# Does Alternating Direction Method of Multipliers Converge for Nonconvex Problems?

Mingyi Hong

IMSE and ECpE Department Iowa State University

ICCOPT, Tokyo, August 2016

#### The Main Content

 M. Hong, Z.-Q. Luo and M. Razaviyayn, "Convergence Analysis of Alternating Direction Method of Multipliers for a Family of Nonconvex Problems", *SIAM Journal on Optimization*, Vol. 26, No. 1, 2016 (first online Oct. 2014)

・ロ・・ (日・・ (日・・ (日・)

- The Alternating Direction Method of Multipliers (ADMM) is a very popular method for dealing with large-scale optimization problems
- Applications to classical problems
  - LP [Boyd 11], [Ye 15]
  - SDP [Wen-Goldfarb-Yin 10], [Sun-Toh-Yang 15]
  - QCQP [Huang-Sidiropoulos 16]
- Applications to emerging areas
  - Social network inference/computing [Baingana et al 15]
  - Training neural networks [Taylor et al 16]
  - Smart grid [Dall'Anese et al 13], [Peng-Low 15]
  - Bioinformatics [Forouzan-Ihler 13]

・ロン ・回 ・ ・ ヨン・

- The Alternating Direction Method of Multipliers (ADMM) is a very popular method for dealing with large-scale optimization problems
- Applications to classical problems
  - LP [Boyd 11], [Ye 15]
  - SDP [Wen-Goldfarb-Yin 10], [Sun-Toh-Yang 15]
  - QCQP [Huang-Sidiropoulos 16]
- Applications to emerging areas
  - Social network inference/computing [Baingana et al 15]
  - Training neural networks [Taylor et al 16]
  - Smart grid [Dall'Anese et al 13], [Peng-Low 15]
  - Bioinformatics [Forouzan-Ihler 13]

・ロ・・ (日・・ 日・・ 日・・

- The Alternating Direction Method of Multipliers (ADMM) is a very popular method for dealing with large-scale optimization problems
- Applications to classical problems
  - LP [Boyd 11], [Ye 15]
  - SDP [Wen-Goldfarb-Yin 10], [Sun-Toh-Yang 15]
  - QCQP [Huang-Sidiropoulos 16]
- Applications to emerging areas
  - Social network inference/computing [Baingana et al 15]
  - Training neural networks [Taylor et al 16]
  - Smart grid [Dall'Anese et al 13], [Peng-Low 15]
  - Bioinformatics [Forouzan-Ihler 13]

#### **Research Question**

• Q: Is ADMM convergent for nonconvex problems?

• A: Yes, for global consensus and sharing problems, and many more

・ロ・・ (日・・ (日・・ (日・)

#### **Research Question**

#### • Q: Is ADMM convergent for nonconvex problems?

#### • A: Yes, for global consensus and sharing problems, and many more

・ロン ・回 ・ ・ ヨン・

#### **Research Question**

• Q: Is ADMM convergent for nonconvex problems?

• A: Yes, for global consensus and sharing problems, and many more

・ロ・・ (日・・ (日・・ (日・)

### Contribution

#### Develop a new framework for analyzing the nonconvex ADMM

Obtain key insights on the behavior of the algorithm

Motivate new research in theory and applications

・ロ・・ (日・・ (日・・ (日・)

### Contribution

- Develop a new framework for analyzing the nonconvex ADMM
- Obtain key insights on the behavior of the algorithm

Motivate new research in theory and applications

・ロン ・回 ・ ・ ヨン・

### Contribution

- Develop a new framework for analyzing the nonconvex ADMM
- Obtain key insights on the behavior of the algorithm
- Motivate new research in theory and applications

・ロン ・回 ・ ・ ヨン・









#### Outline

#### Overview

#### Literature Review

- A New Analysis Framework
  - A Toy Example
  - Nonconvex Consensus Problem
  - Algorithm and Analysis
- Recent Advances

#### Conclusion

・ロト ・回ト ・ヨト ・ヨト

Consider the following problem with *K* blocks of variables  $\{x_k\}_{k=1}^{K}$ :

min 
$$f(x) := \sum_{k=1}^{K} h_k(x_k) + g(x_1, \cdots, x_K)$$
 (P)  
s.t.  $\sum_{k=1}^{K} A_k x_k = q, \ x_k \in X_k, \ \forall \ k = 1, \cdots, K$ 

•  $h_k(\cdot)$ : a convex nonsmooth function

•  $g(\cdot)$ : a smooth, possibly nonconvex function

• Ax = q: linearly coupling constraint,  $A_k \in \mathbb{R}^{M \times N_k}$ ,  $q \in \mathbb{R}^M$ 

•  $X_k \subseteq \mathbb{R}^{N_k}$ : a closed convex set

Consider the following problem with *K* blocks of variables  $\{x_k\}_{k=1}^{K}$ :

min 
$$f(x) := \sum_{k=1}^{K} h_k(x_k) + g(x_1, \cdots, x_K)$$
 (P)  
s.t.  $\sum_{k=1}^{K} A_k x_k = q, \ x_k \in X_k, \ \forall \ k = 1, \cdots, K$ 

•  $h_k(\cdot)$ : a convex nonsmooth function

•  $g(\cdot)$ : a smooth, possibly nonconvex function

• Ax = q: linearly coupling constraint,  $A_k \in \mathbb{R}^{M \times N_k}$ ,  $q \in \mathbb{R}^M$ 

•  $X_k \subseteq \mathbb{R}^{N_k}$ : a closed convex set

・ロト ・回ト ・ヨト ・ヨト … ヨ

Consider the following problem with *K* blocks of variables  $\{x_k\}_{k=1}^{K}$ :

min 
$$f(x) := \sum_{k=1}^{K} h_k(x_k) + g(x_1, \cdots, x_K)$$
 (P)  
s.t.  $\sum_{k=1}^{K} A_k x_k = q, \ x_k \in X_k, \ \forall \ k = 1, \cdots, K$ 

- $h_k(\cdot)$ : a convex nonsmooth function
- $g(\cdot)$ : a smooth, possibly nonconvex function
- Ax = q: linearly coupling constraint,  $A_k \in \mathbb{R}^{M imes N_k}$ ,  $q \in \mathbb{R}^M$
- $X_k \subseteq \mathbb{R}^{N_k}$ : a closed convex set

・ロト ・回ト ・ヨト ・ヨト … ヨ

Consider the following problem with *K* blocks of variables  $\{x_k\}_{k=1}^{K}$ :

min 
$$f(x) := \sum_{k=1}^{K} h_k(x_k) + g(x_1, \cdots, x_K)$$
 (P)  
s.t.  $\sum_{k=1}^{K} A_k x_k = q, \ x_k \in X_k, \ \forall \ k = 1, \cdots, K$ 

