# Balancing Computation and Communication in Distributed Optimization

#### Ermin Wei

#### Department of Electrical Engineering and Computer Science Northwestern University

DIMACS Rutgers Univeristy

Aug 23, 2017

- **→** → **→** 

### Collaborators



Albert S. Berahas



Raghu Bollapragada



Nitish Shirish Keskar

æ

<ロ> <同> <同> < 回> < 回>

#### Overview



- 2 Distributed Gradient Descent Variant
- 3 Communication Computation Decoupled DGD Variants
- 4 Conclusions & Future Work

#### Overview



2 Distributed Gradient Descent Variant

#### 3 Communication Computation Decoupled DGD Variants

4 Conclusions & Future Work

/⊒ > < ∃ >

### **Problem Formulation**

$$\min_{x\in\mathbb{R}^p}f(x)=\sum_{i=1}^n f_i(x)$$

• Applications: Sensor Networks, Robotic Teams, Machine Learning.



Parameter estimation in sensor networks. Communication



Multi-agent cooperative control and coordination. Battery



Large scale computation. Computation

э

(日) (同) (三) (三)

### Algorithm Evaluation

• Typical numerical results (measured in iterations or time or communication rounds).



## Algorithm Evaluation

• Typical numerical results (measured in iterations or time or communication rounds).



• Evaluation framework should reflect features of different applications.

#### **Problem Formulation**

$$\min_{x\in\mathbb{R}^p}f(x)=\sum_{i=1}^n f_i(x)$$

• Applications: Sensor Networks, Robotic Teams, Machine Learning.

- 4 同 6 4 日 6 4 日 6

### **Problem Formulation**

$$\min_{x\in\mathbb{R}^p}f(x)=\sum_{i=1}^n f_i(x)$$

- Applications: Sensor Networks, Robotic Teams, Machine Learning.
- Distributed Setting:

・ 同 ト ・ ヨ ト ・ ヨ ト

## **Problem Formulation**

$$\min_{x\in\mathbb{R}^p}f(x)=\sum_{i=1}^n f_i(x)$$

- Applications: Sensor Networks, Robotic Teams, Machine Learning.
- Distributed Setting: Consensus Optimization problem

$$\min_{x_i \in \mathbb{R}^p} f(x) = \sum_{i=1}^n f_i(x_i)$$
  
s.t.  $x_i = x_j, \quad \forall i, j \in \mathcal{N}_i$ 

- each node *i* has a local copy of the parameter vector  $x_i$
- O optimality consensus is achieved among all the nodes in the network

▲ □ ▶ < □ ▶</p>

7/38

4 3 b

### **Consensus Optimization Problem**

$$\min_{x_i \in \mathbb{R}^p} f(x) = \sum_{i=1}^n f_i(x_i)$$
  
s.t.  $x_i = x_j, \quad \forall i, j \in \mathcal{N}_i$ 

▲□ ► < □ ► </p>

э

### **Consensus Optimization Problem**

$$\min_{x_i \in \mathbb{R}^p} f(x) = \sum_{i=1}^n f_i(x_i)$$
  
s.t.  $x_i = x_j, \quad \forall i, j \in \mathcal{N}_i$ 



▲ 同 ▶ → 三 ▶

æ

∃ >

**Consensus Optimization Problem** 

$$\min_{x_i \in \mathbb{R}^p} f(x) = \sum_{i=1}^n f_i(x_i)$$
s.t.  $\mathbf{Z}\mathbf{x} = \mathbf{x}$ 

- x is a concatenation of all local x<sub>i</sub>'s
- $\bullet~{\bf W}$  is a doubly-stochastic matrix that defines the connections in the network

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{np}, \quad \mathbf{W} = \begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1n} \\ w_{21} & w_{22} & \cdots & w_{2n} \\ \vdots & \ddots & & \vdots \\ w_{n1} & w_{n2} & \cdots & w_{nn} \end{bmatrix} \in \mathbb{R}^{n \times n}, \quad \mathbf{Z} = \mathbf{W} \otimes I_p \in \mathbb{R}^{np \times np}$$

#### Literature Review

#### Sublinearly Converging Methods

- DGD [Tsitsiklis and Bertsekas, 1989; Nedić and Ozdaglar, 2009; Sundhar Ram et. al., 2010; Tsianos and Rabbat, 2012; Yuan et. al., 2015; Zeng and Yin, 2016] ...
- 2 Linearly and Superlinearly Converging Methods
  - EXTRA [Shi et. al., 2015], DIGing [Nedić et. al., 2017], NEXT [Lorenzo and Scutari, 2015], Aug-DGM [Xu et. al., 2015], NN-EXTRA [Mokhtari et. al., 2016], [Qu and Li, 2017], DQN [Eisen et. al., 2017], NN [Mokhtari et. al., 2014, 2015]...

#### Ommunication Efficient Methods

- [Chen and Ozdaglar, 2012], [Shamir et. al., 2014], [Chow et. al., 2016], [Lan et. al., 2017], [Tsianos et. al., 2012], [Zhang and Lin, 2015], ...
- Asynchronous Methods
  - [Tsitsiklis and Bertsekas, 1989], [Tsitsiklis et. al., 1986], [Sundhar Ram et. al., 2009], [Wei and Ozdaglar, 2013], [Mansoori and Wei, 2017], [Zhang and Kwok, 2014], [Wu et. al., 2017], ...

イロン 不同 とくほう イロン

#### Literature Review

#### Sublinearly Converging Methods

- DGD [Tsitsiklis and Bertsekas, 1989; Nedić and Ozdaglar, 2009; Sundhar Ram et. al., 2010; Tsianos and Rabbat, 2012; Yuan et. al., 2015; Zeng and Yin, 2016] ...
- 2 Linearly and Superlinearly Converging Methods

	EXTRA [Sh		NEXT [Lorenzo and
	Scutari, 201		[Mokhtari et. al., 2016],
	[Qu and Li,	Communication	ari et. al., 2014, 2015]
3	Communication	Computation	
	[Chen and		v et. al., 2016], [Lan et.
	al., 2017], [	Tsianos et. al., 2012], [Zhang and Lin, 2015]	,

Asynchronous Methods

• [Tsitsiklis and Bertsekas, 1989], [Tsitsiklis et. al., 1986], [Sundhar Ram et. al., 2009], [Wei and Ozdaglar, 2013], [Mansoori and Wei, 2017], [Zhang and Kwok, 2014], [Wu et. al., 2017], ...

