

Modern Convex Optimization Methods

Mingyi Hong

University Of Minnesota

M. Hong would like to thank Dr. Prashant Khanduri for helping prepare the slides.

Outline

- **Alternating Direction Method of Multipliers (ADMM)**
 - Dual Ascent
 - Dual Decomposition
 - Augmented Lagrangian Method
 - ADMM: Convergence
- **Decentralized ADMM**
 - Assumptions
 - Convergence: Proof Sketch
- **EXTRA**
 - Assumptions
 - Convergence
- **Connection:** EXTRA and Primal-Dual Methods

ADMM Basics: Dual Ascent¹

- **Problem**

$$\begin{array}{ll}\text{minimize}_{\mathbf{x}} & f(\mathbf{x}) \\ \text{subject to} & A\mathbf{x} = \mathbf{b}\end{array}$$

with $\mathbf{x} \in \mathbb{R}^n$, with $A \in \mathbb{R}^{p \times n}$ and $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex

- **Dual Ascent**

$$\begin{aligned}\mathbf{x}^{k+1} &= \underset{\mathbf{x}}{\operatorname{argmin}} \{L(\mathbf{x}, \mathbf{y}^k) = f(\mathbf{x}) + (\mathbf{y}^k)^T (A\mathbf{x} - \mathbf{b})\} \\ \mathbf{y}^{k+1} &= \mathbf{y}^k + \alpha^k (A\mathbf{x}^{k+1} - \mathbf{b})\end{aligned}$$

with $\alpha^k > 0$ being the step-size

¹Boyd, et al., Distributed optimization and statistical learning via the alternating direction method of multipliers. Now Publishers Inc, 2011.

Dual Decomposition

- Problem

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \sum_{i=1}^m f_i(\mathbf{x}_i) \\ & \text{subject to} && A\mathbf{x} = \mathbf{b} \end{aligned}$$

where $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_m) \in \mathbb{R}^n$ with $\mathbf{x}_i \in \mathbb{R}^{n_i}$ with $i \in [m]$ and

$$A = [A_1, \dots, A_m] \text{ with } A_i \in \mathbb{R}^{p \times n_i} \Rightarrow A\mathbf{x} = \sum_{i=1}^m A_i \mathbf{x}_i$$

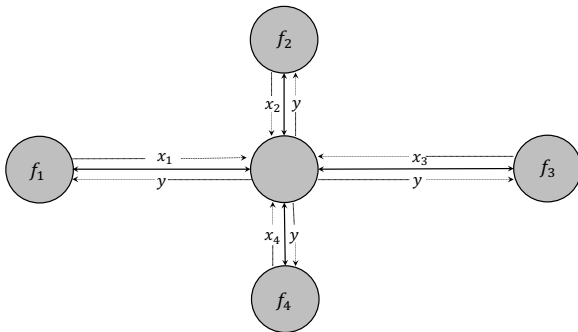
- Dual Decomposition (the Lagrangian function becomes separable over i)

$$\mathbf{x}_i^{k+1} = \underset{\mathbf{x}_i}{\operatorname{argmin}} \{ L_i(\mathbf{x}_i, \mathbf{y}^k) = f_i(\mathbf{x}_i) + (\mathbf{y}^k)^T A_i \mathbf{x}_i - (1/m)(\mathbf{y}^k)^T \mathbf{b} \}$$

$$\mathbf{y}^{k+1} = \mathbf{y}^k + \alpha^k (A\mathbf{x}^{k+1} - \mathbf{b})$$

Properties: Dual Decomposition

- Implementation



Dual decomposition with a central server and $m = 4$.

- Pro:** Decomposes across x for each $i \in [m]$
- Con:** Requires strong convexity to ensure convergence
- Solution:** Use Augmented Lagrangian method

ADMM Basics: Augmented Lagrangian Method

- **Problem:** For $\rho > 0$,

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && f(\mathbf{x}) + \frac{\rho}{2} \|A\mathbf{x} - b\|^2 \\ & \text{subject to} && A\mathbf{x} = b \end{aligned}$$

- **Augmented Lagrangian Method**

$$\mathbf{x}^{k+1} = \underset{\mathbf{x}}{\text{argmin}} \underbrace{\{L_\rho(\mathbf{x}, \mathbf{y}^k) = f(\mathbf{x}) + (\mathbf{y}^k)^T (A\mathbf{x} - b) + \frac{\rho}{2} \|A\mathbf{x} - b\|^2\}}_{\text{Augmented Lagrangian}}$$

$$\mathbf{y}^{k+1} = \mathbf{y}^k + \alpha^k (A\mathbf{x}^{k+1} - b)$$

- **Pro:** Better convergence
- **Con:** Does not decompose for $f(\mathbf{x}) = \sum_{i=1}^m f_i(\mathbf{x}_i)$
- **Solution:** ADMM blends decomposability of dual ascent with superior convergence properties of method of multipliers

Alternating Direction Method of Multipliers (ADMM)

- **Problem** Let $f : \mathbb{R}^{n_x} \rightarrow \mathbb{R}$ and $g : \mathbb{R}^{n_z} \rightarrow \mathbb{R}$ be convex

$$\begin{aligned} & \text{minimize}_{\mathbf{x}, \mathbf{z}} && f(\mathbf{x}) + g(\mathbf{z}) \\ & \text{subject to} && A\mathbf{x} + B\mathbf{z} = \mathbf{c} \end{aligned}$$

with $\mathbf{x} \in \mathbb{R}^{n_x}$ and $\mathbf{z} \in \mathbb{R}^{n_z}$ and $A \in \mathbb{R}^{p \times n_x}$, $B \in \mathbb{R}^{p \times n_z}$

- **Augmented Lagrangian**

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

- **ADMM**

$$\mathbf{x}^{k+1} = \underset{\mathbf{x}}{\operatorname{argmin}} L_\rho(\mathbf{x}, \mathbf{z}^k, \mathbf{y}^k)$$

$$\mathbf{z}^{k+1} = \underset{\mathbf{z}}{\operatorname{argmin}} L_\rho(\mathbf{x}^{k+1}, \mathbf{z}, \mathbf{y}^k)$$

$$\mathbf{y}^{k+1} = \mathbf{y}^k + \rho(A\mathbf{x}^{k+1} + B\mathbf{z}^{k+1} - \mathbf{c})$$