- $h_k(\cdot)$ : a convex nonsmooth function
- $g(\cdot)$ : a smooth, possibly nonconvex function
- Ax = q: linearly coupling constraint,  $A_k \in \mathbb{R}^{M \times N_k}$ ,  $q \in \mathbb{R}^M$
- $X_k \subseteq \mathbb{R}^{N_k}$ : a closed convex set

◆□▶ ◆□▶ ◆臣▶ ★臣▶ 臣 のへの

Consider the following problem with *K* blocks of variables  $\{x_k\}_{k=1}^{K}$ :

min 
$$f(x) := \sum_{k=1}^{K} h_k(x_k) + g(x_1, \cdots, x_K)$$
 (P)  
s.t.  $\sum_{k=1}^{K} A_k x_k = q, \ x_k \in X_k, \ \forall \ k = 1, \cdots, K$ 

- $h_k(\cdot)$ : a convex nonsmooth function
- $g(\cdot)$ : a smooth, possibly nonconvex function
- Ax = q: linearly coupling constraint,  $A_k \in \mathbb{R}^{M \times N_k}$ ,  $q \in \mathbb{R}^M$
- $X_k \subseteq \mathbb{R}^{N_k}$ : a closed convex set

• The augmented Lagrangian (AL) is given by

$$L(x;y) = \sum_{k=1}^{K} h_k(x_k) + g(x_1, \cdots, x_K) + \langle y, q - Ax \rangle + \frac{\rho}{2} ||q - Ax||^2,$$

where  $\rho > 0$  is the penalty parameter; *y* is the dual variable

(日)

- The ADMM performs a block coordinate descent (BCD) on the AL, followed by an (approximate) dual ascent
- Inexactly optimizing the AL often yields closed-form solutions

# **The ADMM Algorithm** At each iteration t + 1: Update the primal variables: $x_k^{t+1} = \arg \min_{x_k \in X_k} L(x_1^{t+1}, \dots, x_{k-1}^{t+1}, x_k, x_{k+1}^t, \dots, x_K^t; y^t), \forall k.$ Update the dual variable: $y^{t+1} = y^t + \rho(q - Ax^{t+1}).$

#### The Convex Case

- ADMM works for convex, separable, 2-block problems
- The  $g(\cdot)$  and  $h_k(\cdot)$ 's convex;  $g(x_1, \cdots, x_k) = \sum_{k=1}^{K} g_k(x_k)$ ; K = 2
- Many classic works on the analysis [Glowinski-Marroco 75], [Gabay-Mercier 76] [Glowinski 83]...
- Equivalence to Douglas-Rachford Splitting and PPA [Gabay 83], [Eckstein-Bertsekas 92]
- Convergence rates and iteration complexity analysis [Eckstein 89] [He-Yuan 12] [Deng-Yin 12] [Hong-Luo 12]
- Extension to multiple-blocks [Sun-Luo-Ye 14] [Chen et al 13] [Ma 12]

#### The Convex Case

- ADMM works for convex, separable, 2-block problems
- The  $g(\cdot)$  and  $h_k(\cdot)$ 's convex;  $g(x_1, \cdots, x_k) = \sum_{k=1}^{K} g_k(x_k)$ ; K = 2
- Many classic works on the analysis [Glowinski-Marroco 75], [Gabay-Mercier 76] [Glowinski 83]...
- Equivalence to Douglas-Rachford Splitting and PPA [Gabay 83], [Eckstein-Bertsekas 92]
- Convergence rates and iteration complexity analysis [Eckstein 89] [He-Yuan 12] [Deng-Yin 12] [Hong-Luo 12]
- Extension to multiple-blocks [Sun-Luo-Ye 14] [Chen et al 13] [Ma 12]

◆□▶ ◆□▶ ◆三▶ ◆三▶ ◆□▶ ◆□

# Solving Nonconvex Problems?

#### All the works mentioned before are for convex problems

- Recently, widely (and wildly) applied to nonconvex problems as well
  - Distributed clustering [Forero-Cano-Giannakis 11]
  - Matrix separation/completion [Xu-Yin-Wen-Zhang 11]
  - Phase retrieval [Wen-Yang-Liu-Marchesini 12]
  - Distributed matrix factorization [Ling-Yin-Wen 12]
  - Manifold optimization [Lai-Osher 12]
  - Asset allocation [Wen-Peng-Liu-Bai-Sun 13]
  - Nonnegative matrix factorization [Sun-Fevotte 14]
  - Polynomial optimization/tensor decomposition [Jiang-Ma-Zhang 13, Livavas-Sidiropoulos 14]

・ロ・・ (日・・ 日・・ 日・・

# Solving Nonconvex Problems?

- All the works mentioned before are for convex problems
- Recently, widely (and wildly) applied to nonconvex problems as well
  - Distributed clustering [Forero-Cano-Giannakis 11]
  - Matrix separation/completion [Xu-Yin-Wen-Zhang 11]
  - Phase retrieval [Wen-Yang-Liu-Marchesini 12]
  - Oistributed matrix factorization [Ling-Yin-Wen 12]
  - Manifold optimization [Lai-Osher 12]
  - Asset allocation [Wen-Peng-Liu-Bai-Sun 13]
  - Nonnegative matrix factorization [Sun-Fevotte 14]
  - Polynomial optimization/tensor decomposition [Jiang-Ma-Zhang 13, Livavas-Sidiropoulos 14]

・ロ・・ (日・・ (日・・ (日・)

### Solving Nonconvex Problems?

- All the works mentioned before are for convex problems
- Recently, widely (and wildly) applied to nonconvex problems as well
  - Distributed clustering [Forero-Cano-Giannakis 11]
  - Matrix separation/completion [Xu-Yin-Wen-Zhang 11]
  - Phase retrieval [Wen-Yang-Liu-Marchesini 12]
  - Distributed matrix factorization [Ling-Yin-Wen 12]
  - Manifold optimization [Lai-Osher 12]
  - Salaria Sun 13] Asset allocation [Wen-Peng-Liu-Bai-Sun 13]
  - Nonnegative matrix factorization [Sun-Fevotte 14]
  - Polynomial optimization/tensor decomposition [Jiang-Ma-Zhang 13, Livavas-Sidiropoulos 14]

・ロン ・四 と ・ 回 と ・ 回 と

# Application 1: Nonnegative Tensor Factorization



Figure: ADMM for solving tensor factorization problem [Liavas-Sidiropoulos 14]

- Pros: Nonconvex ADMM achieves excellent numerical performance
- Cons: A general lack of global performance analysis

#### Convergence claim

#### But this is a big "IF"!

・ロン ・回 ・ ・ ヨン・

- Pros: Nonconvex ADMM achieves excellent numerical performance
- Cons: A general lack of global performance analysis

#### Convergence claim

- IF the successive differences of all the primal and dual variables go to zero (e.g.,  $x^{t+1} x^t \rightarrow 0, y^{t+1} y^t \rightarrow 0$ )
- Provide the second s

#### But this is a big "IF"!

・ロト ・回ト ・ヨト ・ヨト … ヨ

- Pros: Nonconvex ADMM achieves excellent numerical performance
- Cons: A general lack of global performance analysis