(日)

#### Literature Review

#### Sublinearly Converging Methods

- DGD [Tsitsiklis and Bertsekas, 1989; Nedić and Ozdaglar, 2009; Sundhar Ram et. al., 2010; Tsianos and Rabbat, 2012; Yuan et. al., 2015; Zeng and Yin, 2016] ...
- 2 Linearly and Superlinearly Converging Methods

	EXTRA [Sh		NEXT [Lorenzo and
	Scutari, 201	Cost =	[Mokhtari et. al., 2016],
	[Qu and Li,	#Communications	ari et. al., 2014, 2015]
3	Communicatior	+ #Computations	
	[Chen and		v et. al., 2016], [Lan et.
	al., 2017], [ <mark>1</mark>	sianos et. al., 2012], [Zhang and Lin, 2015]	,

Asynchronous Methods

• [Tsitsiklis and Bertsekas, 1989], [Tsitsiklis et. al., 1986], [Sundhar Ram et. al., 2009], [Wei and Ozdaglar, 2013], [Mansoori and Wei, 2017], [Zhang and Kwok, 2014], [Wu et. al., 2017], ...

< ロ > < 同 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

#### Literature Review

#### Sublinearly Converging Methods

- DGD [Tsitsiklis and Bertsekas, 1989; Nedić and Ozdaglar, 2009; Sundhar Ram et. al., 2010; Tsianos and Rabbat, 2012; Yuan et. al., 2015; Zeng and Yin, 2016] ...
- ② Linearly and Superlinearly Converging Methods

	EXTRA [Sh		NEXT [Lorenzo and
	Scutari, 201	Cost =	[Mokhtari et. al., 2016],
	[Qu and Li,	#Communications×c	ari et. al., 2014, 2015]
3	Communicatior	+ #Computations $\times c_c$	
	[Chen and		v et. al., 2016], [Lan et.
	al., 2017], [	[sianos et. al., 2012], [Zhang and Lin, 2015],	

Asynchronous Methods

• [Tsitsiklis and Bertsekas, 1989], [Tsitsiklis et. al., 1986], [Sundhar Ram et. al., 2009], [Wei and Ozdaglar, 2013], [Mansoori and Wei, 2017], [Zhang and Kwok, 2014], [Wu et. al., 2017], ...

< 日 > < 同 > < 三 > < 三 >

### Goal of the Project

• Develop an algorithmic framework that is independent of the method used to balance computation and communication in distributed optimization

同 ト イ ヨ ト イ ヨ ト

### Goal of the Project

- Develop an algorithmic framework that is independent of the method used to balance computation and communication in distributed optimization
- Prove convergence for methods that use the framework

## Goal of the Project

- Develop an algorithmic framework that is independent of the method used to balance computation and communication in distributed optimization
- Prove convergence for methods that use the framework
- Show that the framework can be applied to a many consensus optimization problems (with different communication and computation costs)

## Goal of the Project

- Develop an algorithmic framework that is independent of the method used to balance computation and communication in distributed optimization
- Prove convergence for methods that use the framework
- Show that the framework can be applied to a many consensus optimization problems (with different communication and computation costs)
- Illustrate empirically that methods that utilize the framework outperform their base algorithms for specific applications

### This talk

- First stage of the project
- Multiple consensus in DGD (theoretically and in practice)
- Design a flexible first-order algorithm that decouples the two operations
- Investigate the method theoretically and empirically
- By-product: variants of DGD with exact convergence

### This talk

- First stage of the project
- Multiple consensus in DGD (theoretically and in practice)
- Design a flexible first-order algorithm that decouples the two operations
- Investigate the method theoretically and empirically
- By-product: variants of DGD with exact convergence

#### Not in this talk (ongoing work)

- Multiple gradients
- Extend framework to different algorithms (e.g., EXTRA, NN) or asynchronous methods

- 4 同 6 4 日 6 4 日 6

### Overview



2 Distributed Gradient Descent Variant

#### 3 Communication Computation Decoupled DGD Variants

4 Conclusions & Future Work

## Distributed Gradient Descent (DGD)

**DGD** [Tsitsiklis and Bertsekas, 1989; Nedić and Ozdaglar, 2009; Sundhar Ram et. al., 2010; Tsianos and Rabbat, 2012; Yuan et. al., 2015; Zeng and Yin, 2016]

$$x_{i,k+1} = \sum_{j \in \mathcal{N}_i \cup i} w_{ij} x_{j,k} - \alpha \nabla f_i(x_{i,k}), \qquad \forall i = 1, ..., n$$

・ 同 ト ・ ヨ ト ・ ヨ ト

## Distributed Gradient Descent (DGD)

**DGD** [Tsitsiklis and Bertsekas, 1989; Nedić and Ozdaglar, 2009; Sundhar Ram et. al., 2010; Tsianos and Rabbat, 2012; Yuan et. al., 2015; Zeng and Yin, 2016]

$$x_{i,k+1} = \sum_{j \in \mathcal{N}_i \cup i} w_{ij} x_{j,k} - \alpha \nabla f_i(x_{i,k}), \qquad \forall i = 1, ..., n$$

$$\mathbf{x}_{k+1} = \mathbf{Z}\mathbf{x}_k - lpha 
abla \mathbf{f}(\mathbf{x}_k)$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{np}, \quad \nabla \mathbf{f}(\mathbf{x}_k) = \begin{bmatrix} \nabla f_1(\mathbf{x}_{1,k}) \\ \nabla f_2(\mathbf{x}_{2,k}) \\ \vdots \\ \nabla f_n(\mathbf{x}_{n,k}) \end{bmatrix} \in \mathbb{R}^{np}, \quad \mathbf{Z} = \mathbf{W} \otimes I_p \in \mathbb{R}^{np \times np}$$

・ 同 ト ・ ヨ ト ・ ヨ ト

## Distributed Gradient Descent (DGD)

**DGD** [Tsitsiklis and Bertsekas, 1989; Nedić and Ozdaglar, 2009; Sundhar Ram et. al., 2010; Tsianos and Rabbat, 2012; Yuan et. al., 2015; Zeng and Yin, 2016]

$$x_{i,k+1} = \sum_{j \in \mathcal{N}_i \cup i} w_{ij} x_{j,k} - \alpha \nabla f_i(x_{i,k}), \quad \forall i = 1, ..., n$$

$$\mathbf{x}_{k+1} = \mathbf{Z}\mathbf{x}_k - \alpha \nabla \mathbf{f}(\mathbf{x}_k)$$

- Diminishing  $\alpha$ : Sub-linear convergence to the solution
- Constant  $\alpha$ : Linear convergence to a neighborhood  $\mathcal{O}(\alpha)$

#### DGD – Questions

DGD

$$\mathbf{x}_{k+1} = \mathbf{Z}\mathbf{x}_k - lpha 
abla \mathbf{f}(\mathbf{x}_k)$$

æ

## DGD – Questions

#### DGD

$$\mathbf{x}_{k+1} = \mathbf{Z}\mathbf{x}_k - \alpha \nabla \mathbf{f}(\mathbf{x}_k)$$

• Is it necessary to do one consensus step and one optimization (gradient) step?

