

$$\tilde{\mathbf{x}}_i^{k+1} = \theta \mathbf{x}_i^{k+1} + (1 - \theta) \tilde{\mathbf{x}}_i^k, \quad (64)$$

$$\mathbf{v}_i^{k+1} = \mathbf{v}_i^k + \frac{\beta\theta}{2} \sum_{j \in \mathcal{N}_i} \mathbf{L}_{ij} \mathbf{x}_j^{k+1}. \quad (65)$$

Each node i should do — In each iteration k :

1. Update $\mathbf{y}_i^k, \mathbf{x}_i^{k+1}, \tilde{\mathbf{x}}_i^{k+1}$ by (62) ~ (64), respectively
2. Send \mathbf{x}_i^{k+1} to its neighbors
3. Wait until receiving \mathbf{x}_j^{k+1} from all its neighbors, $j \in \mathcal{N}_i$
4. Update \mathbf{v}_i^{k+1} by (65)

Convergence Rate of Acc-DLADMM

Theorem II.4

Suppose that each $f_i(\mathbf{x}_i)$ is μ -strongly convex and L -smooth, $i \in [m]$. Assume that $\frac{2d_{\max}}{\sigma_{\mathbf{L}}} \leq \frac{L}{\mu}$, where $\sigma_{\mathbf{L}}$ is the smallest non-zero singular value of \mathbf{L} . Let

$$\alpha = \frac{1}{4L}, \beta = \frac{L}{d_{\max}}, \text{ and } \theta = \sqrt{\frac{2\mu d_{\max}}{L\sigma_{\mathbf{L}}}}. \quad (66)$$

Convergence Rate of Acc-DLADMM

Similar to Theorem III.1 of *ADMM Slide 2*, we have the following convergence rate for Acc-DLADMM.

Theorem II.4 (cont'd)

Then for Acc-DLADMM, we have

$$l_{k+1} \leq \mathcal{O}\left(1 - \sqrt{\frac{\mu\sigma_{\mathbf{L}}}{2Ld_{\max}}}\right) l_k, \quad (67)$$

where

$$\begin{aligned} l_k = & (1 - \theta) \left(f(\tilde{\mathbf{X}}^k) - f(\mathbf{X}^*) + \langle \boldsymbol{\Omega}^*, \mathbf{W}\tilde{\mathbf{X}}^k \rangle \right) \\ & + \frac{\theta^2}{2\alpha} \|\mathbf{X}^k - \mathbf{X}^*\|^2 + \frac{1}{2\beta} \|\boldsymbol{\Omega}^k - \boldsymbol{\Omega}^*\|^2. \end{aligned} \quad (68)$$

Complexity Comparisons

For decentralized optimization, the algorithm complexity are listed as follows.

- ▶ DADMM: $\mathcal{O}\left(\left(\sqrt{\frac{Ld_{\max}}{\mu\sigma_{\mathbf{L}}}} + \frac{d_{\max}}{\sigma_{\mathbf{L}}}\right) \log \frac{1}{\epsilon}\right)$
- ▶ DLADMM: $\mathcal{O}\left(\left(\frac{L}{\mu} + \frac{d_{\max}}{\sigma_{\mathbf{L}}}\right) \log \frac{1}{\epsilon}\right)$
- ▶ Acc-DLADMM: $\mathcal{O}\left(\sqrt{\frac{Ld_{\max}}{\mu\sigma_{\mathbf{L}}}} \log \frac{1}{\epsilon}\right)$

Outline

I Centralized Distributed Optimization

II Decentralized Distributed Optimization

III Asynchronous Distributed ADMM

IV Nonconvex Distributed ADMM

V ADMM with Generally Linear Constraints

References

Asynchronous Distributed ADMM

Now let's go back to the centralized optimization. CADMM, CLADMM, and Acc-CLADMM are executed in a synchronous manner. That is, the master needs to wait for all the workers to finish their updates before it can proceed. When the workers have different delays, the master has to wait for the slowest worker before the next iteration, i.e., the system proceeds at the pace of the slowest worker.