#### Convergence claim

- IF the successive differences of all the primal and dual variables go to zero (e.g.,  $x^{t+1} x^t \rightarrow 0, y^{t+1} y^t \rightarrow 0$ )
- Provide the second s

#### But this is a big "IF"!

・ロト ・回ト ・ヨト ・ヨト … ヨ

#### • The assumption on iterates is uncheckable a priori

- "Assume" (without proving) that feasibility holds in the limit
- An exception [Zhang 10]: convergence for certain special QP
  - The AL is strongly convex
  - Only has the linear constraint
  - The dual stepsize is very small

・ロト ・回 ト ・ ヨ ト ・ ヨ ト

- The assumption on iterates is uncheckable a priori
- "Assume" (without proving) that feasibility holds in the limit
- An exception [Zhang 10]: convergence for certain special QP
  - The AL is strongly convex
  - Only has the linear constraint
  - The dual stepsize is very small

・ロン ・回 ・ ・ ヨン・

- The assumption on iterates is uncheckable a priori
- "Assume" (without proving) that feasibility holds in the limit
- An exception [Zhang 10]: convergence for certain special QP
  - The AL is strongly convex
  - Only has the linear constraint
  - The dual stepsize is very small

A D A A B A A B A A B A

#### Rigorously analyzing nonconvex ADMM is challenging

- Cases I: Without the linear constraint, reduces to the classic BCD
- Can diverge for general nonconvex  $g(\cdot)$  with  $K\geq 3$  [Powell 73]
- **Cases II:** With the linear constraint and K = 1
- Can diverge for any fixed  $\rho > 0$  [Wang-Yin-Zeng 16]

・ロ・・ (日・・ 日・・ 日・・
#### **Issues and Challenges**

Rigorously analyzing nonconvex ADMM is challenging

- Cases I: Without the linear constraint, reduces to the classic BCD
- Can diverge for general nonconvex  $g(\cdot)$  with  $K \ge 3$  [Powell 73]
- **Cases II:** With the linear constraint and K = 1
- Can diverge for any fixed  $\rho > 0$  [Wang-Yin-Zeng 16]

#### **Issues and Challenges**

Rigorously analyzing nonconvex ADMM is challenging

- Cases I: Without the linear constraint, reduces to the classic BCD
- Can diverge for general nonconvex  $g(\cdot)$  with  $K \ge 3$  [Powell 73]
- Cases II: With the linear constraint and K = 1
- Can diverge for any fixed ho > 0 [Wang-Yin-Zeng 16]

・ロト ・回ト ・ヨト ・ヨト … ヨ

#### Outline



#### Literature Review

# A New Analysis Framework A Toy Example

- Nonconvex Consensus Problem
- Algorithm and Analysis

#### Recent Advances

#### Conclusion

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

### A Toy Example

#### • First consider the following toy nonconvex example

$$\min_{x,z} \quad \frac{1}{2}x^T A x + bz, \quad \text{s.t.} \quad z \in [1,2], \ z = x$$

where A is a symmetric matrix;  $x \in \mathbb{R}^N$ 

#### ADMM Convergent?

・ロン ・回 ・ ・ ヨン・

### A Toy Example

• First consider the following toy nonconvex example

$$\min_{x,z} \quad \frac{1}{2}x^T A x + bz, \quad \text{s.t.} \quad z \in [1,2], \ z = x$$

where A is a symmetric matrix;  $x \in \mathbb{R}^N$ 

#### **ADMM Convergent?**

・ロン ・回 ・ ・ ヨン・

### A Toy Example (cont.)

- Randomly generate the data matrices A and b with N = 10
- Plot the following
  - Primal feasibility gap: ||z x||
  - 2 The optimality measure: ||x proj[x (Ax + b)]||
  - 3 The *x*-feasibility gap: ||x proj(x)||
- All three quantities go to zero iff a stationary solution has been reached

◆□▶ ◆□▶ ◆臣▶ ★臣▶ 臣 のへの

### First Try: $\rho = 20$



E

・ロト ・四ト ・ヨト ・ヨト

## Second Try: $\rho = 200$



э

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

#### A Toy Example

### A Toy Example (cont.)



E

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・ ・

### A Toy Example (cont.)

- The convergence is  $\rho$ -dependent
- When  $\rho$  is small, the algorithm fails to converge
- Different from the convex case, where any  $\rho > 0$  should work
- Reminiscent to the AL method, careful choice of  $\rho$  in nonconvex case

・ロン ・四 と ・ 回 と ・ 回 と

#### Outline



#### Literature Review

3 A New Analysis Framework

A Toy Example

Nonconvex Consensus Problem

Algorithm and Analysis

#### Recent Advances

#### Conclusion

イロト イヨト イヨト イヨト

- Consider a nonconvex global consensus problem
- A distributed optimization problem defined over a network of K agents



・ロ・・ (日・・ (日・・ (日・)

- Consider a nonconvex global consensus problem
- A distributed optimization problem defined over a network of K agents
- Formally, the problem is given by

min 
$$\sum_{k=1}^{K} g_k(x_k) + h(x_0)$$
, s.t.  $x_k = x_0, \forall k = 1, \dots, K, x_0 \in X$ .

・ロン ・回 ・ ・ ヨン・

- Consider a nonconvex global consensus problem
- A distributed optimization problem defined over a network of K agents
- Formally, the problem is given by

$$\min \sum_{k=1}^{K} g_k(x_k) + h(x_0), \quad \text{s.t.} \quad x_k = x_0, \quad \forall k = 1, \cdots, K, \ x_0 \in X.$$

・ロン ・回 ・ ・ ヨン・

- Consider a nonconvex global consensus problem
- A distributed optimization problem defined over a network of K agents
- Formally, the problem is given by

monconvex part  
min 
$$\sum_{k=1}^{K} \frac{g_k(x_k)}{g_k(x_k)} + h(x_0)$$
, s.t.  $x_k = x_0, \forall k = 1, \cdots, K, x_0 \in X$ .

- Consider a nonconvex global consensus problem
- A distributed optimization problem defined over a network of K agents
- Formally, the problem is given by

min 
$$\sum_{k=1}^{K} g_k(x_k) + \frac{h(x_0)}{h(x_0)}$$
, s.t.  $x_k = x_0, \forall k = 1, \cdots, K, x_0 \in X$ .