## DGD – Questions

#### DGD

$$\mathbf{x}_{k+1} = \mathbf{Z}\mathbf{x}_k - \alpha \nabla \mathbf{f}(\mathbf{x}_k)$$

- Is it necessary to do one consensus step and one optimization (gradient) step?
- If not, what is the interpretation of methods that do more consensus/optimization steps?
- What convergence guarantees can be proven for such methods?
- How do these variants perform in practice?

## DGD – Questions

#### DGD

$$\mathbf{x}_{k+1} = \mathbf{Z}\mathbf{x}_k - \alpha \nabla \mathbf{f}(\mathbf{x}_k)$$

- Is it necessary to do one consensus step and one optimization (gradient) step?
- If not, what is the interpretation of methods that do more consensus/optimization steps?
- What convergence guarantees can be proven for such methods?
- How do these variants perform in practice?

#### <u>DGD<sup>t</sup></u>

$$\mathbf{x}_{k+1} = \mathbf{Z}^{t} \mathbf{x}_{k} - \alpha \nabla \mathbf{f}(\mathbf{x}_{k}), \qquad \mathbf{Z}^{t} = \mathbf{W}^{t} \otimes I_{\rho}$$

・ 同 ト ・ ヨ ト ・ ヨ ト

## DGD – Assumptions & Definitions

#### Assumptions

- Each component function f<sub>i</sub> is µ<sub>i</sub> > 0 strongly convex and has L<sub>i</sub> > 0 Lipschitz continuous gradients
- One mixing matrix W is symmetric and doubly stochastic with β < 1 (β is the second largest eigenvalue)</p>

#### **Definitions**

$$\bar{x}_k = \frac{1}{n} \sum_{i=1}^n x_{i,k}, \quad \nabla f(x_k) = \sum_{i=1}^n \nabla f_i(x_{i,k}), \quad \nabla f(\bar{x}_k) = \sum_{i=1}^n \nabla f_i(\bar{x}_k)$$

・ 同 ト ・ ヨ ト ・ ヨ ト

### DGD – Theory – Bounded distance to minimum

#### Theorem (Bounded distance to minimum) [Yuan et. al., 2015]

Suppose Assumptions 1 & 2 hold, and let the step length satisfy

$$lpha \leq \min\left\{rac{1+\lambda_n(\mathbf{W})}{L_f}, rac{1}{\mu_f+L_f}
ight\}$$

where  $\mu_f$  is the strong convexity parameter of f and  $L_f$  is the Lipschitz constant of the gradient of f. Then, for all k = 0, 1, ...

$$\begin{split} \|\bar{\mathbf{x}}_{k+1} - \mathbf{x}^{\star}\|^{2} &\leq c_{1}^{2} \|\bar{\mathbf{x}}_{k} - \mathbf{x}^{\star}\|^{2} + \frac{c_{3}^{2}}{(1-\beta)^{2}} \\ c_{1}^{2} &= 1 - \alpha c_{2} + \alpha \delta - \alpha^{2} \delta c_{2}, \ c_{2} &= \frac{\mu_{f} L_{f}}{\mu_{f} + L_{f}}, \\ c_{3}^{2} &= \alpha^{3} (\alpha + \delta^{-1}) L^{2} D^{2}, \ D &= \sqrt{2L(\sum_{i=1}^{n} f_{i}(0) - f^{\star})}, \end{split}$$

where  $x^* = \arg \min_x f(x)$  and  $\delta > 0$ .

- 4 同 6 4 日 6 4 日 6

### DGD – Theory – Bounded distance to minimum

#### Theorem (Bounded distance to minimum) [Yuan et. al., 2015]

Suppose Assumptions 1 & 2 hold, and let the step length satisfy

$$lpha \leq \min\left\{rac{1+\lambda_n(\mathbf{W})}{L_f},rac{1}{\mu_f+L_f}
ight\}$$

where  $\mu_f$  is the strong convexity parameter of f and  $L_f$  is the Lipschitz constant of the gradient of f. Then, for all k = 0, 1, ...

$$\begin{split} \|\bar{\mathbf{x}}_{k+1} - \mathbf{x}^{\star}\|^{2} &\leq c_{1}^{2} \|\bar{\mathbf{x}}_{k} - \mathbf{x}^{\star}\|^{2} + \frac{c_{3}^{2}}{(1-\beta)^{2}}\\ c_{1}^{2} &= 1 - \alpha c_{2} + \alpha \delta - \alpha^{2} \delta c_{2}, \ c_{2} &= \frac{\mu_{f} L_{f}}{\mu_{f} + L_{f}},\\ c_{3}^{2} &= \alpha^{3} (\alpha + \delta^{-1}) L^{2} D^{2}, \ D &= \sqrt{2L(\sum_{i=1}^{n} f_{i}(0) - f^{\star})}, \end{split}$$

where  $x^* = \arg \min_x f(x)$  and  $\delta > 0$ .

A (1) > A (2) > A

## DGD<sup>t</sup> – Theory – Bounded distance to minimum

Theorem (Bounded distance to minimum) [ASB, RB, NSK and EW, 2017]

Suppose Assumptions 1 & 2 hold, and let the step length satisfy

$$\alpha \leq \min\left\{\frac{1+\lambda_n(\mathbf{W}^t)}{L_f}, \frac{1}{\mu_f + L_f}\right\}$$

where  $\mu_f$  is the strong convexity parameter of f and  $L_f$  is the Lipschitz constant of the gradient of f. Then, for all k = 0, 1, ...

$$\begin{split} \|\bar{x}_{k+1} - x^{\star}\|^{2} &\leq c_{1}^{2} \|\bar{x}_{k} - x^{\star}\|^{2} + \frac{c_{3}^{2}}{(1 - \beta^{t})^{2}} \\ c_{1}^{2} &= 1 - \alpha c_{2} + \alpha \delta - \alpha^{2} \delta c_{2}, \ c_{2} &= \frac{\mu_{f} L_{f}}{\mu_{f} + L_{f}}, \\ c_{3}^{2} &= \alpha^{3} (\alpha + \delta^{-1}) L^{2} D^{2}, \ D &= \sqrt{2L(\sum_{i=1}^{n} f_{i}(0) - f^{\star})}, \end{split}$$

where  $x^* = \arg \min_x f(x)$  and  $\delta > 0$ .