Performance of ADMM

LASSO²: Given $b \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times p}$ solve:

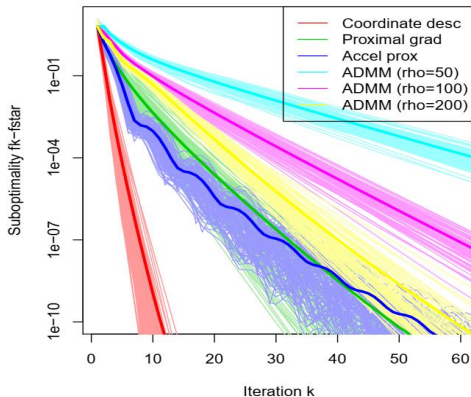
$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{b} - A\mathbf{x}\|^2 + \lambda \|\mathbf{x}\|_1$$

Rephrased in form suitable for ADMM as:

$$\min_{\mathbf{x}, \mathbf{z}} \underbrace{\frac{1}{2} \|\mathbf{b} - A\mathbf{x}\|^2}_{f(\mathbf{x})} + \underbrace{\lambda \|\mathbf{z}\|_1}_{g(\mathbf{z})} \quad \text{subject to} \quad \underbrace{\mathbf{x} = \mathbf{z}}_{\text{linear constraint}}$$

ADMM Comparison³

Comparison of various algorithms for lasso regression: 100 random instances with $n = 200$, $p = 50$



³Ryan Tibshirani, Convex Optimization, Lecture Slides, Fall 2018

Properties of ADMM

- **Slow** to converge to high accuracy
- Converges to **modest accuracy** within a few iterations
 - Sufficient for many applications
 - Well suited for large scale problems in machine learning and statistical estimation
- Computations can be **distributed** across multiple nodes

Assumptions

Assumption 1

The (extended-real-valued) functions $f : \mathbb{R}^{n_x} \rightarrow \mathbb{R} \cup \{+\infty\}$ and $g : \mathbb{R}^{n_z} \rightarrow \mathbb{R} \cup \{+\infty\}$ are closed, proper, and convex.

**A function f satisfies Assumption 1 \iff
 $\text{epi} f = \{(x, t) \in \mathbb{R}^n \times \mathbb{R} : f(x) \leq t\}$ closed, non-empty, convex set**

For simplicity, we will assume that both f and g **differentiable**.
Otherwise, subgradient notation will be used.

Assumptions

Assumption 2

The unaugmented Lagrangian defined as:

$$L_0(\mathbf{x}, \mathbf{z}, \mathbf{y}) = f(\mathbf{x}) + g(\mathbf{z}) + \mathbf{y}^T (A\mathbf{x} + B\mathbf{z} - \mathbf{c})$$

has a saddle point.

There exists $(\mathbf{x}^*, \mathbf{z}^*, \mathbf{y}^*)$ such that (possibly non-unique)

$$L_0(\mathbf{x}^*, \mathbf{z}^*, \mathbf{y}) \leq L_0(\mathbf{x}^*, \mathbf{z}^*, \mathbf{y}^*) \leq L_0(\mathbf{x}, \mathbf{z}, \mathbf{y}^*) \text{ for all } \mathbf{x}, \mathbf{z}, \mathbf{y}$$

Convergence

Under **Assumption 1** and **Assumption 2**, ADMM iterates satisfy:

- **Residual convergence:** The iterates approach feasibility, i.e., the residuals satisfy

$$r^k = A\mathbf{x}^k + B\mathbf{z}^k - c \rightarrow 0 \text{ as } k \rightarrow \infty$$

- **Objective convergence:** $f(\mathbf{x}^k) + g(\mathbf{z}^k) \rightarrow p^*$ as $k \rightarrow \infty$, where p^* is

$$p^* = \inf\{f(\mathbf{x}) + g(\mathbf{z}) : A\mathbf{x} + B\mathbf{z} - c = 0\}$$

- **Dual variable convergence:** $\mathbf{y}^k \rightarrow \mathbf{y}^*$ as $k \rightarrow \infty$ where \mathbf{y}^* is a dual optimal value

Optimality of ADMM

Necessary and sufficient conditions for optimality (KKT conditions)

- **Primal Feasibility**

$$Ax^* + Bz^* - c = 0 \quad (1.1)$$

- **Stationarity**

$$0 = \nabla f(x^*) + A^T y^* \quad (1.2)$$

$$0 = \nabla g(z^*) + B^T y^* \quad (1.3)$$

- Using $z^{k+1} = \operatorname{argmin}_z L_\rho(x^{k+1}, z, y^k)$ **check!** that (1.3) is **always satisfied** by z^{k+1} and y^{k+1}
- **To achieve optimality:** (1.1) and (1.2) need to be satisfied

Optimality of ADMM

- **Dual Residual:** Stationarity condition (1.2)

Using $\mathbf{x}^{k+1} = \operatorname{argmin}_x L_\rho(\mathbf{x}, \mathbf{z}^k, \mathbf{y}^k)$, we get (**Check!**)

$$\underbrace{\rho A^T B(\mathbf{z}^{k+1} - \mathbf{z}^k)}_{\mathbf{s}^{k+1}} = \partial f(\mathbf{x}^{k+1}) + A^T \mathbf{y}^{k+1}$$

\mathbf{s}^{k+1} is referred to a *dual residual*

- **Primal Residual:** Denote

$$\mathbf{r}^{k+1} = A\mathbf{x}^{k+1} + B\mathbf{z}^{k+1} - \mathbf{c}$$

as the *primal residual*

Goal:

$\mathbf{r}^k, \mathbf{s}^k \rightarrow 0 \Rightarrow (1.1) \text{ and } (1.2) \text{ hold} \Rightarrow \text{Optimality of ADMM}$