In the asynchronous ADMM, the master does not wait for all the workers, but proceeds as long as it receives information from a partial set of workers instead.

Asynchronous Distributed ADMM

In the following, we describe how to adapt CADMM to the asynchronous version and gives its convergence analysis.

We denote the partial set at iteration k as \mathcal{A}^k , and \mathcal{A}_c^k as the complementary set, which means the set of workers whose information does not arrive at iteration k . We use α to lower bound the size of \mathcal{A}^k . In the asynchronous ADMM, we often require that the master has to receive the updates from every worker at least once in every τ iterations. That is, we do not allow some workers to be absent for a long time. So we make the following bounded delay assumption.

Assumption III.1

The maximum tolerable delay for all i and k is upper bounded. Denote the upper bound as τ , then it must be that for every i ,

$$i \in \mathcal{A}^k \cup \mathcal{A}^{k-1} \cup \dots \cup \mathcal{A}^{\max\{k-\tau+1, 0\}}. \quad (69)$$

Async-ADMM-M

The master works as follows.

1. Initialize $\tilde{d}_1^1 = \dots = \tilde{d}_m^1 = 0$

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

2.1 Wait until receiving $\hat{\mathbf{x}}_i^k$ and $\hat{\boldsymbol{\lambda}}_i^k$ from workers $i \in \mathcal{A}^k$ such that $|\mathcal{A}^k| \geq \alpha$ and $\tilde{d}_j^k < \tau - 1$ for all $j \in \mathcal{A}_c^k$

2.2 $\mathbf{x}_i^{k+1} = \begin{cases} \hat{\mathbf{x}}_i^k & \forall i \in \mathcal{A}^k \\ \mathbf{x}_i^k & \forall i \in \mathcal{A}_c^k \end{cases}, \boldsymbol{\lambda}_i^{k+1} = \begin{cases} \hat{\boldsymbol{\lambda}}_i^k & \forall i \in \mathcal{A}^k \\ \boldsymbol{\lambda}_i^k & \forall i \in \mathcal{A}_c^k \end{cases}, \tilde{d}_i^{k+1} = \begin{cases} 0 & \forall i \in \mathcal{A}^k \\ \tilde{d}_i^k + 1 & \forall i \in \mathcal{A}_c^k \end{cases}$

2.3 $\mathbf{z}^{k+1} = \operatorname{argmin}_{\mathbf{z}} \left[\sum_{i=1}^m \left(\langle \boldsymbol{\lambda}_i^{k+1}, \mathbf{x}_i^{k+1} - \mathbf{z} \rangle + \frac{\beta}{2} \|\mathbf{x}_i^{k+1} - \mathbf{z}\|^2 \right) + \frac{\rho}{2} \|\mathbf{z} - \mathbf{z}^k\|^2 \right] = \frac{1}{\rho + m\beta} \left[\rho \mathbf{z}^k + \sum_{i=1}^m (\boldsymbol{\lambda}_i^{k+1} + \beta \mathbf{x}_i^{k+1}) \right]$

2.4 Broadcast \mathbf{z}^{k+1} to the workers in \mathcal{A}^k

Async-ADMM-W

The i -th worker work as follows.

1. Initialize $\hat{\mathbf{x}}_i^0$ and $\hat{\boldsymbol{\lambda}}_i^0$, $i \in [m]$
2. **for** $k_i = 1, 2, \dots$ **do**
 - 2.1 Wait until receiving \mathbf{z} from the master
 - 2.2 $\hat{\mathbf{x}}_i^{k_i+1} = \operatorname{argmin}_{\mathbf{x}_i} \left(f_i(\mathbf{x}_i) + \langle \hat{\boldsymbol{\lambda}}_i^{k_i}, \mathbf{x}_i - \mathbf{z} \rangle + \frac{\beta}{2} \|\mathbf{x}_i - \mathbf{z}\|^2 \right) = \operatorname{prox}_{\beta^{-1}f_i} \left(\mathbf{z} - \frac{1}{\beta} \hat{\boldsymbol{\lambda}}_i^{k_i} \right)$
 - 2.3 $\hat{\boldsymbol{\lambda}}_i^{k_i+1} = \hat{\boldsymbol{\lambda}}_i^{k_i} + \beta(\hat{\mathbf{x}}_i^{k_i+1} - \mathbf{z})$
 - 2.4 Send $\hat{\mathbf{x}}_i^{k_i+1}$ and $\hat{\boldsymbol{\lambda}}_i^{k_i+1}$ to the master