### DGD – Theory – Comments

- Can show similar error neighborhood for  $||x_{i,k} x^*||$ .
- Theoretical results are similar to DGD in nature (constant  $\alpha$ )
  - Linear convergence to a neighborhood of the solution
  - Improved neighborhood

$$\mathcal{O}\left(rac{1}{(1-eta)^2}
ight)$$
 v.s.  $\mathcal{O}\left(rac{1}{(1-eta^{t})^2}
ight)$ 

• But, cannot kill the neighborhood with increased communication

(人間) ト く ヨ ト く ヨ ト
# DGD – Theory – Comments

- Can show similar error neighborhood for  $||x_{i,k} x^*||$ .
- Theoretical results are similar to DGD in nature (constant  $\alpha$ )
  - Linear convergence to a neighborhood of the solution
  - Improved neighborhood

$$\mathcal{O}\left(rac{1}{(1-eta)^2}
ight)$$
 v.s.  $\mathcal{O}\left(rac{1}{(1-eta^t)^2}
ight)$ 

- But, cannot kill the neighborhood with increased communication
- Drawback: requires extra communication

(人間) ト く ヨ ト く ヨ ト

# DGD – Theory – Comments

- Can show similar error neighborhood for  $||x_{i,k} x^*||$ .
- Theoretical results are similar to DGD in nature (constant  $\alpha$ )
  - Linear convergence to a neighborhood of the solution
  - Improved neighborhood

$$\mathcal{O}\left(rac{1}{(1-eta)^2}
ight)$$
 v.s.  $\mathcal{O}\left(rac{1}{(1-eta^t)^2}
ight)$ 

- But, cannot kill the neighborhood with increased communication
- Drawback: requires extra communication
- Effectively, DGD<sup>t</sup> is DGD with a different underlying graph (different weights in **W**)

- 4 同 2 4 日 2 4 日 2 4

3

# DGD – Theory – Comments

• Effectively, DGD<sup>t</sup> is DGD with a different underlying graph (different weights in **W**)

DGD – Theory – Comments

• Effectively, DGD<sup>t</sup> is DGD with a different underlying graph (different weights in **W**)



$$W = egin{bmatrix} 1/2 & 1/4 & 1/4 & 0 \ 1/4 & 1/2 & 0 & 1/4 \ 1/4 & 0 & 1/2 & 1/4 \ 0 & 1/4 & 1/4 & 1/2 \ \end{bmatrix},$$

伺 ト く ヨ ト く ヨ ト

DGD – Theory – Comments

• Effectively, DGD<sup>t</sup> is DGD with a different underlying graph (different weights in **W**)



$$W = \begin{bmatrix} 1/2 & 1/4 & 1/4 & 0\\ 1/4 & 1/2 & 0 & 1/4\\ 1/4 & 0 & 1/2 & 1/4\\ 0 & 1/4 & 1/4 & 1/2 \end{bmatrix}, \quad W^2 = \begin{bmatrix} 3/8 & 1/4 & 1/4 & \frac{1}{8}\\ 1/4 & 3/8 & 1/8 & 1/4\\ 1/4 & 1/8 & 3/8 & 1/4\\ 1/8 & 1/4 & 1/4 & 3/8 \end{bmatrix}$$

伺 ト イヨト イヨト

DGD – Theory – Comments

• Effectively, DGD<sup>t</sup> is DGD with a different underlying graph (different weights in **W**)



$$W = \begin{bmatrix} 1/2 & 1/4 & 1/4 & 0\\ 1/4 & 1/2 & 0 & 1/4\\ 1/4 & 0 & 1/2 & 1/4\\ 0 & 1/4 & 1/4 & 1/2 \end{bmatrix}, \quad W^{10} \approx \begin{bmatrix} 1/4 & 1/4 & 1/4 & 1/4\\ 1/4 & 1/4 & 1/4 & 1/4\\ 1/4 & 1/4 & 1/4 & 1/4\\ 1/4 & 1/4 & 1/4 & 1/4 \end{bmatrix}$$

・ 同 ト ・ ヨ ト ・ ヨ ト

DGD – Theory – Comments

• Effectively, DGD<sup>t</sup> is DGD with a different underlying graph (different weights in **W**)



$$W = \begin{bmatrix} 1/2 & 1/4 & 1/4 & 0\\ 1/4 & 1/2 & 0 & 1/4\\ 1/4 & 0 & 1/2 & 1/4\\ 0 & 1/4 & 1/4 & 1/2 \end{bmatrix}, \quad \mathbf{W}^{10} \approx \begin{bmatrix} 1/4 & 1/4 & 1/4 & 1/4\\ 1/4 & 1/4 & 1/4 & 1/4\\ 1/4 & 1/4 & 1/4 & 1/4\\ 1/4 & 1/4 & 1/4 & 1/4 \end{bmatrix}$$



$$W = \begin{bmatrix} 2/3 & 1/3 & 0 & 0 \\ 1/3 & 1/3 & 1/3 & 0 \\ 0 & 1/3 & 1/3 & 1/3 \\ 0 & 0 & 1/3 & 2/3 \end{bmatrix},$$

伺 ト く ヨ ト く ヨ ト

DGD – Theory – Comments

• Effectively, DGD<sup>t</sup> is DGD with a different underlying graph (different weights in **W**)