Proof Sketch: ADMM I

Define: Lyupanov Function:

$$V^k = \frac{1}{\rho} \|\mathbf{y}^k - \mathbf{y}^*\|^2 + \rho \|B(\mathbf{z}^k - \mathbf{z}^*)\|^2$$

The proof relies on three basic inequalities:

- **Inequality I:**

$$V^{k+1} \leq V^k - \rho \|r^{k+1}\|^2 - \rho \|B(\mathbf{z}^{k+1} - \mathbf{z}^k)\|^2$$

- **Inequality II:**

$$p^{k+1} - p^* \leq -(\mathbf{y}^{k+1})^T r^{k+1} - \rho (B(\mathbf{z}^{k+1} - \mathbf{z}^k))^T (-r^{k+1} + B(\mathbf{z}^{k+1} - \mathbf{z}^*))$$

- **Inequality III:**

$$p^* - p^{k+1} \leq \mathbf{y}^{*T} r^{k+1}$$

Proof Sketch: ADMM II

- Iterating **Inequality I** above, we get

$$\rho \sum_{k=0}^{\infty} (\|r^{k+1}\|^2 + \|B(z^{k+1} - z^k)\|^2) \leq V_0$$

this implies $r^k \rightarrow 0$ and $B(z^{k+1} - z^k) \rightarrow 0$ as $k \rightarrow \infty$

- $B(z^{k+1} - z^k) \rightarrow 0$ further implies that the *dual residual* $s^{k+1} \rightarrow 0$ (follows from definition of s^k)

- Inequality II** and **Inequality III** imply $p^k \rightarrow p^*$ as $k \rightarrow \infty$, i.e., objective convergence

Inequalities II and III are used to derive Inequality I

Distributed ADMM

- Problem

$$\text{minimize}_{\mathbf{x}_1, \dots, \mathbf{x}_m, \mathbf{z}} \sum_{i=1}^m f_i(\mathbf{x}_i) \quad \text{subject to} \quad \mathbf{x}_i = \mathbf{z}, i = 1, \dots, m$$

- ADMM steps can be **distributed** across m nodes:

$$\mathbf{x}_i^{k+1} = \underset{\mathbf{x}_i}{\operatorname{argmin}} f_i(\mathbf{x}_i) + \frac{\rho}{2} \|\mathbf{x}_i - \mathbf{z}^k + \mathbf{u}_i^k\|^2$$

$$\mathbf{z}^{k+1} = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i^{k+1} + \mathbf{u}_i^k)$$

$$\mathbf{u}_i^{k+1} = \mathbf{u}_i^k + \mathbf{x}_i^{k+1} - \mathbf{z}^{k+1}$$

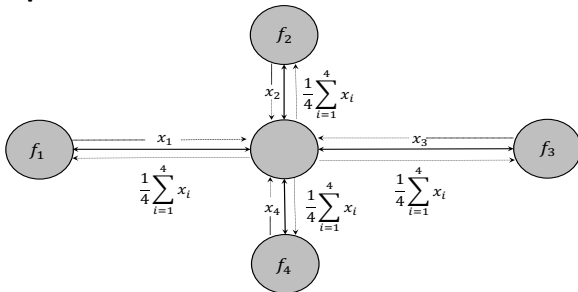
Distributed Implementation

- Simple manipulation yields (**Check!**):

$$\mathbf{x}_i^{k+1} = \underset{\mathbf{x}_i}{\operatorname{argmin}} f_i(\mathbf{x}_i) + \frac{\rho}{2} \left\| \mathbf{x}_i - \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^k + \mathbf{u}_i^k \right\|^2$$

$$\mathbf{u}_i^{k+1} = \mathbf{u}_i^k + \mathbf{x}_i^{k+1} - \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^{k+1}$$

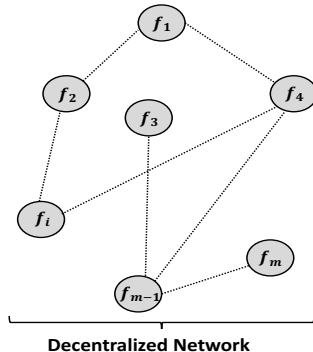
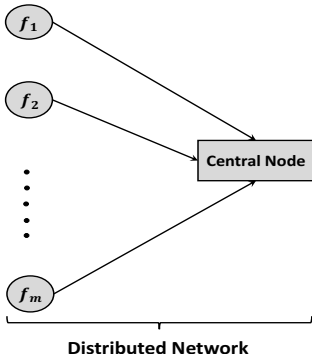
- Parallel Updates



Distributed ADMM with $m = 4$.

Decentralized Setup

- **Distributed ADMM:** Central server for sharing of iterates
 - **Congestion, Privacy concerns, Single point of failure**
- **Decentralized ADMM:** No central server, nodes exchange information with their immediate neighbors
 - **No risk of network congestion, Better privacy, Robust to individual node failures**



Decentralized ADMM⁴

- **Problem**

$$\text{minimize}_x \sum_{i=1}^m f_i(x)$$

- **Model**

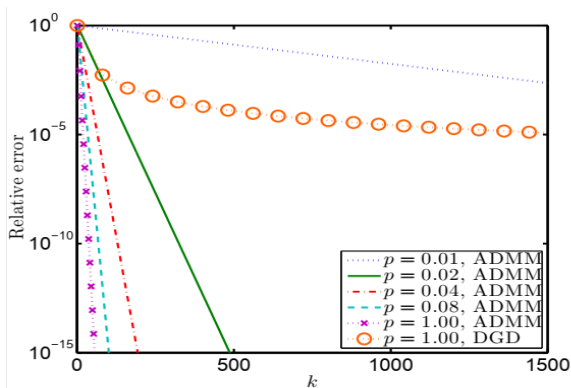
- m agents/nodes present in the network
- Each agent has access to local function $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$

- **Network**

- m agents connected via E edges ($2E$ arcs)
- **Symmetric directed graph:** $\mathcal{G} = \{\mathcal{V}, \mathcal{A}\}$
- **Set of vertices:** \mathcal{V} with $|\mathcal{V}| = m$, **Set of Arcs:** \mathcal{A} with $|\mathcal{A}| = 2E$

⁴Shi, et al., On the linear convergence of the ADMM in decentralized consensus optimization, IEEE Transactions on Signal Processing 62.7 (2014): 1750-1761.