- Wide applications in distributed signal and information processing, parallel optimization, etc [Boyd et al 11]
- For example, in the distributed sparse PCA problem [H.-Luo-Razaviyayn 14]
  - $\bigcirc g_k(x_k) = -x_k^T A_k^T A_k x_k$ :  $A_k^T A_k$  is the covariance matrix for local data

2)  $h(\cdot)$ : some sparsity promoting nonsmooth regularizer





・ロト ・回 ト ・ヨト ・ヨト

- Wide applications in distributed signal and information processing, parallel optimization, etc [Boyd et al 11]
- For example, in the distributed sparse PCA problem [H.-Luo-Razaviyayn 14]

•  $g_k(x_k) = -x_k^T A_k^T A_k x_k$ :  $A_k^T A_k$  is the covariance matrix for local data

2  $h(\cdot)$ : some sparsity promoting nonsmooth regularizer





・ロト ・回 ト ・ヨト ・ヨト

### The Algorithm

• The AL function is given by

$$L(\{x_k\}, x_0; y) = \sum_{k=1}^{K} g_k(x_k) + h(x_0) + \sum_{k=1}^{K} \langle y_k, x_k - x_0 \rangle + \sum_{k=1}^{K} \frac{\rho_k}{2} \|x_k - x_0\|^2.$$

#### Algorithm 1. The Consensus ADMM

At each iteration t + 1, compute:

$$x_0^{t+1} = \operatorname*{argmin}_{x_0 \in X} L(\{x_k^t\}, x_0; y^t).$$

Each node k computes  $x_k$  by solving:

$$x_k^{t+1} = rg\min_{x_k} g_k(x_k) + \langle y_k^t, x_k - x_0^{t+1} \rangle + rac{
ho_k}{2} \|x_k - x_0^{t+1}\|^2.$$

Each node k updates the dual variable:

$$y_k^{t+1} = y_k^t + \rho_k \left( x_k^{t+1} - x_0^{t+1} \right)$$

(日)

#### Illustration: $x_0$ update

 $x_0$  solves:  $x_0^{t+1} = \operatorname{argmin}_{x_0 \in X} L(\{x_k^t\}, x_0; y^t)$  (often with closed-form)



#### Illustration: broadcast

#### Broadcasts the most recent $x_0$



Mingyi Hong (Iowa State University)

æ

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・ ・

### Illustration: $(x_k, \lambda_k)$ update

$$x_k$$
 solves:  $x_k^{t+1} = \arg\min_{x_k} g_k(x_k) + \langle y_k^t, x_k - x_0^{t+1} \rangle + \frac{\rho_k}{2} \|x_k - x_0^{t+1}\|^2$ .



æ

< □ > < □ > < □ > < □ > < □ > .

#### Illustration: aggregate

#### Aggregate $(x_k, y_k)$ to the central node



æ

### Main Assumptions

#### Assumption A

A1. Each  $g_k$  has Lipschitz continuous gradient:

 $\|\nabla_k g_k(x_k) - \nabla_k g_k(z_k)\| \leq \underline{L}_k \|x_k - z_k\|, \forall x_k, z_k, k = 1, \cdots, K.$ 

Moreover, h is convex (possible nonsmooth); X is a closed convex set.

A2.  $\rho_k$  is large enough such that:

If or all k, the  $x_k$  subproblem is strongly convex with modulus  $\gamma_k(\rho_k)$ ;

For all *k*, the following is satisfied

$$\rho_k > \max\left\{\frac{2L_k^2}{\gamma_k(\rho_k)}, L_k\right\}.$$

A3. f(x) is bounded from below over X.

・ロト ・回ト ・ヨト ・ヨト … ヨ

### Main Assumptions

#### Assumption A

A1. Each  $g_k$  has Lipschitz continuous gradient:

 $\|\nabla_k g_k(x_k) - \nabla_k g_k(z_k)\| \leq \underline{L}_k \|x_k - z_k\|, \forall x_k, z_k, k = 1, \cdots, K.$ 

Moreover, h is convex (possible nonsmooth); X is a closed convex set.

- A2.  $\rho_k$  is large enough such that:
  - For all *k*, the  $x_k$  subproblem is strongly convex with modulus  $\gamma_k(\rho_k)$ ;
  - 2 For all k, the following is satisfied

$$\rho_k > \max\left\{\frac{2L_k^2}{\gamma_k(\rho_k)}, L_k\right\}.$$

A3. f(x) is bounded from below over X.

◆□▶ ◆□▶ ◆臣▶ ★臣▶ 臣 のへの

### Main Assumptions

#### Assumption A

A1. Each  $g_k$  has Lipschitz continuous gradient:

 $\|\nabla_k g_k(x_k) - \nabla_k g_k(z_k)\| \leq \underline{L}_k \|x_k - z_k\|, \forall x_k, z_k, k = 1, \cdots, K.$ 

Moreover, h is convex (possible nonsmooth); X is a closed convex set.

- A2.  $\rho_k$  is large enough such that:
  - For all *k*, the  $x_k$  subproblem is strongly convex with modulus  $\gamma_k(\rho_k)$ ;
  - 2 For all *k*, the following is satisfied

$$\rho_k > \max\left\{\frac{2L_k^2}{\gamma_k(\rho_k)}, L_k\right\}.$$

A3. f(x) is bounded from below over X.

◆□▶ ◆□▶ ◆三▶ ◆三▶ ◆□▶ ◆□

- Question: Can we leverage the existing analysis for the convex case?
- $\,$  Unfortunately no, because most existing analysis relies on showing  $\|{\bf x}^t-{\bf x}^*\|^2+\|{\bf y}^t-{\bf y}^*\|^2\to 0$ 
  - where  $(\mathbf{x}^*, \mathbf{y}^*)$  are the globally optimal primal-dual pair
    - How to measure the progress of the algorithm?

・ロン ・回 と ・ 回 と

#### Question: Can we leverage the existing analysis for the convex case?

Unfortunately no, because most existing analysis relies on showing

$$\|\mathbf{x}^{t} - \mathbf{x}^{*}\|^{2} + \|\mathbf{y}^{t} - \mathbf{y}^{*}\|^{2} \to 0$$

where  $(\mathbf{x}^*, \mathbf{y}^*)$  are the globally optimal primal-dual pair

How to measure the progress of the algorithm?

・ロト ・回 ト ・ヨト ・ヨト

- Question: Can we leverage the existing analysis for the convex case?
- Unfortunately no, because most existing analysis relies on showing

$$\|\mathbf{x}^t - \mathbf{x}^*\|^2 + \|\mathbf{y}^t - \mathbf{y}^*\|^2 \to 0$$

where  $(\mathbf{x}^*, \mathbf{y}^*)$  are the globally optimal primal-dual pair

How to measure the progress of the algorithm?