Eq. (



$$W = \begin{bmatrix} 1/2 & 1/4 & 1/4 & 0\\ 1/4 & 1/2 & 0 & 1/4\\ 1/4 & 0 & 1/2 & 1/4\\ 0 & 1/4 & 1/4 & 1/2 \end{bmatrix}, \quad W^{10} \approx \begin{bmatrix} 1/4 & 1/4 & 1/4 & 1/4\\ 1/4 & 1/4 & 1/4 & 1/4\\ 1/4 & 1/4 & 1/4 & 1/4\\ 1/4 & 1/4 & 1/4 & 1/4 \end{bmatrix}$$

$$W = \begin{bmatrix} 2/3 & 1/3 & 0 & 0\\ 1/3 & 1/3 & 1/3 & 0\\ 0 & 1/3 & 1/3 & 1/3\\ 0 & 0 & 1/3 & 2/3 \end{bmatrix}, \quad W^2 = \begin{bmatrix} 5/9 & 1/3 & 1/9 & 0\\ 1/3 & 1/3 & 2/9 & 1/9\\ 1/9 & 2/9 & 1/3 & 1/3\\ 0 & 1/9 & 1/3 & 5/9 \end{bmatrix}$$

/⊒ > < ∃ >

- ∢ ⊒ →

DGD – Theory – Comments

• Effectively, DGD<sup>t</sup> is DGD with a different underlying graph (different weights in **W**)



$$W = \begin{bmatrix} 1/2 & 1/4 & 1/4 & 0\\ 1/4 & 1/2 & 0 & 1/4\\ 1/4 & 0 & 1/2 & 1/4\\ 0 & 1/4 & 1/4 & 1/2 \end{bmatrix}, \quad W^{10} \approx \begin{bmatrix} 1/4 & 1/4 & 1/4 & 1/4\\ 1/4 & 1/4 & 1/4 & 1/4\\ 1/4 & 1/4 & 1/4 & 1/4\\ 1/4 & 1/4 & 1/4 & 1/4 \end{bmatrix}$$

$$W = \begin{bmatrix} 2/3 & 1/3 & 0 & 0 \\ 1/3 & 1/3 & 1/3 & 0 \\ 0 & 1/3 & 1/3 & 1/3 \\ 0 & 0 & 1/3 & 2/3 \end{bmatrix}, \quad W^3 = \begin{bmatrix} 0.48 & 0.33 & 0.15 & 0.04 \\ 0.33 & 0.30 & 0.22 & 0.15 \\ 0.15 & 0.22 & 0.30 & 0.33 \\ 0.04 & 0.15 & 0.33 & 0.48 \end{bmatrix}$$

伺 ト く ヨ ト く ヨ ト

DGD<sup>(1)</sup> – Numerical Results

• Problem: Quadratic

$$f(x) = \frac{1}{2} \sum_{i=1}^{n} x^{T} A_{i} x + b_{i}^{T} x$$

each node  $i = \{1, ..., n\}$  has local data  $A_i \in \mathbb{R}^{n_i imes p}$  and  $b_i \in \mathbb{R}^{n_i}$ 

- **Parameters**: n = 10, p = 10,  $n_i = 10$ ,  $\kappa = 10^2$
- Methods: DGD (1,1), DGD (1,2), DGD (1,5), DGD (1,10)

• **Graph**: 4-cyclic graph, 
$$w_{ii} = \frac{1}{5}$$
,  $w_{ij} = \begin{cases} \frac{1}{5} & \text{if } j \in \mathcal{N}_i \\ 0 & \text{otherwise} \end{cases}$ 

-

DGD<sup>(1)</sup> – Numerical Results

#### • Problem: Quadratic

$$f(x) = \frac{1}{2} \sum_{i=1}^{n} x^{T} A_{i} x + b_{i}^{T} x$$

each node  $i = \{1, ..., n\}$  has local data  $A_i \in \mathbb{R}^{n_i imes p}$  and  $b_i \in \mathbb{R}^{n_i}$ 

- **Parameters**: n = 10, p = 10,  $n_i = 10$ ,  $\kappa = 10^2$
- Methods: DGD (1,1), DGD (1,2), DGD (1,5), DGD (1,10)

• **Graph**: 4-cyclic graph, 
$$w_{ii} = \frac{1}{5}$$
,  $w_{ij} = \begin{cases} \frac{1}{5} & \text{if } j \in \mathcal{N}_i \\ 0 & \text{otherwise} \end{cases}$ 

Show the effect of multiple consensus steps per gradient step

- 4 同 2 4 日 2 4 H

# DGD<sup>(†)</sup> – Numerical Results





 $Cost = #Communications \times 1 + #Computations \times 1$ 

(日)

# Overview



2 Distributed Gradient Descent Variant

#### 3 Communication Computation Decoupled DGD Variants

#### 4 Conclusions & Future Work

→ < ∃→

# Operators

## DGD

$$\mathbf{x}_{k+1} = \mathbf{Z}\mathbf{x}_k - lpha 
abla \mathbf{f}(\mathbf{x}_k)$$

æ

# Operators

## DGD

$$\mathbf{x}_{k+1} = \mathbf{Z}\mathbf{x}_k - lpha 
abla \mathbf{f}(\mathbf{x}_k)$$

#### Operators

- $\bullet \ \mathcal{W}[\mathbf{x}] = \mathbf{Z}\mathbf{x}$
- $\mathcal{T}[\mathbf{x}] = \mathbf{x} \alpha \nabla \mathbf{f}(\mathbf{x})$

# Operators

## DGD

$$\mathbf{x}_{k+1} = \mathbf{Z}\mathbf{x}_k - \alpha \nabla \mathbf{f}(\mathbf{x}_k)$$

#### Operators

•  $\mathcal{W}[\mathbf{x}] = \mathbf{Z}\mathbf{x}$ 

• 
$$\mathcal{T}[\mathbf{x}] = \mathbf{x} - \alpha \nabla \mathbf{f}(\mathbf{x})$$

#### Methods

• DGD: 
$$(\mathcal{T} - I + \mathcal{W})[\mathbf{x}_k] = \mathbf{Z}\mathbf{x}_k - \alpha \nabla \mathbf{f}(\mathbf{x}_k)$$

æ

# Operators

## DGD

$$\mathbf{x}_{k+1} = \mathbf{Z}\mathbf{x}_k - \alpha \nabla \mathbf{f}(\mathbf{x}_k)$$

#### Operators

•  $\mathcal{W}[\mathbf{x}] = \mathbf{Z}\mathbf{x}$ 

• 
$$\mathcal{T}[\mathbf{x}] = \mathbf{x} - \alpha \nabla \mathbf{f}(\mathbf{x})$$

#### <u>Methods</u>

- DGD:  $(\mathcal{T} I + \mathcal{W})[\mathbf{x}_k] = \mathbf{Z}\mathbf{x}_k \alpha \nabla \mathbf{f}(\mathbf{x}_k)$
- TW:  $\mathcal{T}[\mathcal{W}[\mathbf{x}_k]] = \mathbf{Z}\mathbf{x}_k \alpha \nabla \mathbf{f}(\mathbf{Z}\mathbf{x}_k)$
- WT:  $\mathcal{W}[\mathcal{T}[\mathbf{x}_k]] = \mathbf{Z}\mathbf{x}_k \alpha \mathbf{Z} \nabla \mathbf{f}(\mathbf{x}_k)$

A special case of the algorithm appeared as CTA (Combined then Adapt) and ATC in [Sayed, 13] for quadratic problems

3

# Operators

## DGD

$$\mathbf{x}_{k+1} = \mathbf{Z}\mathbf{x}_k - \alpha \nabla \mathbf{f}(\mathbf{x}_k)$$

#### Operators

- $\bullet \ \mathcal{W}[\textbf{x}] = \textbf{Z}\textbf{x}$
- $\mathcal{T}[\mathbf{x}] = \mathbf{x} \alpha \nabla \mathbf{f}(\mathbf{x})$

#### <u>Methods</u>

- DGD:  $(\mathcal{T} I + \mathcal{W})[\mathbf{x}_k] = \mathbf{Z}\mathbf{x}_k \alpha \nabla \mathbf{f}(\mathbf{x}_k)$
- TW:  $\mathcal{T}[\mathcal{W}[\mathbf{x}_k]] = \mathbf{Z}\mathbf{x}_k \alpha \nabla \mathbf{f}(\mathbf{Z}\mathbf{x}_k)$
- WT:  $\mathcal{W}[\mathcal{T}[\mathbf{x}_k]] = \mathbf{Z}\mathbf{x}_k \alpha \mathbf{Z} \nabla \mathbf{f}(\mathbf{x}_k)$

向 ト イヨ ト イヨ ト

3

# Operators

#### TW

$$\mathbf{y}_k = \mathbf{Z}\mathbf{x}_k$$
 $\mathbf{x}_{k+1} = \mathbf{y}_k - lpha 
abla \mathbf{f}(\mathbf{y}_k)$ 

#### <u>Methods</u>

- DGD:  $(\mathcal{T} I + \mathcal{W})[\mathbf{x}_k] = \mathbf{Z}\mathbf{x}_k \alpha \nabla \mathbf{f}(\mathbf{x}_k)$
- TW:  $\mathcal{T}[\mathcal{W}[\mathbf{x}_k]] = \mathbf{Z}\mathbf{x}_k \alpha \nabla \mathbf{f}(\mathbf{Z}\mathbf{x}_k)$
- WT:  $\mathcal{W}[\mathcal{T}[\mathbf{x}_k]] = \mathbf{Z}\mathbf{x}_k \alpha \mathbf{Z} \nabla \mathbf{f}(\mathbf{x}_k)$

伺 ト く ヨ ト く ヨ ト

# TW – Methods, Assumptions & Definitions

#### <u>Methods</u>

- TW<sup>t</sup>: t (predetermined) consensus steps for every gradient step
- $\bullet~TW^+\colon$  increasing number of consensus steps

#### Assumptions

# • Assumptions 1 & 2, same as before <u>Definitions</u>

$$\bar{x}_k = \frac{1}{n} \sum_{i=1}^n x_{i,k}, \quad \nabla f(y_k) = \sum_{i=1}^n \nabla f_i(y_{i,k}), \quad \nabla f(\bar{x}_k) = \sum_{i=1}^n \nabla f_i(\bar{x}_k)$$

・ 同 ト ・ ヨ ト ・ ヨ ト

# TW<sup>t</sup> – Theory – Bounded distance to minimum

#### Theorem (Bounded distance to minimum) [ASB, RB, NSK and EW, 2017]

Suppose Assumptions 1 & 2 hold, and let the step length satisfy

$$lpha \leq \min\left\{rac{1+\lambda_n(\mathbf{W}^t)}{L_f},rac{1}{\mu_f+L_f}
ight\}$$

where  $\mu_f$  is the strong convexity parameter of f and  $L_f$  is the Lipschitz constant of the gradient of f. Then, for all k = 0, 1, ...

$$\|x_{i,k} - x^{\star}\| \le c_i^k \|x^{\star}\| + \frac{c_3}{\sqrt{1 - c_1^2}} \beta^t + \beta^t \alpha D$$

where  $x^* = \arg \min_x f(x)$  and  $\delta > 0$ .

- ・ 同 ト ・ ヨ ト - - ヨ

# TW Theory (increasing Consensus)

• Can we increase the number of consensus steps and converge to the solution?

同 ト イ ヨ ト イ ヨ ト

# TW Theory (increasing Consensus)

- Can we increase the number of consensus steps and converge to the solution?
- Increase t(k) accordingly so that we kill the error term

 $O(\beta^{t(k)})$ 

(人間) ト く ヨ ト く ヨ ト

# TW Theory (increasing Consensus)

- Can we increase the number of consensus steps and converge to the solution?
- Increase t(k) accordingly so that we kill the error term

 $O(\beta^{t(k)})$ 

• Similar idea appeared in [Chen and Ozdaglar, 2012] for nonsmooth problems

伺 ト イ ヨ ト イ ヨ ト

# TW Theory (increasing Consensus)

- Can we increase the number of consensus steps and converge to the solution?
- Increase t(k) accordingly so that we kill the error term

 $O(\beta^{t(k)})$ 

- Similar idea appeared in [Chen and Ozdaglar, 2012] for nonsmooth problems
- Resulting in TW<sup>t(k)</sup> algorithm with exact convergence: As long as we keep increasing the number of consensus

(4月) イヨト イヨト

# TW<sup>+</sup> – Theory – Bounded distance to minimum

#### Theorem (Bounded distance to minimum)[ASB, RB, NSK and EW, 2017]

Suppose Assumptions 1 & 2 hold, t(k) = k and let the step length satisfy

$$\alpha \leq \min\left\{\frac{1}{L_f}, \frac{1}{\mu_{\overline{f}} + L_{\overline{f}}}\right\}$$

where  $\mu_f$  is the strong convexity parameter of f and  $L_f$  is the Lipschitz constant of the gradient of f. Then, for all k = 0, 1, ...

$$\|x_{i,k}-x^{\star}\|\leq C\rho^{k},$$

where  $x^* = \arg \min_x f(x)$  and some constants C,  $\rho$ .