Note that we introduce d_i , the amount of delay, for each worker such that the bounded delay assumption holds. The master must wait if there exists one worker with $\tilde{d}_i = \tau - 1$.

We name the two subprocedures as *Async-ADMM-M* and *Async-ADMM-W*, respectively.

Existence of Convergence

To simplify the analysis, we rewrite the method from the master's point of view:

$$\mathbf{x}_i^{k+1} = \begin{cases} \operatorname{argmin}_{\mathbf{x}_i} \left(f_i(\mathbf{x}_i) + \langle \boldsymbol{\lambda}_i^{\bar{k}_i+1}, \mathbf{x}_i \rangle + \frac{\beta}{2} \|\mathbf{x}_i - \mathbf{z}^{\bar{k}_i+1}\|^2 \right) & i \in \mathcal{A}^k \\ \mathbf{x}_i^k & i \in \mathcal{A}_c^k, \end{cases} \quad (70)$$

$$\boldsymbol{\lambda}_i^{k+1} = \begin{cases} \boldsymbol{\lambda}_i^{\bar{k}_i+1} + \beta(\mathbf{x}_i^{k+1} - \mathbf{z}^{\bar{k}_i+1}) & i \in \mathcal{A}^k \\ \boldsymbol{\lambda}_i^k & i \in \mathcal{A}_c^k, \end{cases} \quad (71)$$

$$\mathbf{z}^{k+1} = \operatorname{argmin}_{\mathbf{z}} \left[\sum_{i=1}^m \left(\langle \boldsymbol{\lambda}_i^{k+1}, \mathbf{x}_i^{k+1} - \mathbf{z} \rangle + \frac{\beta}{2} \|\mathbf{x}_i^{k+1} - \mathbf{z}\|^2 \right) + \frac{\rho}{2} \|\mathbf{z} - \mathbf{z}^k\|^2 \right], \quad (72)$$

where \bar{k}_i is the last iteration before iteration k for which worker $i \in \mathcal{A}^k$ arrives, i.e., $i \in \mathcal{A}^{\bar{k}_i}$.

Existence of Convergence

Thus, for all workers $i \in \mathcal{A}^k$, we have

$$\begin{aligned} \mathbf{x}_i^{\bar{k}_i+1} &= \mathbf{x}_i^{\bar{k}_i+2} = \dots = \mathbf{x}_i^k, \\ \boldsymbol{\lambda}_i^{\bar{k}_i+1} &= \boldsymbol{\lambda}_i^{\bar{k}_i+2} = \dots = \boldsymbol{\lambda}_i^k, \text{ and} \end{aligned} \tag{73}$$

$$\max\{k - \tau, 0\} \leq \bar{k}_i \leq k. \tag{74}$$

For each $i \in \mathcal{A}_c^k$ we denote \tilde{k}_i as the last iteration before iteration k for which worker i arrives, i.e., $i \in \mathcal{A}^{\tilde{k}_i}$. Under the bounded delay assumption, we have

$$\max\{k - \tau + 1, 0\} \leq \tilde{k}_i < k. \tag{75}$$

Thus, for all workers $i \in \mathcal{A}_c^k$, we have

$$\mathbf{x}_i^{\tilde{k}_i+1} = \mathbf{x}_i^{\tilde{k}_i+2} = \dots = \mathbf{x}_i^k = \mathbf{x}_i^{k+1}, \tag{76}$$

$$\boldsymbol{\lambda}_i^{\tilde{k}_i+1} = \boldsymbol{\lambda}_i^{\tilde{k}_i+2} = \dots = \boldsymbol{\lambda}_i^k = \boldsymbol{\lambda}_i^{k+1}. \tag{77}$$