・ロン ・四 と ・ 回 と ・ 回 と

- Question: Can we leverage the existing analysis for the convex case?
- Unfortunately no, because most existing analysis relies on showing

$$\|\mathbf{x}^{t} - \mathbf{x}^{*}\|^{2} + \|\mathbf{y}^{t} - \mathbf{y}^{*}\|^{2} \to 0$$

where  $(x^*, y^*)$  are the globally optimal primal-dual pair

#### How to measure the progress of the algorithm?

### Proof Ideas: A Key Step

#### • **Solution**: Use *L*(*x*; *y*) as the merit function to guide the progress

#### • Challenge: The behavior of *L*(*x*; *y*) is difficult to characterize

- Decreases after each primal update
- Increases after each dual update

#### • Technique: Bound the change of the dual update by that of the primal

・ロ・・ (日・・ (日・・ (日・)

### Proof Ideas: A Key Step

- **Solution**: Use *L*(*x*; *y*) as the merit function to guide the progress
- Challenge: The behavior of *L*(*x*; *y*) is difficult to characterize
  - Decreases after each primal update
  - Increases after each dual update

• Technique: Bound the change of the dual update by that of the primal

### Proof Ideas: A Key Step

- **Solution**: Use *L*(*x*; *y*) as the merit function to guide the progress
- Challenge: The behavior of *L*(*x*; *y*) is difficult to characterize
  - Decreases after each primal update
  - Increases after each dual update

Technique: Bound the change of the dual update by that of the primal

・ロン ・四 と ・ 回 と ・ 回 と

#### **Proof Steps**

• We develop a three-step analysis framework

Step 1: Show "sufficient descent"

$$L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^t\}, x_0^t; y^t) \\ \leq \sum_{k=1}^K \left(\frac{L_k^2}{\rho_k} - \frac{\gamma_k(\rho_k)}{2}\right) \|x_k^{t+1} - x_k^t\|^2 - \frac{\sum_{k=1}^K \rho_k}{2} \|x_0^{t+1} - x_0^t\|^2$$

Step 2: Show the following is "lower bounded"

$$L(x_0^{t+1}, \{x_k^{t+1}\}; y^{t+1}) \ge -\infty$$

Step 3: Show convergence to the set of stationary solutions

#### **Proof Steps**

- We develop a three-step analysis framework
- Step 1: Show "sufficient descent"

$$\begin{split} & L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^t\}, x_0^t; y^t) \\ & \leq \sum_{k=1}^K \left(\frac{L_k^2}{\rho_k} - \frac{\gamma_k(\rho_k)}{2}\right) \|x_k^{t+1} - x_k^t\|^2 - \frac{\sum_{k=1}^K \rho_k}{2} \|x_0^{t+1} - x_0^t\|^2 \end{split}$$

Step 2: Show the following is "lower bounded"

$$L(x_0^{t+1}, \{x_k^{t+1}\}; y^{t+1}) \ge -\infty$$

• Step 3: Show convergence to the set of stationary solutions

#### **Proof Steps**

- We develop a three-step analysis framework
- Step 1: Show "sufficient descent"

$$L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^t\}, x_0^t; y^t) \\ \leq \sum_{k=1}^K \left(\frac{L_k^2}{\rho_k} - \frac{\gamma_k(\rho_k)}{2}\right) \|x_k^{t+1} - x_k^t\|^2 - \frac{\sum_{k=1}^K \rho_k}{2} \|x_0^{t+1} - x_0^t\|^2$$

Step 2: Show the following is "lower bounded"

$$L(x_0^{t+1}, \{x_k^{t+1}\}; y^{t+1}) \geq -\infty$$

Step 3: Show convergence to the set of stationary solutions
### **Proof Steps**

- We develop a three-step analysis framework
- Step 1: Show "sufficient descent"

$$\begin{split} & L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^t\}, x_0^t; y^t) \\ & \leq \sum_{k=1}^K \left(\frac{L_k^2}{\rho_k} - \frac{\gamma_k(\rho_k)}{2}\right) \|x_k^{t+1} - x_k^t\|^2 - \frac{\sum_{k=1}^K \rho_k}{2} \|x_0^{t+1} - x_0^t\|^2 \end{split}$$

• Step 2: Show the following is "lower bounded"

$$L(x_0^{t+1}, \{x_k^{t+1}\}; y^{t+1}) \ge -\infty$$

• Step 3: Show convergence to the set of stationary solutions

### **Proof Steps**

- We develop a three-step analysis framework
- Step 1: Show "sufficient descent"

$$L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^t\}, x_0^t; y^t)$$

$$\leq \sum_{k=1}^{K} \left(\frac{L_k^2}{\rho_k} - \frac{\gamma_k(\rho_k)}{2}\right) \|x_k^{t+1} - x_k^t\|^2 - \frac{\sum_{k=1}^{K} \rho_k}{2} \|x_0^{t+1} - x_0^t\|^2$$

• Step 2: Show the following is "lower bounded"

$$L(x_0^{t+1}, \{x_k^{t+1}\}; y^{t+1}) \ge -\infty$$

• Step 3: Show convergence to the set of stationary solutions

・ロン ・四 と ・ 回 と ・ 回 と

### **Proof Steps**

- We develop a three-step analysis framework
- Step 1: Show "sufficient descent"

$$\begin{split} & L(\{x_k^{t+1}\}, x_0^{t+1}; y^{t+1}) - L(\{x_k^t\}, x_0^t; y^t) \\ & \leq \sum_{k=1}^K \left(\frac{L_k^2}{\rho_k} - \frac{\gamma_k(\rho_k)}{2}\right) \|x_k^{t+1} - x_k^t\|^2 - \frac{\sum_{k=1}^K \rho_k}{2} \|x_0^{t+1} - x_0^t\|^2 \end{split}$$

• Step 2: Show the following is "lower bounded"

$$L(x_0^{t+1}, \{x_k^{t+1}\}; y^{t+1}) \ge -\infty$$

• Step 3: Show convergence to the set of stationary solutions

#### Algorithm and Analysis

# The Convergence Claim

#### **Convergence of ADMM for Nonconvex Global Consensus**

Claim: Suppose Assumption A is satisfied. Then we have

The linear constraint is satisfied eventually:

$$\lim_{t \to \infty} \|x_k^{t+1} - x_0^{t+1}\| = 0, \ \forall \ k$$

Any limit point of the sequence generated by Algorithm 1 is a stationary solution of the consensus problem

# The Iteration Complexity Analysis

Need new gap function to measure the gap to stationarity

$$P(x^{t}, y^{t}) := \frac{\|\tilde{\nabla}L(\{x_{k}^{t}\}, x_{0}^{t}, y^{t})\|^{2}}{\|\tilde{\nabla}L(\{x_{k}^{t}\}, x_{0}^{t}, y^{t})\|^{2}} + \sum_{k=1}^{K} \frac{\|u_{k}^{t} - u_{0}^{t}\|^{2}}{\|x_{k}^{t} - x_{0}^{t}\|^{2}}$$

•  $P(x,y) = 0 \Leftrightarrow (x,y)$  is a stationary solution

**Claim:** Suppose Assumption A is satisfied,  $\epsilon > 0$  be some constant. Let  $T(\epsilon)$  denote an iteration index which satisfies