ADMM vs DGD



Relative error for least squares consensus ADMM vs DGD with step-size $\frac{1}{k^{1/3}}$

Problem Formulation

- Reformulation suitable for ADMM

$$\begin{aligned} & \text{minimize}_{\{\mathbf{x}_i\}, \{\mathbf{z}_{ij}\}} \sum_{i=1}^m f_i(\mathbf{x}_i) \\ & \text{subject to } \mathbf{x}_i = \mathbf{z}_{ij}, \mathbf{x}_j = \mathbf{z}_{ij}, \forall (i, j) \in \mathcal{A} \end{aligned}$$

- Standard ADMM form

$$\begin{aligned} & \text{minimize}_{\mathbf{x}, \mathbf{z}} f(\mathbf{x}) + g(\mathbf{z}) \\ & \text{subject to } A\mathbf{x} + B\mathbf{z} = 0 \end{aligned}$$

with $\mathbf{x} \in \mathbb{R}^{mn}$, $\mathbf{z} \in \mathbb{R}^{2En}$ and $g(\mathbf{z}) = 0$

◦ **Matrices:** $A = [A_1; A_2]$ and $B = [-I_{2En}; -I_{2En}]$

Check! Dimensions and structure of A ?

ADMM Updates

- **Augmented Lagrangian:** Dual variable $\mathbf{y} \in \mathbb{R}^{4En}$

$$L_\rho(\mathbf{x}, \mathbf{z}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^T(A\mathbf{x} + B\mathbf{z}) + \frac{\rho}{2}\|A\mathbf{x} + B\mathbf{z}\|^2$$

- **ADMM updates**

$$\mathbf{x} - \text{update} : \quad \nabla f(\mathbf{x}^{k+1}) + A^T \mathbf{y}^k + \rho A^T (A\mathbf{x}^{k+1} + B\mathbf{z}^k) = 0$$

$$\mathbf{z} - \text{update} : \quad B^T \mathbf{y}^k + \rho B^T (A\mathbf{x}^{k+1} + B\mathbf{z}^{k+1}) = 0$$

$$\mathbf{y} - \text{update} : \quad \mathbf{y}^{k+1} - \mathbf{y}^k - \rho(A\mathbf{x}^{k+1} + B\mathbf{z}^{k+1}) = 0$$

Goal: Manipulate update equations above to yield decentralized implementation of ADMM!

Manipulating ADMM Updates I

- Multiplying \mathbf{y} update by A^T and adding with the \mathbf{x} update
- Multiplying \mathbf{y} update by B^T and adding with the \mathbf{z} update
- We get

$$\nabla f(\mathbf{x}^{k+1}) + A^T \mathbf{y}^{k+1} + \rho A^T B(\mathbf{z}^k - \mathbf{z}^{k+1}) = 0$$

$$B^T \mathbf{y}^{k+1} = 0$$

$$\mathbf{y}^{k+1} - \mathbf{y}^k - \rho(A\mathbf{x}^{k+1} + B\mathbf{z}^{k+1}) = 0$$

- Take $\mathbf{y} = [\beta; \gamma]$ with $\beta, \gamma \in \mathbb{R}^{2En}$
- **Recall:** $B = [-I_{2En}; -I_{2En}]$, therefore the second equation implies

$$\beta^{k+1} = -\gamma^{k+1}$$

Manipulating ADMM Updates II

Using $\beta^{k+1} = -\gamma^{k+1}$ from previous slide:

- With **initialization** $\beta^0 = -\gamma^0$ and $z^0 = \frac{1}{2}A_+^T x^0$
- **With some manipulation (check!),** the update equations can be equivalently written as:

$$\begin{aligned}\nabla f(x^{k+1}) + A_- \beta^{k+1} - \rho A_+(z^k - z^{k+1}) &= 0 \\ \beta^{k+1} - \beta^k - \frac{\rho}{2} A_-^T x^{k+1} &= 0 \\ \frac{1}{2} A_+^T x^k - z^k &= 0\end{aligned}\tag{1.4}$$

with, $A_+ = A_1^T + A_2^T$ and $A_- = A_1^T - A_2^T$ (**Recall:**
 $A = [A_1; A_2]$)

Equation (1.4) will be used in the analysis of Decentralized ADMM!

Finally: Algorithm Updates

- With **initialization** $\beta^0 = -\gamma^0$ and $\mathbf{z}^0 = \frac{1}{2}A_+^T\mathbf{x}^0$
- **With some more manipulation (check!)**, the update equations can be further simplified to

$$\mathbf{x} - \text{update} : \nabla f(\mathbf{x}^{k+1}) + \alpha^k + 2\rho W\mathbf{x}^{k+1} - \rho L_+\mathbf{x}^k = 0$$

$$\alpha - \text{update} : \alpha^{k+1} - \alpha^k - \rho L_-\mathbf{x}^{k+1} = 0$$

$$L_+ = \frac{1}{2}A_+A_+^T, L_- = \frac{1}{2}A_-A_-^T, W = \frac{1}{2}(L_+ + L_-) \text{ and } \alpha = A_-\beta.$$

- Matrices A_+ , A_- , L_+ , L_- and W capture the **network topology**
 - A_+ and A_- : Unoriented and oriented **incidence** matrices
 - L_+ and L_- : Signless and signed **Laplacian** matrices
 - W : **Degree** matrix

Algorithm: Fully Decentralized Implementation

Algorithm updates translate to the following updates for the i th node

Just using definitions of W , L_+ and L_-

Decentralized Consensus Optimization Based on ADMM

- **Input:** Functions f_i ; initialize $\mathbf{x}_i^0 = 0$, $\alpha_i^0 = 0$, set $\rho > 0$
 - **For** $k = 0, 1, \dots$, every agent i **do**
 - Update \mathbf{x}_i^{k+1} by solving
$$\nabla f_i(\mathbf{x}_i^{k+1}) + \alpha_i^k + 2\rho|\mathcal{N}_i|\mathbf{x}_i^{k+1} - \rho\left(|\mathcal{N}_i|\mathbf{x}_i^k - \sum_{j \in \mathcal{N}_i} \mathbf{x}_j^k\right) = 0$$
 - Update $\alpha_i^{k+1} = \alpha_i^k + \rho\left(|\mathcal{N}_i|\mathbf{x}_i^{k+1} - \sum_{j \in \mathcal{N}_i} \mathbf{x}_j^{k+1}\right)$
 - **End For**
-

Assumptions

Assumption 3 (Strong Convexity)

Local objective functions are strongly convex, i.e, for each agent $i \in [m]$ given any $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$:

$$\langle \nabla f_i(\mathbf{x}_1) - \nabla f_i(\mathbf{x}_2), \mathbf{x}_1 - \mathbf{x}_2 \rangle \geq m_{f_i} \|\mathbf{x}_1 - \mathbf{x}_2\|^2 \text{ with } m_{f_i} > 0$$

Assumption 3 implies f is also strongly convex, with parameter
 $m_f = \min_i m_{f_i}$

Assumption 4 (Lipschitz Continuity)

The gradients of the local objective functions are Lipschitz continuous, i.e. for each agent $i \in [m]$ given any $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$:

$$\|\nabla f_i(\mathbf{x}_1) - \nabla f_i(\mathbf{x}_2)\| \leq M_{f_i} \|\mathbf{x}_1 - \mathbf{x}_2\| \text{ with } M_{f_i} > 0$$

Assumption 4 implies f is also Lipschitz continuous with
 $M_f = \max_i M_{f_i}$

Definitions

Definition 1.1 (Q-Linear Convergence)

A sequence \mathbf{y}^k , Q-linearly converges to a point \mathbf{y}^* if there exist a number $\sigma \in (0, 1)$ such that

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{y}^{k+1} - \mathbf{y}^*\|}{\|\mathbf{y}^k - \mathbf{y}^*\|} = \sigma$$

Main Result

Define:

$$\mathbf{u}^k = [\mathbf{z}^k; \beta^k], \mathbf{u}^* = [\mathbf{z}^*; \beta^*] \text{ and } G = [\rho I_{2En} \ 0_{2En}; 0_{2En} \ \frac{1}{\rho} I_{2En}]$$

Theorem 1.2

*Under Assumptions 3 and 4, and **proper initialization**, for any $\mu > 0$, \mathbf{u}^k is Q-linearly convergence to \mathbf{u}^* w.r.t. G-norm:*

$$\|\mathbf{u}^{k+1} - \mathbf{u}^*\|_G^2 \leq \frac{1}{1 + \delta} \|\mathbf{u}^k - \mathbf{u}^*\|_G^2$$

for $\delta > 0$, with δ dependent on the network and function parameters. Also, \mathbf{x}^k converges as:

$$\|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 \leq \frac{1}{m_f} \|\mathbf{u}^k - \mathbf{u}^*\|_G^2$$

Proper initialization: Dual variable β is in the column space of A_-^T

Proof Sketch: I

Step 1: Use update equations (1.4), KKT conditions, and strong convexity to show:

$$\|\mathbf{u}^k - \mathbf{u}^{k+1}\|_G^2 + m_f \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 \leq \|\mathbf{u}^k - \mathbf{u}^*\|_G^2 - \|\mathbf{u}^{k+1} - \mathbf{u}^*\|_G^2$$

- ❶ First, subtract update equations (1.4) and the KKT conditions
- ❷ Use strong convexity, the equations derived in (I) above, and the definition of G to reach the conclusion

Proof Sketch II

Step 2: Prove: $\|\mathbf{u}^k - \mathbf{u}^{k+1}\|_G^2 + m_f \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 \geq \delta \|\mathbf{u}^{k+1} - \mathbf{u}^*\|_G^2$

Note: Combining this with the conclusion of Step I leads to the main result of Theorem 1.2

❶ Definition of \mathbf{u} and G implies equation above is equivalent to:

$$\begin{aligned} \rho \|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2 + \frac{1}{\rho} \|\beta^{k+1} - \beta^k\|^2 + m_f \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 \\ \geq \delta \rho \|\mathbf{z}^{k+1} - \mathbf{z}^*\|^2 + \frac{\delta}{\rho} \|\beta^{k+1} - \beta^*\|^2 \end{aligned}$$

❷ Bound $\|\mathbf{z}^{k+1} - \mathbf{z}^*\|^2$: Use (1.4) & KKT conditions (**Simple!**)

❸ Bound $\|\beta^{k+1} - \beta^*\|^2$: Next!

Proof Sketch III

Step 3: Finally, upper bound: $\|\beta^{k+1} - \beta^*\|^2$

- ❶ Use Lipschitz continuity of local functions, (1.4) and KKT conditions
 - ❷ Finally, utilize initialization β **lies in the column space of A_-^T**
- Crucial Step

Conclude, Step II and combine with Step I to retrieve the statement of Theorem 1.2.

Note: R-linear convergence of sequence x^k follows from the statement of **Step I**

Takeaway from Decentralized ADMM

- A general consensus optimization problem can be **reformulated** in ADMM framework
- Global **linear convergence** for strongly convex objectives
- Convergence guarantees capture the effect of
 - **Topology** of the network
 - **Condition number** of the objective function
- **Other variations**
 - Inexact Consensus⁵
 - Asynchronous with convex objectives⁶
 - Linearize the objective⁷

⁵Chang et al., Multi-agent distributed optimization via inexact consensus ADMM, IEEE Transactions on Signal Processing 2014

⁶Wei et al., On the $O(1/K)$ convergence of asynchronous distributed alternating direction method of multipliers, IEEE GlobalSIP, 2013.

⁷Ling et al., DLM: Decentralized Linearized Alternating Direction Method of Multipliers, IEEE Transactions on Signal Processing 2015

- **Problem:**

$$\text{minimize}_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x}) \quad (1.5)$$

with $\mathbf{x} \in \mathbb{R}^n$ and $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$

- Set of m nodes
- Each node $i \in [m]$ has access to **convex** function f_i

Goal: Solve the problem in a decentralized fashion!