Existence of Convergence

We also denote \hat{k}_i as the last iteration before \tilde{k}_i for which $i \in \mathcal{A}^{\tilde{k}_i}$ arrives, i.e., $i \in \mathcal{A}^{\hat{k}_i}$. We also have

$$\max\{\tilde{k}_i - \tau, 0\} \leq \hat{k}_i < \tilde{k}_i. \quad (78)$$

Thus, for all workers $i \in \mathcal{A}_c^k$, we have

$$\mathbf{x}_i^{k+1} = \mathbf{x}_i^{\tilde{k}_i+1} = \underset{\mathbf{x}_i}{\operatorname{argmin}} \left(f_i(\mathbf{x}_i) + \langle \boldsymbol{\lambda}_i^{\hat{k}_i+1}, \mathbf{x}_i \rangle + \frac{\beta}{2} \|\mathbf{x}_i - \mathbf{z}^{\hat{k}_i+1}\|^2 \right), \quad (79)$$

$$\boldsymbol{\lambda}_i^{k+1} = \boldsymbol{\lambda}_i^{\tilde{k}_i+1} = \boldsymbol{\lambda}_i^{\hat{k}_i+1} + \beta(\mathbf{x}_i^{\tilde{k}_i+1} - \mathbf{z}^{\hat{k}_i+1}), \quad (80)$$

and

$$\mathbf{x}_i^{\hat{k}_i+1} = \mathbf{x}_i^{\hat{k}_i+2} = \dots = \mathbf{x}_i^{\tilde{k}_i}, \quad (81)$$

$$\boldsymbol{\lambda}_i^{\hat{k}_i+1} = \boldsymbol{\lambda}_i^{\hat{k}_i+2} = \dots = \boldsymbol{\lambda}_i^{\tilde{k}_i}. \quad (82)$$

Existence of Convergence

Denote $(\mathbf{x}_1^*, \dots, \mathbf{x}_m^*, \mathbf{z}^*, \boldsymbol{\lambda}_1^*, \dots, \boldsymbol{\lambda}_m^*)$ to be a KKT point. We have

$$\sum_{i=1}^m \boldsymbol{\lambda}_i^* = \mathbf{0}, \mathbf{z}^* = \mathbf{x}_i^*, \text{ and } \nabla f_i(\mathbf{x}_i^*) + \boldsymbol{\lambda}_i^* = \mathbf{0}, i \in [m]. \quad (83)$$

Also denote $f^* = \sum_{i=1}^m f_i(\mathbf{z}^*)$.

Theorem III.1

Suppose that each $f_i(\mathbf{x}_i)$ is μ -strongly convex and L -smooth, $i \in [m]$, and Assumption III.1 holds. Let

$$\begin{aligned} \beta &> \frac{1 + L^2 + \sqrt{(1 + L^2)^2 + 8L^2}}{2}, \\ \rho &> \frac{1}{2} [m(1 + \beta^2)(\tau - 1)^2 - m\beta]. \end{aligned} \quad (84)$$

Existence of Convergence

Theorem III.1 (cont'd)

Suppose that $(\mathbf{x}_1^k, \dots, \mathbf{x}_m^k, \mathbf{z}^k, \boldsymbol{\lambda}_1^k, \dots, \boldsymbol{\lambda}_m^k)$ generated by (70) \sim (72) are bounded, then $(\mathbf{x}_1^k, \dots, \mathbf{x}_m^k, \mathbf{z}^k, \boldsymbol{\lambda}_1^k, \dots, \boldsymbol{\lambda}_m^k)$ converge to the set of KKT points of (2) in the sense of

$$\sum_{i=1}^m \boldsymbol{\lambda}_i^k \rightarrow \mathbf{0}, \quad (85)$$

$$\mathbf{x}_i^{k+1} - \mathbf{z}^{k+1} \rightarrow \mathbf{0}, \quad (86)$$

$$\nabla f_i(\mathbf{x}_i^{k+1}) + \boldsymbol{\lambda}_i^{k+1} = \mathbf{0}, i \in [m]. \quad (87)$$