When t(k) = k to reach an  $\epsilon$ -accurate solution, we need  $O(\log(\frac{1}{\epsilon}))$  number of gradient evaluation and  $O((\log(\frac{1}{\epsilon}))^2)$  rounds of communication.

イロト 不得 とくほ とくほ とうほう

## Numerical Experiments

 Methods: DGD, TW (1,1,-), TW (10,1,-), TW (1,10,-), TW (1,1,k), TW (1,1,500), TW (1,1,1000)

同 ト イ ヨ ト イ ヨ ト

## Numerical Experiments

- Methods: DGD, TW (1,1,-), TW (10,1,-), TW (1,10,-), TW (1,1,k), TW (1,1,500), TW (1,1,1000)
- Problem: Quadratic

$$f(x) = \frac{1}{2} \sum_{i=1}^{n} x^{T} A_{i} x + b_{i}^{T} x$$

each node  $i = \{1, ..., n\}$  has local data  $A_i \in \mathbb{R}^{n_i imes p}$  and  $b_i \in \mathbb{R}^{n_i}$ 

- **Parameters**: n = 10, p = 10,  $n_i = 10$ ,  $\kappa = \frac{L}{\mu} = 10^4$
- Graph: 4-cyclic graph

- \* 同 \* \* ヨ \* \* ヨ \* - ヨ

## Numerical Experiments – Quadratic Problems



Quadratic. n = 10, p = 10,  $n_i = 10$ ,  $\kappa = 10^4$ .

 $Cost = #Communications \times 1 + #Computations \times 1$ 

э

## Numerical Experiments – Quadratic Problems



Quadratic. n = 10, p = 10,  $n_i = 10$ ,  $\kappa = 10^4$ .

 $Cost = #Communications \times 1 + #Computations \times 1$ 

э

-

## Numerical Experiments – Quadratic Problems



Quadratic. n = 10, p = 10,  $n_i = 10$ ,  $\kappa = 10^4$ .

 $Cost = #Communications \times 1 + #Computations \times 1$ 

э

∃ >

## Numerical Experiments – Quadratic Problems



Quadratic. n = 10, p = 10,  $n_i = 10$ ,  $\kappa = 10^4$ .

 $Cost = #Communications \times 1 + #Computations \times 1$ 

э

## Numerical Experiments – Quadratic Problems



Quadratic. n = 10, p = 10,  $n_i = 10$ ,  $\kappa = 10^4$ .

 $Cost = #Communications \times 1 + #Computations \times 1$ 

## Numerical Experiments – Quadratic Problems



Quadratic. n = 10, p = 10,  $n_i = 10$ ,  $\kappa = 10^4$ .

 $Cost = #Communications \times 1 + #Computations \times 1$ 

## Numerical Experiments – Quadratic Problems



Quadratic. n = 10, p = 10,  $n_i = 10$ ,  $\kappa = 10^4$ .

 $Cost = #Communications \times 1 + #Computations \times 1$ 

## Experiments – Quadratic Problems – Different Costs



Quadratic. n = 10, p = 10,  $n_i = 10$ ,  $\kappa = 10^4$ .

**Center**:  $c_g = 1$ ,  $c_c = 1$ ;

 $Cost = \#Communications \times c_{c} + \#Computations \times c_{g}$ 

(日)

3
#### Experiments – Quadratic Problems – Different Costs



 $Cost = \#Communications \times c_c + \#Computations \times c_{\sigma}$ 

3

#### Experiments – Quadratic Problems – Different Costs



 $Cost = \#Communications \times c_{c} + \#Computations \times c_{g}$ 

イロト イポト イヨト イヨト

3

## Numerical Experiments – Logistic Regression

• **Problem**: Logistic Regression - Binary Classification (Mushroom Dataset)

$$f(x) = \frac{1}{n \cdot n_i} \sum_{i=1}^n \sum_{j=1}^{n_i} \log(1 + e^{-(b_i)_j(x^T(A_i)_{j\cdot})})$$

where  $A \in \mathbb{R}^{n \cdot n_i \times p}$  and  $b \in \{-1, 1\}^{n \cdot n_i}$ , and each node i = 1, ..., n has a portion of A and b,  $A_i \in \mathbb{R}^{n_i \times p}$  and  $b_i \in \mathbb{R}^{n_i}$ 

- Parameters: n = 10, p = 114,  $n_i = 812$ ,  $\kappa = 10^4$
- Graph: 4-cyclic graph

A (1) A (2) A (2) A

## Numerical Experiments – Logistic Regression



Logistic Regression - mushroom. n = 10. p = 114.  $n_i = 812$ .

 $Cost = #Communications \times 1 + #Computations \times 1$ 

э

# Overview



2 Distributed Gradient Descent Variant

#### 3 Communication Computation Decoupled DGD Variants



## **Final Remarks**

 Most distributed optimization algorithms do one communication and one computation per iteration

- **→** → **→** 

э

# **Final Remarks**

- Most distributed optimization algorithms do one communication and one computation per iteration
- Showed the effect (theoretically and empirically) of doing multiple consensus in DGD

/₽ ► < ∃ ►

э

# **Final Remarks**

- Most distributed optimization algorithms do one communication and one computation per iteration
- Showed the effect (theoretically and empirically) of doing multiple consensus in DGD
- Proposed a variant of DGD, TW, that decouples the two operations (consensus and computation) and converges to the solution by performing multiple consensus steps

# **Final Remarks**

- Most distributed optimization algorithms do one communication and one computation per iteration
- Showed the effect (theoretically and empirically) of doing multiple consensus in DGD
- Proposed a variant of DGD, TW, that decouples the two operations (consensus and computation) and converges to the solution by performing multiple consensus steps

DGD: 
$$x_{k+1} = Zx_k - \alpha \nabla f(x_k)$$
  
TW:  $x_{k+1} = Zx_k - \alpha \nabla f(Zx_k)$ 

# **Final Remarks**

- Most distributed optimization algorithms do one communication and one computation per iteration
- Showed the effect (theoretically and empirically) of doing multiple consensus in DGD
- Proposed a variant of DGD, TW, that decouples the two operations (consensus and computation) and converges to the solution by performing multiple consensus steps
- Important to balance communication and computation in order to get best performance in terms of cost — right balance depends on the application (e.g., cost of communication and cost of computation)

・ロト ・得ト ・ヨト ・ヨト

# Future Work

- Apply framework to other algorithms (exact, second-order, asynchronous, ...)
- Construct framework to do multiple gradient steps
- Adapt number of gradient and communication steps in algorithmic way
- Other considerations: memory access, partial blocks, quantization effects, dynamic environment

# Backup Slides

▶ ∢ ≣ ▶

# DGD – Theory – Bounded gradients

#### Lemma (Bounded gradients) [Yuan et. al., 2015]

Suppose Assumption 1 holds, and let the step size satisfy

$$lpha \leq rac{1 + \lambda_n(\mathbf{W})}{L}$$

where  $\lambda_n(\mathbf{W})$  is the smallest eigenvalue of  $\mathbf{W}$  and  $L = \max_i L_i$ . Then, starting from  $x_{i,0} = 0$  (i = 1, 2, ..., n), the sequence  $x_{i,k}$  generated by DGD converges. In addition, we also have

$$\|\nabla \mathbf{f}(\mathbf{x}_k)\| \leq D = \sqrt{2L(\frac{1}{n}\sum_{i=1}^n f_i(0) - f_i^{\star})}$$
(1)

for all k = 1, 2, ..., where  $f_i^* = f_i(x_i^*)$  and  $x_i^* = \arg \min_x f_i(x)$ .

# DGD<sup>t</sup> – Theory – Bounded gradients

#### Lemma (Bounded gradients) [ASB, RB, NSK and EW, 2017]

Suppose Assumption 1 holds, and let the step size satisfy

$$\alpha \leq \frac{1 + \lambda_n(\mathbf{W}^t)}{L}$$

where  $\lambda_n(\mathbf{W}^t)$  is the smallest eigenvalue of  $\mathbf{W}^t$  and  $L = \max_i L_i$ . Then, starting from  $x_{i,0} = 0$  (i = 1, 2, ..., n), the sequence  $x_{i,k}$  generated by DGD converges. In addition, we also have

$$\|\nabla \mathbf{f}(\mathbf{x}_k)\| \le D = \sqrt{2L(\frac{1}{n}\sum_{i=1}^n f_i(0) - f_i^{\star})}$$
 (1)

for all k = 1, 2, ..., where  $f_i^* = f_i(x_i^*)$  and  $x_i^* = \arg \min_x f_i(x)$ .

イロト 不得 トイヨト イヨト 二日

# DGD – Theory – Bounded deviation from mean

#### Lemma (Bounded deviation from mean) [Yuan et. al., 2015]

If (1) and Assumption 2 hold, then the total deviation from the mean is bounded, namely,

$$\|x_{i,k} - \bar{x}_k\| \le \frac{\alpha D}{1 - \beta}$$

for all k and i. Moreover, if in addition Assumption 1 holds, then

$$\begin{split} \|\nabla f_i(x_{i,k}) - \nabla f_i(\bar{x}_k)\| &\leq \frac{\alpha DL_i}{1-\beta} \\ \|\nabla f(x_k) - \nabla f(\bar{x}_k)\| &\leq \frac{\alpha DL}{1-\beta} \end{split}$$

for all k and i.

э

# DGD<sup>t</sup> – Theory – Bounded deviation from mean

#### Lemma (Bounded deviation from mean) [ASB, RB, NSK and EW, 2017]

If (1) and Assumption 2 hold, then the total deviation from the mean is bounded, namely,

$$\|x_{i,k} - \bar{x}_k\| \le \frac{\alpha D}{1 - \beta^t}$$

for all k and i. Moreover, if in addition Assumption 1 holds, then

$$egin{aligned} \|
abla f_i(x_{i,k}) - 
abla f_i(ar{x}_k)\| &\leq rac{lpha DL_i}{1-eta^t} \ \|
abla f(x_k) - 
abla f(ar{x}_k)\| &\leq rac{lpha DL}{1-eta^t} \end{aligned}$$

for all k and i.

- 4 同 6 4 日 6 4 日 6

# TW<sup>t</sup> – Theory – Bounded deviation from mean

#### Lemma (Bounded deviation from mean) [ASB, RB, NSK and EW, 2017]

If Assumption 2 holds, then the total deviation from the mean is bounded, namely,

$$\|y_{i,k}-\bar{x}_k\|\leq rac{eta^tlpha Dk(k+1)}{2}$$

for all k and i. Moreover, if in addition Assumption 1 holds, then

$$egin{aligned} \|
abla f_i(y_{i,k}) - 
abla f_i(ar{x}_k)\| &\leq rac{eta^t lpha DL_i k(k+1)}{2} \ \|
abla f(y_k) - 
abla f(ar{x}_k)\| &\leq rac{eta^t lpha DL k(k+1)}{2} \end{aligned}$$

for all k and i.

(人間) ト く ヨ ト く ヨ ト

## DGD<sup>t</sup> Numerical Results



Quadratic. n = 10, p = 10,  $n_i = 10$ ,  $\kappa = 10^2$ .

・ロト ・回ト ・ヨト ・ヨト

≣ 36/38

#### Experiments – Quadratic Problems



Quadratic.  $n = 10, d = 10, n_i = 10, \kappa = 10^4$ .

・ロン ・部 と ・ ヨ と ・ ヨ と …

#### Experiments – Quadratic Problems



Quadratic. n = 10, d = 10,  $n_i = 10$ ,  $\kappa = 10^4$ .

・ロン ・部 と ・ ヨ と ・ ヨ と …

#### Experiments – Quadratic Problems



#### Experiments – Quadratic Problems



Quadratic. n = 10, d = 10,  $n_i = 10$ ,  $\kappa = 10^4$ .

#### Experiments – Quadratic Problems



Quadratic. n = 10, d = 10,  $n_i = 10$ ,  $\kappa = 10^4$ .

・ロン ・部 と ・ ヨ と ・ ヨ と …

#### Experiments – Quadratic Problems



Quadratic. n = 10, d = 10,  $n_i = 10$ ,  $\kappa = 10^4$ .

・ロン ・部 と ・ ヨ と ・ ヨ と …

# Experiments – Quadratic Problems



Quadratic. n = 10, d = 10,  $n_i = 10$ ,  $\kappa = 10^4$ .

・ロン ・部 と ・ ヨ と ・ ヨ と …

#### Experiments – Logistic Regression



Logistic Regression - mushroom.  $n = 10, d = 114, n_i = 812.$ 

・ロン ・部 と ・ ヨ と ・ ヨ と …