⁸Shi et al., Extra: An exact first-order algorithm for decentralized consensus optimization, SIAM Journal on Optimization 25.2 (2015): 944-966.

DGD vs ADMM vs EXTRA

- Recall from last lecture: **DGD** to solve (1.5)
 - **Fixed** step-size: **Neighborhood** convergence
 - **Decreasing** step-size: **Slow** convergence
- **ADMM**
 - Can be **computationally expensive**
 - Requires **successive minimizations**
- **EXTRA**: Exact First-Order Algorithm
 - **Exact** convergence with **fixed** step-size
 - **Faster** than **DGD**
 - Convex objectives: $O(1/k)$, Strongly convex objectives: **Linear**

EXTRA: Utilizes the gradient estimate of the previous iteration.

Assumptions

Assumption 5 (Mixing Matrix)

Consider a connected network $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ with $\mathcal{V} = [m]$ agents and a set of undirected edges \mathcal{E} . The mixing matrices $W = [w_{ij}] \in \mathbb{R}^{m \times m}$ and $\tilde{W} = [\tilde{w}_{ij}] \in \mathbb{R}^{m \times m}$ satisfy

- ① (Decentralized) If $i \neq j$ and $(i, j) \notin \mathcal{E}$, then $w_{ij} = 0$.
- ② (Symmetry) $W = W^T$,
- ③ (Null Space) $\text{null}\{I - W\} \supseteq \text{span}\{\mathbf{1}\}$

- Eigenvalues of W lie in $(-1, 1]$
- Slightly simplified assumption, for full details see the original paper

Assumption

Assumption 6 (Convex objective with Lipschitz continuous gradient)

Objective functions f_i are proper closed convex and Lipschitz differentiable:

$$\|\nabla f_i(\mathbf{x}_1) - \nabla f_i(\mathbf{x}_2)\| \leq M_{f_i} \|\mathbf{x}_1 - \mathbf{x}_2\| \text{ for all } \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$$

where $M_{f_i} > 0$

- Function $f(\mathbf{x}) = \sum_{i=1}^m f_i(\mathbf{x}_i)$ is M_f -Lipschitz with $M_f = \max_i \{M_{f_i}\}$

Assumption 7 (Solution Existence)

Problem (1.5) has a non-empty set of optimal solutions: $\mathcal{X}^ \neq \emptyset$*

Algorithm: EXTRA

Algorithm: EXTRA

- Choose $\alpha > 0$ and mixing matrices $W \in \mathbb{R}^{m \times m}$ and $\tilde{W} \in \mathbb{R}^{m \times m}$
 - **For** $i = 1, 2, \dots, m$
 - Pick any $\mathbf{x}_i^0 \in \mathbb{R}^n$
 - $\mathbf{x}_i^1 = \sum_{j=1}^m w_{ij} \mathbf{x}_j^0 - \alpha \nabla f_i(\mathbf{x}_i^0)$
 - **For** $k = 0, 1, \dots$ **do**
 - $$\mathbf{x}_i^{k+2} = \mathbf{x}_i^{k+1} + \sum_{j=1}^m w_{ij} \mathbf{x}_j^{k+1} - \sum_{j=1}^m (w_{ij} + 1)/2 \mathbf{x}_j^k - \alpha [\nabla f_i(\mathbf{x}_i^{k+1}) - \nabla f_i(\mathbf{x}_i^k)]$$
 - **End For**
 - **End For**
-

Algorithm: EXTRA

Stacking the iterates together, we have the following compact form

$$\mathbf{x}^{k+2} = \mathbf{x}^{k+1} + W\mathbf{x}^{k+1} - \frac{I + W}{2}\mathbf{x}^k - \alpha(\nabla f(\mathbf{x}^{k+1}) - \nabla f(\mathbf{x}^k)) \quad (1.6)$$

Convergence

- Denote $\mathbf{x} \in \mathbb{R}^{m \times n}$ and $\mathbf{x}^* \in \mathbb{R}^{m \times n}$ as a matrices with i th rows \mathbf{x}_i^T and \mathbf{x}_i^{*T} , respectively
- Introduce auxiliary sequence:

$$\mathbf{q}^k = \sum_{t=0}^k U \mathbf{x}^t \quad \text{with} \quad U = (\tilde{W} - W)^{1/2}$$

and for each k

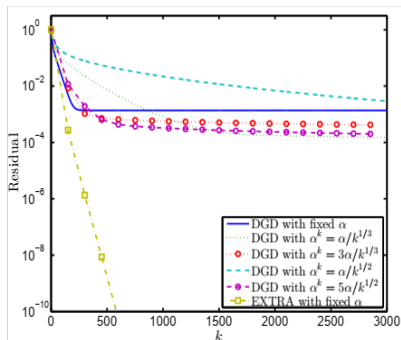
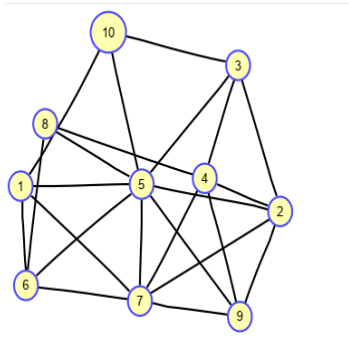
$$\mathbf{z}^k = (\mathbf{q}^k; \mathbf{x}^k), \quad \mathbf{z}^* = (\mathbf{q}^k; \mathbf{x}^*), \quad G = \begin{pmatrix} I & \mathbf{0}; \mathbf{0} & \tilde{W} \end{pmatrix}$$

Theorem 1.3

Under Assumptions 5, 6 and 7 if α satisfies $0 < \alpha < \frac{2\lambda_{\min}(\tilde{W})}{M_f}$, then

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

EXTRA vs DGD: Performance



Performance of EXTRA vs DGD for least squares consensus problem

Connection: Extra and Primal-Dual Methods^{9 10}

- **Notations:** Assume $n = 1$, denote $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_m)^T \in \mathbb{R}^m$
- $W = \mathbb{R}^{n \times n}$ mixing matrix
- **Goal** is to solve (1.5), repeated here for convenience

$$\text{minimize}_{\mathbf{x}} f(\mathbf{x}) = \sum_{i=1}^m f_i(\mathbf{x}) \quad (1.7)$$

- **EXTRA** update rule

$$\begin{aligned} \mathbf{x}^1 &= W\mathbf{x}^0 - \alpha \nabla f(\mathbf{x}^0) \\ \mathbf{x}^{k+2} &= W\mathbf{x}^{k+1} - \alpha \nabla f(\mathbf{x}^{k+1}) - (I + W)/2 \mathbf{x}^k + \alpha \nabla f(\mathbf{x}^k) \end{aligned} \quad (1.8)$$

⁹Jakovetić, A unification and generalization of exact distributed first-order methods, IEEE Transactions on Signal and Information Processing over Networks, 2018.