Linear Convergence Rate of Async-ADMM

Theorem III.2

Suppose that each $f_i(\mathbf{x}_i)$ is μ -strongly convex and L -smooth, $i \in [m]$, and Assumption III.1 holds. Let β and ρ be large enough such that

$$8m(\beta - \mu) \leq \rho, \quad (88)$$

$$\frac{m\beta + 2\rho}{2} - 1 - \tau 2^{2\tau} - \left(\frac{1 + \beta^2}{2} + \frac{1}{2m}\right)m\tau 2^\tau \geq 0, \quad (89)$$

$$\frac{\beta}{2} - \frac{L^2}{\beta} - \frac{L^2}{2} - \frac{1}{2} - \frac{L^2}{4m\beta^2} - \frac{L^2}{4m\beta^2} 2^{\tau-1}\tau > 0. \quad (90)$$

Linear Convergence Rate of Async-ADMM

Theorem III.2 (cont'd)

Then we have

$$\begin{aligned} L_{\beta}(\mathbf{x}^{K+1}, \mathbf{z}^{K+1}, \boldsymbol{\lambda}^{K+1}) - f^* \\ \leq \left(1 + \frac{1}{\delta\rho}\right)^{-(K+1)} (L_{\beta}(\mathbf{x}^0, \mathbf{z}^0, \boldsymbol{\lambda}^0) - f^*), \end{aligned} \quad (91)$$

where $\delta \geq \max\{1, \frac{1}{\rho}, \frac{\rho+m\beta}{m\mu} - 1\}$.

We believe that in general the asynchronous ADMM needs more iterations than synchronous ADMM. It is unclear whether the time saved per iteration of the asynchronous ADMM can offset the cost of more iterations in theory, although it shows great advantages in practice.

Outline

I Centralized Distributed Optimization

II Decentralized Distributed Optimization

III Asynchronous Distributed ADMM

IV Nonconvex Distributed ADMM

V ADMM with Generally Linear Constraints

References

Nonconvex Distributed ADMM

Async-ADMM can also be used to solve nonconvex problems. In this case, we only assume that each f_i is L -smooth. Then, $L_\beta(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda})$ is $(\beta - L)$ -strongly convergence w.r.t. \mathbf{x} , and we have the following intermediate result:

$$\begin{aligned} & L_\beta(\mathbf{x}^{k+1}, \mathbf{z}^k, \boldsymbol{\lambda}^k) - L_\beta(\mathbf{x}^k, \mathbf{z}^k, \boldsymbol{\lambda}^k) \\ & \leq \sum_{i \in \mathcal{A}^k} \left(\beta \langle \mathbf{z}^{\bar{k}_i+1} - \mathbf{z}^k, \mathbf{x}_i^{k+1} - \mathbf{x}_i^k \rangle - \frac{\beta - L}{2} \|\mathbf{x}_i^{k+1} - \mathbf{x}_i^k\|^2 \right). \end{aligned} \tag{92}$$

Existence of Convergence

For nonconvex distributed optimization, Async-ADMM also converges.

Theorem IV.1

Suppose that each $f_i(\mathbf{x}_i)$ is L -smooth, $i \in [m]$, and Assumption III.1 holds. Let

$$\begin{aligned}\beta &> \frac{1 + L + L^2 + \sqrt{(1 + L + L^2)^2 + 8L^2}}{2}, \\ \rho &> \frac{m(1 + \beta^2)(\tau - 1)^2 - m\beta}{2}.\end{aligned}\tag{93}$$

Suppose that $(\mathbf{x}_1^k, \dots, \mathbf{x}_m^k, \mathbf{z}^k, \boldsymbol{\lambda}_1^k, \dots, \boldsymbol{\lambda}_m^k)$ generated by (70) \sim (72) are bounded, then $(\mathbf{x}_1^k, \dots, \mathbf{x}_m^k, \mathbf{z}^k, \boldsymbol{\lambda}_1^k, \dots, \boldsymbol{\lambda}_m^k)$ converge to the set of KKT points of (2) in the sense of (85) \sim (87) hold.