 $T(\epsilon) := \min\left\{t \mid P(x^t, y^t) \le \epsilon, t \ge 0\right\}$ 

for some  $\epsilon > 0$ . Then there exists some constant C > 0 such that

$$T(\epsilon) \leq \frac{C}{\epsilon}.$$

・ロン ・四 と ・ 回 と ・ 回 と

## The Iteration Complexity Analysis

• Need new gap function to measure the gap to stationarity

$$P(x^{t}, y^{t}) := \frac{\|\tilde{\nabla}L(\{x_{k}^{t}\}, x_{0}^{t}, y^{t})\|^{2}}{\|\tilde{\nabla}L(\{x_{k}^{t}\}, x_{0}^{t}, y^{t})\|^{2}} + \sum_{k=1}^{K} \frac{\|x_{k}^{t} - x_{0}^{t}\|^{2}}{\|x_{k}^{t} - x_{0}^{t}\|^{2}}$$

•  $P(x,y) = 0 \Leftrightarrow (x,y)$  is a stationary solution

**Claim:** Suppose Assumption A is satisfied,  $\epsilon > 0$  be some constant. Let  $T(\epsilon)$  denote an iteration index which satisfies

 $T(\epsilon) := \min\left\{t \mid P(x^t, y^t) \le \epsilon, t \ge 0\right\}$ 

for some  $\epsilon > 0$ . Then there exists some constant C > 0 such that

$$T(\epsilon) \leq \frac{C}{\epsilon}.$$

# The Iteration Complexity Analysis

Need new gap function to measure the gap to stationarity

$$P(x^{t}, y^{t}) := \frac{\|\tilde{\nabla}L(\{x_{k}^{t}\}, x_{0}^{t}, y^{t})\|^{2}}{\|\tilde{\nabla}L(\{x_{k}^{t}\}, x_{0}^{t}, y^{t})\|^{2}} + \sum_{k=1}^{K} \frac{\|x_{k}^{t} - x_{0}^{t}\|^{2}}{\|x_{k}^{t} - x_{0}^{t}\|^{2}}$$

•  $P(x,y) = 0 \Leftrightarrow (x,y)$  is a stationary solution

**Claim:** Suppose Assumption A is satisfied,  $\epsilon > 0$  be some constant. Let  $T(\epsilon)$  denote an iteration index which satisfies

$$T(\epsilon) := \min\left\{t \mid P(x^t, y^t) \le \epsilon, t \ge 0\right\}$$

for some  $\epsilon > 0$ . Then there exists some constant C > 0 such that

$$T(\epsilon) \leq rac{C}{\epsilon}.$$

- Use proximal gradient to update xk for cheap iterations
- The *x<sub>k</sub>* step is replaced by

$$\begin{aligned} x_k^{t+1} &= \arg\min_{x_k} \ \langle \nabla g_k(x_0^{t+1}), x_k - x_0^{t+1} \rangle + \langle y_k^t, x_k - x_0^{t+1} \rangle \\ &+ \frac{\rho_k + L_k}{2} \| x_k - x_0^{t+1} \|^2. \end{aligned}$$

- Gradient evaluated at the most recent *x*<sub>0</sub>!
- We can also use stochastic node sampling
- Similar convergence guarantee as Algorithm 1

・ロト ・回 ト ・ヨト ・ヨト

- Use proximal gradient to update *x<sub>k</sub>* for cheap iterations
- The xk step is replaced by

$$\begin{aligned} x_k^{t+1} &= \arg\min_{x_k} \ \langle \nabla g_k(x_0^{t+1}), x_k - x_0^{t+1} \rangle + \langle y_k^t, x_k - x_0^{t+1} \rangle \\ &+ \frac{\rho_k + L_k}{2} \| x_k - x_0^{t+1} \|^2. \end{aligned}$$

- Gradient evaluated at the most recent *x*<sub>0</sub>!
- We can also use stochastic node sampling
- Similar convergence guarantee as Algorithm 1

・ロト ・回 ト ・ヨト ・ヨト

- Use proximal gradient to update *x<sub>k</sub>* for cheap iterations
- The xk step is replaced by

$$\begin{aligned} x_k^{t+1} &= \arg\min_{x_k} \ \langle \nabla g_k(x_0^{t+1}), x_k - x_0^{t+1} \rangle + \langle y_k^t, x_k - x_0^{t+1} \rangle \\ &+ \frac{\rho_k + L_k}{2} \| x_k - x_0^{t+1} \|^2. \end{aligned}$$

- Gradient evaluated at the most recent *x*<sub>0</sub>!
- We can also use stochastic node sampling
- Similar convergence guarantee as Algorithm 1

・ロ・・ (日・・ (日・・ (日・)

- Use proximal gradient to update *x<sub>k</sub>* for cheap iterations
- The xk step is replaced by

$$\begin{aligned} x_k^{t+1} &= \arg\min_{x_k} \ \langle \nabla g_k(x_0^{t+1}), x_k - x_0^{t+1} \rangle + \langle y_k^t, x_k - x_0^{t+1} \rangle \\ &+ \frac{\rho_k + L_k}{2} \| x_k - x_0^{t+1} \|^2. \end{aligned}$$

- Gradient evaluated at the most recent *x*<sub>0</sub>!
- We can also use stochastic node sampling
- Similar convergence guarantee as Algorithm 1

・ロン ・回 ・ ・ ヨン・

- Use proximal gradient to update *x<sub>k</sub>* for cheap iterations
- The x<sub>k</sub> step is replaced by

$$\begin{aligned} x_k^{t+1} &= \arg\min_{x_k} \ \langle \nabla g_k(x_0^{t+1}), x_k - x_0^{t+1} \rangle + \langle y_k^t, x_k - x_0^{t+1} \rangle \\ &+ \frac{\rho_k + L_k}{2} \| x_k - x_0^{t+1} \|^2. \end{aligned}$$

- Gradient evaluated at the most recent *x*<sub>0</sub>!
- We can also use stochastic node sampling
- Similar convergence guarantee as Algorithm 1

イロト イポト イヨト イヨト

## The Nonconvex Sharing Problem

• Our analysis also works for the well-known sharing problem [Boyd-Parikh-Chu-Peleato-Eckstein 11]

min 
$$\sum_{k=1}^{K} h_k(x_k) + g(x_0)$$
  
s.t.  $\sum_{k=1}^{K} A_k x_k = x_0, \quad x_k \in X_k, \ k = 1, \cdots, K.$ 

- $x_k \in \mathbb{R}^{N_k}$  is the variable associated with agent k
- (K+1)-block problem, convergence unknown for the convex case
- Apply our analysis to show convergence

## The Nonconvex Sharing Problem

• Our analysis also works for the well-known sharing problem [Boyd-Parikh-Chu-Peleato-Eckstein 11]

min 
$$\sum_{k=1}^{K} h_k(x_k) + g(x_0)$$
  
s.t.  $\sum_{k=1}^{K} A_k x_k = x_0, \quad x_k \in X_k, \ k = 1, \cdots, K.$ 

- $x_k \in \mathbb{R}^{N_k}$  is the variable associated with agent k
- (K+1)-block problem, convergence unknown for the convex case
- Apply our analysis to show convergence