¹⁰Mokhtari et al., DSA: Decentralized double stochastic averaging gradient algorithm, Journal of Machine Learning Research, vol. 17, pp. 1-35, 2016.

Connection: Extra and Primal-Dual Methods

- **Equivalent Reformulation:** Let us denote $\mathcal{L} = I - \mathcal{W}$
- **Goal** is to minimize

$$\text{minimize}_{\mathbf{x}} f(\mathbf{x}) = \sum_{i=1}^m f_i(\mathbf{x}_i) \quad \text{subject to} \quad \underbrace{\frac{1}{\alpha} \mathcal{L}^{1/2} \mathbf{x} = 0}_{\text{Ensures consensus}} \quad (1.9)$$

- $\mathcal{L}^{1/2} \mathbf{x} = 0 \iff \mathbf{x}_1 = \mathbf{x}_2 = \dots = \mathbf{x}_m$ (**Check!**)

This implies that (1.7) is **equivalent** to (1.9)

Connection: Extra and Primal-Dual Methods

- **Augmented Lagrangian:** From (1.9), penalty parameter $\rho = \alpha$,

$$L_{\alpha}(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \frac{1}{\alpha} \mathbf{y}^T \mathcal{L}^{1/2} \mathbf{x} + \frac{1}{2\alpha} \mathbf{x}^T \mathcal{L} \mathbf{x}$$

- Using the notation: $\mathbf{u}^k = \frac{1}{\alpha} \mathcal{L}^{1/2} \mathbf{y}^k$, we can write the **primal-dual** update as (**Check!**)

$$\begin{aligned} \mathbf{x}^{k+1} &= \mathbf{x}^k - \alpha \left(\frac{1}{\alpha} \mathcal{L} \mathbf{x}^k + \nabla f(\mathbf{x}^k) + \mathbf{u}^k \right) \\ \mathbf{u}^{k+1} &= \mathbf{u}^k + \frac{1}{\alpha} \mathcal{L} \mathbf{x}^{k+1} \end{aligned} \tag{1.10}$$

Lemma 1.4

*The sequence $\{\mathbf{x}^k\}$ generated by **EXTRA** (1.8), with initialization $\mathbf{x}_i^0 = \mathbf{x}_j^0, \forall i, j$, is same as the sequence generated by **primal-dual** iterations (1.10) with the same initialization $\mathbf{x}_i^0 \forall i \in [m]$ and $\mathbf{u}^0 = 0$.*

Story Thus Far

- We studied **Primal-Dual Methods: ADMM**
 - Distributed/Decentralized implementations of ADMM
 - Assumptions and Convergence Guarantees
- We studied **EXTRA**
 - Limitations of DGD
 - Assumptions and Performance
- Finally, we studied **connection** between EXTRA and Primal-Dual methods
 - Specifically, we noted that EXTRA can be developed with a specific Primal-Dual construction

Next: Gradient Tracking, Push-Sum methods

Gradient Tracking¹¹

- Alternate approach to **EXTRA**
- **Goal:** To solve (1.7) in a decentralized fashion
- **Gradient Tracking Based Algorithm**
 - **Idea:** Maintain an iterative estimate of the true gradient
- Static undirected networks
- **Convex** and **Strongly** convex objectives
 - **Sublinear** convergence of $O(1/k)$ for convex objectives
 - **Linear** convergence for strongly-convex objectives

¹¹Guannan et al., Harnessing smoothness to accelerate distributed optimization, IEEE Transactions on Control of Network Systems, 2017.

Assumptions

Assumption 8 (Mixing Matrix, W)

Graph $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$, consensus weight matrix $W = [w_{ij}] \in \mathbb{R}^{m \times m}$:

- 1 For $i, j \in \mathcal{E}$, we have $w_{ij} > 0$ other wise $w_{ij} = 0$
- 2 W is doubly stochastic

Assumption 9 (Convexity and Lipschitz smoothness)

Local functions f_i for all $i \in [m]$ are convex and M_f -smooth

Assumption 10 (Strong Convexity)

Local functions f_i for all $i \in [m]$ are m_f -strongly convex

Recall: Assumptions 3, 4 and 6 used earlier

Algorithm: Gradient Tracking

Algorithm: Gradient Tracking for Decentralized Optimization

- **Initialize:** \mathbf{x}_i^0 and $\mathbf{g}_i^0 = \nabla f_i(\mathbf{x}_i^0)$
- **For** $k = 0, 1, \dots$ **do**
- Iterate update

$$\mathbf{x}_i^{k+1} = \sum_{j=1}^m w_{ij} \mathbf{x}_j^k - \alpha \mathbf{g}_i^k$$

- Descent direction update

$$\mathbf{g}_i^{k+1} = \sum_{j=1}^m w_{ij} \mathbf{g}_j^k + \nabla f_i(\mathbf{x}_i^{k+1}) - \nabla f_i(\mathbf{x}_i^k)$$

- **End For**
-

Convergence of Gradient Tracking Algorithm

Theorem 1.5

Under Assumptions 8 and 9, we have for all $i \in [m]$

$$f(\hat{\mathbf{x}}_i^{k+1}) - f^* \leq O\left(\frac{1}{k}\right)$$

where $\hat{\mathbf{x}}_i^{k+1}$ is a running average of iterates for i th agent

Choice of step-size α depends on m_f and W

Comparison to EXTRA

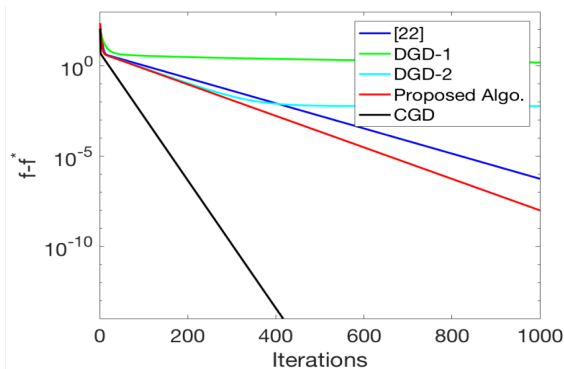
Pros:

- Easy to extended to many centralized methods
- Achieves optimality for both residuals and the objective error which is a more direct measure of optimality
 - **EXTRA** without strong convexity, achieved convergence in terms of the optimality residuals (see Theorem 1.3)

Con:

- Step-size depends on W , whereas for **EXTRA** it is independent of W

Performance: Gradient Tracking vs EXTRA and DGD



Performance of **Gradient tracking** vs **EXTRA** [22] and DGD for linear regression problem. **DGD 1** is **DGD** with fixed step-size and **DGD 2** is **DGD** with vanishing step-size; **CGD** is the centralized Gradient Descent

(Sub) Gradient Push Methods¹²

- Based on **Push-Sum** algorithm
- **Goal:** To solve (1.7) in a decentralized fashion
 - **Time-Varying** and **Directed** Networks
 - Does not assume smoothness: Only **bounded subgradients**

Assumption 11 (Uniform Strong Connectivity)

Graph $\mathcal{G}^k = \{\mathcal{V}, \mathcal{E}^k\}$

- Vertex set: \mathcal{V} with $|\mathcal{V}| = m$
- Edge set at time k : \mathcal{E}^k
- \mathcal{G}^k is **uniformly strongly connected**, i.e., $\mathcal{E}_B^k = \bigcup_{i=kB}^{(k+1)B-1} \mathcal{E}^i$ is strongly connected for every $k \geq 0$

¹²Nedić et al., Distributed optimization over time-varying directed graphs." IEEE Transactions on Automatic Control, 2014.

Algorithm

- **Notation**

- **Neighbourhoods:** $\mathcal{N}_{i,\text{in}}^k = \{j : (j, i) \in \mathcal{E}^k\} \cup \{i\}$ is the in, and $\mathcal{N}_{i,\text{out}}^k = \{j : (i, j) \in \mathcal{E}^k\} \cup \{i\}$ is the out neighbourhood
- **Out-Degree:** $d_i^k = |\mathcal{N}_{i,\text{out}}^k|$

- **Subgradient:** g_i^k is the (sub)gradient of the function f_i at x_i^{k+1}

Assumption 12 (Convex Bounded Subgradient)

- 1 Each function f_i is convex over \mathbb{R}^n and the set $\mathcal{X}^* = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ is non-empty
- 2 Subgradients g_i are bounded, i.e., $\|g_i\| \leq L_i$

Algorithm: (Sub) Gradient Push

Each node $i \in [m]$ maintains variables $\mathbf{y}_i^k, \mathbf{z}_i^k \in \mathbb{R}^n$ and scalar u_i^k

Algorithm: (Sub) Gradient Push

- **Initialize:** \mathbf{z}_i^0 and $u_i^0 = 1$ for all $i \in [m]$
- **For** $k = 0, 1, \dots$, and for all $i \in [m]$ **do**

$$\mathbf{y}_i^{k+1} = \sum_{j \in \mathcal{N}_{i,\text{in}}^k} \frac{\mathbf{z}_j^k}{d_j^k}$$

$$u_i^{k+1} = \sum_{j \in \mathcal{N}_{i,\text{in}}^k} \frac{u_j^k}{d_j^k}$$

$$\mathbf{x}_i^{k+1} = \frac{\mathbf{y}_i^{k+1}}{y_i^{k+1}}$$

$$\mathbf{z}_i^{k+1} = \mathbf{y}_i^{k+1} - \alpha^{k+1} \mathbf{g}_i^{k+1}$$

- **End For**

Convergence: (Sub) Gradient Push

Theorem 1.6

Under Assumptions 11 and 12 and with non-increasing step-sizes $\alpha^k > 0$ satisfying: $\sum_{k=1}^{\infty} \alpha^k = \infty$ and $\sum_{k=1}^{\infty} (\alpha^k)^2 < \infty$:

$$\lim_{k \rightarrow \infty} \mathbf{x}_i^k = \mathbf{x}^*$$

for some $\mathbf{x}^ \in \mathcal{X}^*$*

Theorem 1.7

With $\alpha^k = \frac{1}{\sqrt{k}}$ and each node maintaining $\tilde{\mathbf{x}}_i^{k+1} = \frac{\alpha^{k+1} \mathbf{x}_i^{k+1} + S^k \tilde{\mathbf{x}}^k}{S^{k+1}}$ for $k > 0$ with $S^0 = 0$ and $S^k = \sum_{s=1}^k \alpha^s$:

$$f(\tilde{\mathbf{x}}_i^{k+1}) - f(\mathbf{x}^*) \leq O\left(\frac{\ln k}{\sqrt{k}}\right)$$

Conclusion

- **Primal-Dual** algorithms for decentralized optimization
- **EXTRA**, **Gradient tracking** and **Push-(Sub)Gradient method**
- Many algorithms exist based on combination of similar ideas for decentralized optimization over **directed** and **dynamic** networks
 - **EXTRA** and **Push-Sum**¹³
 - **Push-Sum** and **Gradient Tracking (DIGing based methods)**¹⁴

Next: Non-convex decentralized optimization

¹³Zeng et al., Extrapush for convex smooth decentralized optimization over directed networks." arXiv preprint arXiv, 2015.

¹⁴Nedic et al., Achieving geometric convergence for distributed optimization over time-varying graphs." SIAM Journal on Optimization 2017.