Outline

I Centralized Distributed Optimization

II Decentralized Distributed Optimization

III Asynchronous Distributed ADMM

IV Nonconvex Distributed ADMM

V ADMM with Generally Linear Constraints

References

Distributed Optimization with Generally Linear Constraints

In *ADMM Slide: Part 3*, we introduce the following multi-block problem:

$$\min_{\mathbf{x}_i} \sum_{i=1}^m f_i(\mathbf{x}_i), \quad s.t. \quad \sum_{i=1}^m \mathbf{A}_i \mathbf{x}_i = \mathbf{b}. \quad (94)$$

In the following, we present the distributed version of LADMM-PS in a centralized network. We name it *CLADMM-PS-M* and *CLADMM-PS-W*, for the master and the workers, respectively.

CLADMM-PS-M

The master should do:

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

1.1 Wait until receiving \mathbf{y}_i^{k+1} from all the workers $i \in [m]$

1.2 $\mathbf{s}^{k+1} = \sum_{i=1}^m \mathbf{y}_i^{k+1}$

1.3 $\boldsymbol{\lambda}^{k+1} = \boldsymbol{\lambda}^k + \beta(\mathbf{s}^{k+1} - \mathbf{b})$

1.4 Send \mathbf{s}^{k+1} and $\boldsymbol{\lambda}^{k+1}$ to all the workers

CLADMM-PS-W

Each worker i should do:

1. Initialize \mathbf{x}_i^0 and $\boldsymbol{\lambda}_i^0$, $i \in [m]$
2. $\mathbf{y}_i^0 = \mathbf{A}_i \mathbf{x}_i^0$
3. Send \mathbf{y}_i^0 to the master
4. Wait until receiving \mathbf{s}^0 and $\boldsymbol{\lambda}^0$ from the master
5. for $k = 0, 1, 2, \dots$ do
 - 5.1 $\mathbf{x}_i^{k+1} = \operatorname{argmin}_{\mathbf{x}_i} \left(f_i(\mathbf{x}_i) + \langle \boldsymbol{\lambda}^k, \mathbf{A}_i \mathbf{x}_i \rangle + \beta \langle \mathbf{A}_i^T (\mathbf{s}^k - \mathbf{b}), \mathbf{x}_i - \mathbf{x}_i^k \rangle + \frac{m\beta \|\mathbf{A}_i\|^2}{2} \|\mathbf{x}_i - \mathbf{x}_i^k\|^2 \right) =$
 $\operatorname{prox}_{(m\beta \|\mathbf{A}_i\|^2)^{-1} f_i} \left(\mathbf{x}_i^k - \frac{1}{m\beta \|\mathbf{A}_i\|^2} \mathbf{A}_i^T (\boldsymbol{\lambda}^k + \beta (\mathbf{s}^k - \mathbf{b})) \right)$
 - 5.2 $\mathbf{y}_i^{k+1} = \mathbf{A}_i \mathbf{x}_i^{k+1}$
 - 5.3 Send \mathbf{y}_i^{k+1} to the master
 - 5.4 Wait until receiving \mathbf{s}^{k+1} and $\boldsymbol{\lambda}^{k+1}$ from the master

The convergence result of LADMM-PS can be kept.

Outline

I Centralized Distributed Optimization

II Decentralized Distributed Optimization

III Asynchronous Distributed ADMM

IV Nonconvex Distributed ADMM

V ADMM with Generally Linear Constraints

References

References

1. Lin, Zhouchen, Huan Li, and Cong Fang. *Alternating Direction Method of Multipliers for Machine Learning*. Springer Nature, 2022.
2. Boyd, Stephen, Stephen P. Boyd, and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.