### Remarks

- The first analysis framework for iteration complexity of nonconvex ADMM
- A major departure from the classic analysis for convex problems
- The AL guides the convergence of the algorithm
- The  $\rho_k$ 's should be large enough, with computable lower bounds

・ロン ・回 ・ ・ ヨン・

### Outline

#### Overview

#### Literature Review

- A New Analysis Framework
  - A Toy Example
  - Nonconvex Consensus Problem
  - Algorithm and Analysis

#### Recent Advances

#### Conclusion

・ロト ・回ト ・ヨト ・ヨト

### **Recent Advances**

Many exciting recent works have been built upon our results

• New analysis, new algorithms and new connections

・ロン ・回 ・ ・ ヨン・

New analysis for weaker conditions

- Work [Li-Pong 14]: h, nonconvex, coercive; more general A<sub>k</sub>; whole sequence convergence under Kurdyka- Lojasiewicz (KL) property
- Work [Kumar et al 16]: different update schedules
- Work [Bai-Scheinberg 15]: different characterization of iteration complexity
- Works [Jiang et al 16, Wang-Yin-Zeng 16]: both relax conditions for the *K*-agent sharing problem
- Work [Yang-Pong-Chen 15]: enlarges the dual stepsize by  $\frac{\sqrt{5}+1}{2} \approx 1.618...$

• ...

◆□▶ ◆□▶ ◆三▶ ◆三▶ ◆□▶ ◆□

New analysis for weaker conditions

- Work [Li-Pong 14]: h, nonconvex, coercive; more general A<sub>k</sub>; whole sequence convergence under Kurdyka- Lojasiewicz (KL) property
- Work [Kumar et al 16]: different update schedules
- Work [Bai-Scheinberg 15]: different characterization of iteration complexity
- Works [Jiang et al 16, Wang-Yin-Zeng 16]: both relax conditions for the *K*-agent sharing problem
- Work [Yang-Pong-Chen 15]: enlarges the dual stepsize by  $\frac{\sqrt{5}+1}{2} \approx 1.618...$

◆□▶ ◆□▶ ◆臣▶ ★臣▶ 臣 のへの

New analysis for weaker conditions

- Work [Li-Pong 14]: h, nonconvex, coercive; more general A<sub>k</sub>; whole sequence convergence under Kurdyka- Lojasiewicz (KL) property
- Work [Kumar et al 16]: different update schedules
- Work [Bai-Scheinberg 15]: different characterization of iteration complexity
- Works [Jiang et al 16, Wang-Yin-Zeng 16] : both relax conditions for the *K*-agent sharing problem
- Work [Yang-Pong-Chen 15]: enlarges the dual stepsize by  $\frac{\sqrt{5}+1}{2} \approx 1.618...$

• ...

◆□▶ ◆□▶ ◆臣▶ ★臣▶ 臣 のへの

New analysis for weaker conditions

- Work [Li-Pong 14]: h, nonconvex, coercive, more general A<sub>k</sub>; whole sequence convergence under Kurdyka- Lojasiewicz (KL) property
- Work [Kumar et al 16]: different update schedules
- Work [Bai-Scheinberg 15]: different characterization of iteration complexity
- Works [Jiang et al 16, Wang-Yin-Zeng 16]: both relax conditions for the *K*-agent sharing problem
- Work [Yang-Pong-Chen 15] : enlarges the dual stepsize by  $\frac{\sqrt{5}+1}{2} \approx 1.618...$

• ...

◆□▶ ◆□▶ ◆臣▶ ★臣▶ 臣 のへの

New analysis for weaker conditions

- Work [Li-Pong 14]: h, nonconvex, coercive, A<sub>k</sub>'s full row rank; whole sequence convergence under Kurdyka- Lojasiewicz (KL) property
- Work [Kumar et al 16]: different update schedules
- Work [Bai-Scheinberg 15]: different characterization of iteration complexity
- Works [Jiang et al 16, Wang-Yin-Zeng 16]: both relax conditions for the *K*-agent sharing problem
- Work [Yang-Pong-Chen 15] : enlarges the dual stepsize by  $\frac{\sqrt{5}+1}{2} \approx 1.618...$

#### All based upon our analysis framework

◆□▶ ◆□▶ ◆三▶ ◆三▶ ◆□▶ ◆□

#### New applications in FE, SP, ML, Comm etc.

- Risk parity portfolio selection [Bai-Scheinberg 15]
- Solving certain Hamilton-Jacobi equations and differential games [Chow-Darbon-Osher-Yin-16]
- Distributed radio interference calibration [Yatawatta 16]
- Non-convex background/foreground extraction [Yang-Pong-Chen 15]
- Solving QCQP problems [Huang-Sidiropoulos 16]
- Distributed and asynchronous optimization over networks [Chang et al 16]
- Denoising using tight frame regularization [Parekh-Selesnick 15]
- Beamforming design in wireless communications [Kaleva-Tolli-Juntti 15]
- Penalized zero-variance discriminant analysis [Ames-H. 16]

• ...

◆□▶ ◆□▶ ◆三▶ ◆三▶ ◆□▶ ◆□

#### New applications in FE, SP, ML, Comm etc.

- Risk parity portfolio selection [Bai-Scheinberg 15]
- Solving certain Hamilton-Jacobi equations and differential games [Chow-Darbon-Osher-Yin-16]
- Distributed radio interference calibration [Yatawatta 16]
- Non-convex background/foreground extraction [Yang-Pong-Chen 15]
- Solving QCQP problems [Huang-Sidiropoulos 16]
- Distributed and asynchronous optimization over networks [Chang et al 16]
- Denoising using tight frame regularization [Parekh-Selesnick 15]
- Beamforming design in wireless communications [Kaleva-Tolli-Juntti 15]
- Penalized zero-variance discriminant analysis [Ames-H. 16]

o ...

#### New applications in FE, SP, ML, Comm etc.

- Risk parity portfolio selection [Bai-Scheinberg 15]
- Solving certain Hamilton-Jacobi equations and differential games [Chow-Darbon-Osher-Yin-16]
- Distributed radio interference calibration [Yatawatta 16]
- Non-convex background/foreground extraction [Yang-Pong-Chen 15]
- Solving QCQP problems [Huang-Sidiropoulos 16]
- Distributed and asynchronous optimization over networks [Chang et al 16]
- Denoising using tight frame regularization [Parekh-Selesnick 15]
- Beamforming design in wireless communications [Kaleva-Tolli-Juntti 15]
- Penalized zero-variance discriminant analysis [Ames-H. 16]

• ...

#### New applications in FE, SP, ML, Comm etc.

- Risk parity portfolio selection [Bai-Scheinberg 15]
- Solving certain Hamilton-Jacobi equations and differential games [Chow-Darbon-Osher-Yin-16]
- Distributed radio interference calibration [Yatawatta 16]
- Non-convex background/foreground extraction [Yang-Pong-Chen 15]
- Solving QCQP problems [Huang-Sidiropoulos 16]
- Distributed and asynchronous optimization over networks [Chang et al 16]
- Denoising using tight frame regularization [Parekh-Selesnick 15]
- Beamforming design in wireless communications [Kaleva-Tolli-Juntti 15]
- Penalized zero-variance discriminant analysis [Ames-H. 16]

• ...

#### New applications in FE, SP, ML, Comm etc.

- Risk parity portfolio selection [Bai-Scheinberg 15]
- Solving certain Hamilton-Jacobi equations and differential games [Chow-Darbon-Osher-Yin-16]
- Distributed radio interference calibration [Yatawatta 16]
- Non-convex background/foreground extraction [Yang-Pong-Chen 15]
- Solving QCQP problems [Huang-Sidiropoulos 16]
- Distributed and asynchronous optimization over networks [Chang et al 16]
- Denoising using tight frame regularization [Parekh-Selesnick 15]
- Beamforming design in wireless communications [Kaleva-Tolli-Juntti 15]
- Penalized zero-variance discriminant analysis [Ames-H. 16]

• ...

#### Connections of variants of ADMM with algorithm for convex problems

Nonconvex ADMM analysis

Generalize algorithm to nonconvex problems

Mingyi Hong (Iowa State University)

#### Connections of variants of ADMM with algorithm for convex problems

+

Nonconvex ADMM analysis

Generalize algorithm to nonconvex problems

Mingyi Hong (Iowa State University)

・ロン ・回 ・ ・ ヨン・

Connections of variants of ADMM with algorithm for convex problems

+

#### Nonconvex ADMM analysis

Generalize algorithm to nonconvex problems

Mingyi Hong (Iowa State University)

Connections of variants of ADMM with algorithm for convex problems

+

#### Nonconvex ADMM analysis

Generalize algorithm to nonconvex problems

Mingyi Hong (Iowa State University)

Connections of variants of ADMM with algorithm for convex problems

+

Nonconvex ADMM analysis

Generalize algorithm to nonconvex problems

Apply the Prox-ADMM to consensus over a general network [H. 16],

$$\min_{\mathbf{x}} f(\mathbf{x}) := \sum_{i=1}^{N} f_i(x_i) \quad \text{s.t.} \quad x_i = x_j \text{ if } i, j \text{ are neighbors}$$



Apply the Prox-ADMM to consensus over a general network [H. 16],

$$\min_{\mathbf{x}} f(\mathbf{x}) := \sum_{i=1}^{N} f_i(x_i) \quad \text{s.t.} \quad x_i = x_j \text{ if } i, j \text{ are neighbors}$$



The resulting algorithm is equivalent to the following primal-only iteration

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \frac{1}{2\rho} \mathbf{D}^{-1} \left( \nabla f(\mathbf{x}^t) - \nabla f(\mathbf{x}^{t-1}) \right) + \mathbf{W} \mathbf{x}^t - \frac{1}{2} (\mathbf{I} + \mathbf{W}) \mathbf{x}^{t-1}$$

where D, W are some network-related matrices

• The above iteration is precisely the EXTRA algorithm [Shi-Ling-Wu-Yin 14] for convex network consensus optimization

New Claim. EXTRA converges sublinearly for nonconvex problems

The resulting algorithm is equivalent to the following primal-only iteration

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \frac{1}{2\rho} \mathbf{D}^{-1} \left( \nabla f(\mathbf{x}^t) - \nabla f(\mathbf{x}^{t-1}) \right) + \mathbf{W} \mathbf{x}^t - \frac{1}{2} (I + W) \mathbf{x}^{t-1}$$

where D, W are some network-related matrices

 The above iteration is precisely the EXTRA algorithm [Shi-Ling-Wu-Yin 14] for convex network consensus optimization

New Claim. EXTRA converges sublinearly for nonconvex problems

・ロト ・回 ト ・ ヨト ・ ヨト
# $\mathsf{Prox}-\mathsf{ADMM} = \mathsf{EXTRA}$

The resulting algorithm is equivalent to the following primal-only iteration

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \frac{1}{2\rho} \mathbf{D}^{-1} \left( \nabla f(\mathbf{x}^t) - \nabla f(\mathbf{x}^{t-1}) \right) + \mathbf{W} \mathbf{x}^t - \frac{1}{2} (I + W) \mathbf{x}^{t-1}$$

where D, W are some network-related matrices

 The above iteration is precisely the EXTRA algorithm [Shi-Ling-Wu-Yin 14] for convex network consensus optimization

New Claim. EXTRA converges sublinearly for nonconvex problems

Consider the following convex finite sum problem:

$$\min_{x \in X} \quad f(x) := \frac{1}{N} \sum_{i=1}^{N} g_i(x),$$

where  $g_i$ ,  $i = 1, \dots N$  are cost functions; N is # of data points

- Many popular fast learning algorithms, like SAG [Le Roux-Schmidt-Bach 12], IAG [Blatt et al 07], SAGA [Defazio et al 14]:
  - In Stochastically/deterministically pick one component function  $g_i$
  - Compute its gradient
  - Update x<sup>t+1</sup> by using an average of the past gradients

Consider the following convex finite sum problem:

$$\min_{x \in X} \quad f(x) := \frac{1}{N} \sum_{i=1}^{N} g_i(x),$$

where  $g_i$ ,  $i = 1, \dots, N$  are cost functions; N is # of data points

- Many popular fast learning algorithms, like SAG [Le Roux-Schmidt-Bach 12], IAG [Blatt et al 07], SAGA [Defazio et al 14]:
  - Stochastically/deterministically pick one component function g<sub>i</sub>
  - Compute its gradient
  - Update x<sup>t+1</sup> by using an average of the past gradients

◆□▶ ◆□▶ ◆臣▶ ★臣▶ 臣 のへの



Equivalent to some variants of prox-ADMM [Hajinezhad et al 16]

New Claim. SAG/IAG/SAGA converge sublinearly for nonconvex problems



Equivalent to some variants of prox-ADMM [Hajinezhad et al 16]

New Claim. SAG/IAG/SAGA converge sublinearly for nonconvex problems



Equivalent to some variants of prox-ADMM [Hajinezhad et al 16]

New Claim. SAG/IAG/SAGA converge sublinearly for nonconvex problems

## Outline

#### Overview

#### Literature Review

- A New Analysis Framework
  - A Toy Example
  - Nonconvex Consensus Problem
  - Algorithm and Analysis

#### Recent Advances

### 5 Conclusion

・ロト ・回 ・ ・ ヨ ・ ・ ヨ ・

## Summary

- Quesiton: Whether ADMM converges for nonconvex problems?
- Yes, for a class of consensus and sharing problems, and many more

#### Key insights

- The penalty parameters are required to be large enough
- Provide a constraint of the augmented Lagrangian measures the algorithm progress

• Key technique: AL as merit function, leading to a three-step analysis

・ロン ・四 と ・ 回 と ・ 回 と

Conclusion

## Summary



# **Thank You!**

Mingyi Hong (Iowa State University)

38/38

æ

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト