

# Distributed intelligence in multi-agent systems

Usman Khan

Department of Electrical and Computer Engineering

Tufts University

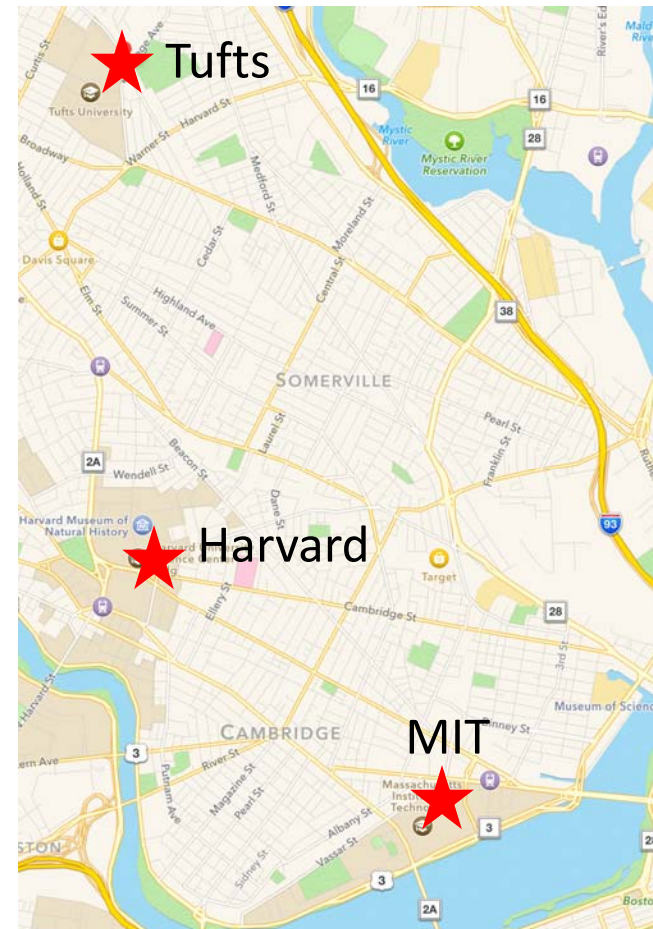
Workshop on Distributed Optimization, Information Processing, and Learning

Rutgers University

August 21, 2017

# Who am I

- **Usman A. Khan**
  - Associate Professor, Tufts
- **Postdoc**
  - U-Penn
- **Education**
  - PhD, Carnegie Mellon
  - MS, UW-Madison
  - BS, Pakistan



# My Research Lab: Projects and demos

## Research Team

### PhD Students

- Current**
- Ran Yin, Sep. 2016 to date
- Fakhteh Sadjadnaki, Sep. 2014 to date
- Sam Safavi, Jan. 2013 to date

### Alumni

- Chengjiao
- Distributed
- Mohammed
- Distributed

### MS Students

#### Current

- Noel Hwa
- Diane Lisk

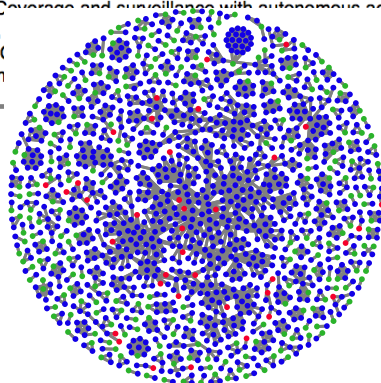
#### Alumni

- Steven Sa
- Christoph
- Dong Park

- Alexander Henry, Aug. 2015, Adaptive methods for robotic path planning
- Michael Tran, Aug. 2015, Distributed target tracking in a sensor network
- Dibeyandu Das, Aug. 2014, Consensus with non-replicating agents
- Anders Simpson-Wolf, Dec. 2014, Privacy and differentially private methods
- Luke Grymek, Aug. 2013, Coverage and surveillance with autonomous agents
- Gerald Solimini, May 2012,
- Syed S. Akbar, Dec. 2011, (
- Qiong Wu (Applied Mathem



## Inference in Social Networks



wind turbines

### Undergraduates

#### Current

- Mathias Barth, Class of 2019, Jun. 2017 to date
- Isaac Collins, Class of 2020, Jun. 2017 to date
- Ashton Stevens, Class of 2019, Jun. 2017 to date
- Danie
- Matec

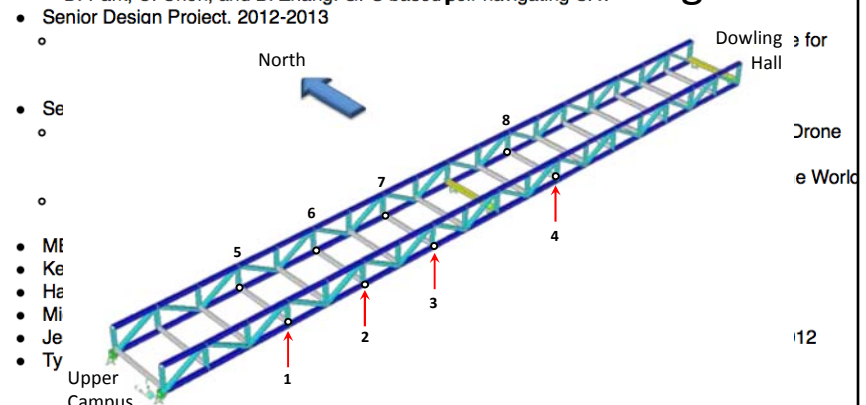
#### Alum

- Kathr
- Zach
- Anuth
- Ryan
- Syed
- Craig
- Ryan
- Terrer
- Dong
- Cody
- Oghe
- Pratik
- Corne

## Aerial Formation Flying



## SHM over a campus footbridge



# Trailer

**SPARTN**—Signal Processing and RoboTic Networks Lab at Tufts



# My Research Lab: Theory



**Reza (2011-15):**  
Graph-theoretic estimation

Best paper  
Journal cover



**Xi (2012-16):**  
Optimization over directed graphs

4 TAC papers



**Sam (2013-):**  
Fusion in non-deterministic graphs

2 Best papers  
6 IEEE journal papers



**Fakhteh (2014-):**  
Distributed estimation cont...d



**Xin (2016-):**  
Optimization, Graph theory

# My Research: In depth

---

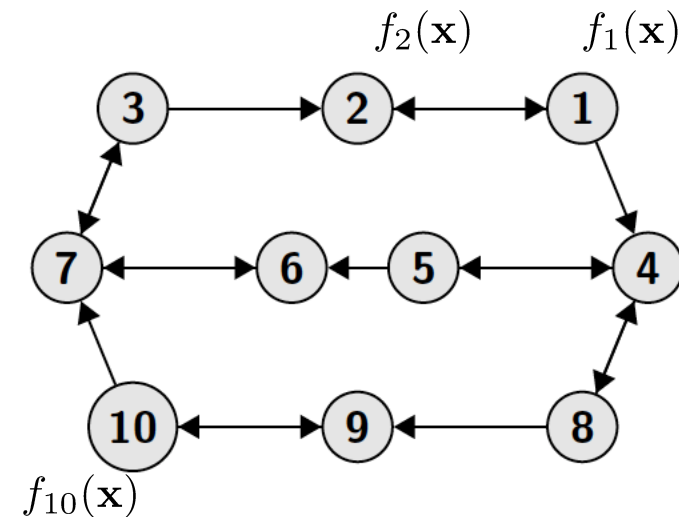
- Distributed Intelligence in **multi-agent systems**
  - Estimation, optimization, and control over **graphs (networks)**
- **Mobile** → **Dynamic**
- **Heterogeneous** → **Directed**
- **Autonomous** → **Non-deterministic**
- **Applications:**
  - Cyber-physical systems, IoTs, Big Data
  - Aerial SHM, Power grid, Personal exposome
  - Distributed Optimization: Path planning and Formation control

# Optimization over directed graphs

# Problem

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

- Agents interact over a graph
  - Directional informational flow
- No center with all information





# A nice solution

- **Gradient Descent**

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

- No one knows the function  $f$

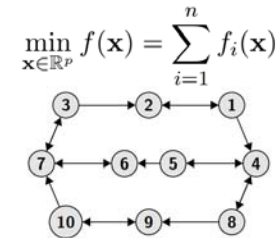
- **Local Gradient Descent**

$$\mathbf{x}_{k+1}^i = \mathbf{x}_k^i - \alpha_k \nabla f_i(\mathbf{x}_k^i)$$

- Converges to only to a local optimal

- **Distributed Gradient Descent [Nedich et al., 2009]: Fuse Information**

$$\mathbf{x}_{k+1}^i = \sum_{j \in \mathcal{N}_i} w_{ij} \mathbf{x}_k^j - \alpha_k \nabla f_i(\mathbf{x}_k^i)$$

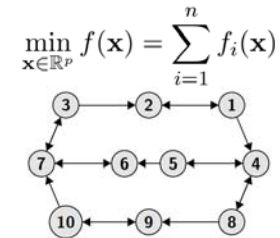


# Distributed Gradient Descent

- Distributed Gradient Descent**

$$\text{Local: } \mathbf{x}_{k+1}^i = \sum_{j \in \mathcal{N}_i} w_{ij} \mathbf{x}_k^j - \alpha_k \nabla f_i(\mathbf{x}_k^i)$$

$$\text{Network: } \mathbf{x}_{k+1} = W \mathbf{x}_k - \alpha_k \nabla \mathbf{f}(\mathbf{x}_k)$$



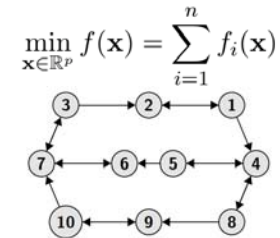
- $W = \{w_{ij}\}$  is a doubly-stochastic matrix (underlying graph is balanced)
- Step-size goes to zero (but not too fast)
- Agreement:**  $W\mathbf{1} = \mathbf{1}$
- Optimality:**  $\mathbf{1}_n^\top \nabla \mathbf{f}_\infty = 0$
- Lets do a simple analysis...

# Distributed Gradient Descent

- Distributed Gradient Descent**

$$\text{Local: } \mathbf{x}_{k+1}^i = \sum_{j \in \mathcal{N}_i} w_{ij} \mathbf{x}_k^j - \alpha_k \nabla f_i(\mathbf{x}_k^i)$$

$$\text{Network: } \mathbf{x}_{k+1} = W \mathbf{x}_k - \alpha_k \nabla \mathbf{f}(\mathbf{x}_k)$$



- Assume the corresponding sequences converge to their limits**

$$\begin{aligned} \mathbf{x}_\infty &= W \mathbf{x}_\infty - \alpha_\infty \nabla \mathbf{f}(\mathbf{x}_\infty) \\ &= W \mathbf{x}_\infty - 0 \cdot \nabla \mathbf{f}(\mathbf{x}_\infty) \end{aligned} \longrightarrow (I_n - W) \mathbf{x}_\infty = 0$$

- Let  $W$  be CS but not RS**

$$\mathbf{1}^\top W = \mathbf{1}^\top \text{ and } W \boldsymbol{\pi} = \boldsymbol{\pi} \neq \mathbf{1}$$

- Then  $\mathbf{x}_\infty = c \boldsymbol{\pi}$ , no agreement!**

- Let  $W$  be RS but not CS**

$$W \mathbf{1} = \mathbf{1} \text{ and } \boldsymbol{\pi}^\top W = \boldsymbol{\pi}^\top$$

- Then  $\mathbf{x}_\infty = c \mathbf{1}$ , i.e., agreement**

- But suboptimal!**

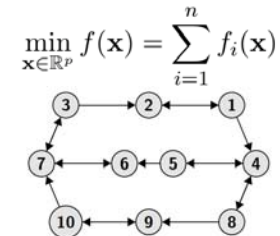
$$\boldsymbol{\pi}^\top \nabla \mathbf{f}(\mathbf{x}_\infty) = \sum_{i=1}^n \pi_i \nabla f_i(c) = 0$$

# Distributed Gradient Descent

- Distributed Gradient Descent**

$$\text{Local: } \mathbf{x}_{k+1}^i = \sum_{j \in \mathcal{N}_i} w_{ij} \mathbf{x}_k^j - \alpha_k \nabla f_i(\mathbf{x}_k^i)$$

$$\text{Network: } \mathbf{x}_{k+1} = W \mathbf{x}_k - \alpha_k \nabla \mathbf{f}(\mathbf{x}_k)$$



- If  $W$  is RS but not CS (unbalanced directed graphs), agents agree on a suboptimal solution

$$W\mathbf{1} = \mathbf{1} \text{ and } \pi^\top W = \pi^\top \quad \pi^\top \nabla \mathbf{f}(\mathbf{x}_\infty) = \sum_{i=1}^n \pi_i \nabla f_i(c) = 0$$

- Consider a modification (Nedich 2013, similar in spirit but with different execution):

$$\mathbf{x}_{k+1}^i = \sum_{j \in \mathcal{N}_i} w_{ij} \mathbf{x}_k^j - \alpha_k \frac{\nabla f_i(\mathbf{x}_k^i)}{\underbrace{\left[ \mathbf{y}_k^i \right]_i}_{\rightarrow \pi}}$$

- Row-stochasticity guarantees agreement, scaling ensures optimality
- Estimate the left eigenvector?

# Estimating the left eigenvector

- $A = \{a_{ij}\}$  is row-stochastic with  $\pi^\top A = \pi^\top$
- Consider the following iteration:

$$y_{k+1,i} = \sum_{j=1}^n a_{ij} y_{k,j} \quad y_{0,i} = e_i$$

$$Y_{k+1} = \begin{bmatrix} y_{k+1,1}^\top \\ \vdots \\ y_{k+1,n}^\top \end{bmatrix} = \begin{bmatrix} \sum_{j=1}^n a_{1j} y_{k,j}^\top \\ \vdots \\ \sum_{j=1}^n a_{nj} y_{k,j}^\top \end{bmatrix} = \begin{bmatrix} a_{11} y_{k,1}^\top + \cdots + a_{1n} y_{k,n}^\top \\ \vdots \\ a_{n1} y_{k,1}^\top + \cdots + a_{nn} y_{k,n}^\top \end{bmatrix} = AY_k$$

$$Y_\infty \triangleq \lim_{k \rightarrow \infty} Y_{k+1} = A^\infty Y_0 = A^\infty I_n = \boxed{A^\infty = \mathbf{1}_n \pi^\top}$$

- Every agent learns the entire left eigenvector asymptotically
- Similar method learns the right eigenvector for CS matrices

# Optimization over directed graphs: Recipe

- 1. Design row- or column-stochastic weights
  - 2. Estimate the non-1 eigenvector for the eval of 1
  - 3. Scale to remove the imbalance
- 
- Side note: Push-sum algorithm (Gehrke et al., 2003; Vetterli et al., 2010)

# Related work (a very small sample)

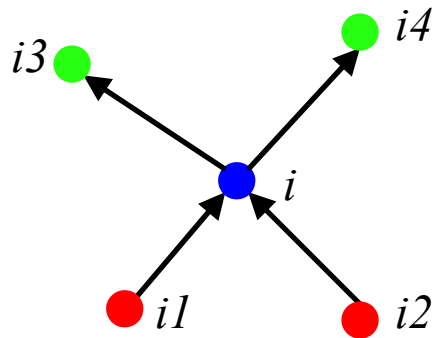
- Algorithms over undirected graphs:
  - Distributed Gradient Descent (Nedich et al., 2009)
    - Non-smooth
  - EXTRA (Yin et al., Apr. 2014)
    - Fuses information over past two iterates
    - Use gradient information over past two iterates
    - Smooth, Strong-convexity, Linear convergence
  - NEXT (Scutari et al., Dec. 2015)
    - Functions are smooth non-convex + non-smooth convex
  - Harnessing smoothness ... (Li et al., May 2016)
    - Some similarities to EXTRA

## Related work (a small sample)

- Add push-sum to the previous obtain algorithms for directed graphs:
  - Gradient Push (Nedich et al., 2013)
    - Sub-linear convergence
  - DEXTRA (Khan et al., Oct. 2015)
    - Strong-convexity, Linear convergence
    - Difficult to compute step-size interval
  - SONATA (Scutari et al., Jul. 2016)
    - Functions are (smooth non-convex + non-smooth convex)
    - Sub-linear convergence
  - ADD-OPT (Khan et al., Jun. 2016) and PUSH-DIGing (Nedich et al., Jul. 2016)
    - Strong-convexity, Linear convergence
    - Step-size interval lower bound is 0
- All these algorithms employ column-stochastic matrices



# Column- vs. Row-stochastic Weights



$$A = \begin{bmatrix} & a_{i_3 i} & & & \\ & 0 & & & \\ & 0 & & & \\ 0 & a_{i i_1} & 0 & a_{i i} & a_{i i_2} \\ & & & a_{i_4 i} & \\ & & & \uparrow & \\ & & & \text{outgoing} & \end{bmatrix} \quad \leftarrow \text{incoming weights at } i$$

- Incoming weights are simpler to design
- For column sum to be 1, agent  $i$  cannot design the incoming weights as it does not know the neighbors of  $i1$  and  $i2$ 
  - Column-stochastic weights thus are designed at outgoing edges
  - Requires the knowledge of out-neighbors or out-degree

# Optimization with Row-stochastic weights

- $A = \{a_{ij}\}$  is row-stochastic

Left Eigenvector:  $\mathbf{y}_{k+1}^i = \sum_{j \in \mathcal{N}_i^{\text{in}}} a_{ij} \mathbf{y}_k^j, \quad (\text{vector in } \mathbb{R}^n)$

Update:  $\mathbf{x}_{k+1}^i = \sum_{j \in \mathcal{N}_i^{\text{in}}} a_{ij} \mathbf{x}_k^j - \alpha \mathbf{z}_k^i$

$$\mathbf{z}_{k+1}^i = \sum_{j \in \mathcal{N}_i^{\text{in}}} a_{ij} \mathbf{z}_k^j + \frac{\nabla f_i(\mathbf{x}_{k+1}^i)}{y_{k+1}^{ii}} - \frac{\nabla f_i(\mathbf{x}_k^i)}{y_k^{ii}}$$

- Row-stochastic weight design is simple
- However, in contrast to CS methods:
  - Agents run an  $n$ th order consensus for the left eigenvector
  - Agents need unique identifiers

# Optimization with Row-stochastic weights

- $A = \{a_{ij}\}$  is row-stochastic
- **Vector form of the algorithm:** arbitrary  $\mathbf{x}_0$ ,  $\tilde{Y}_0 = Y_0 = I_n$ , and  $\mathbf{z}_0 = \nabla f_0$

$$Y_{k+1} = AY_k, \quad Y_\infty = \mathbf{1}_n \boldsymbol{\pi}^\top$$

$$\mathbf{x}_{k+1} = 2A\mathbf{x}_k - A^2\mathbf{x}_{k-1} - \alpha \left( \tilde{Y}_k^{-1} \nabla \mathbf{f}_k - \tilde{Y}_{k-1} \nabla \mathbf{f}_{k-1} \right)$$

- In contrast, with a column-stochastic  $B$ , ADDOPT/PUSH-DIGing is:

$$\mathbf{x}_{k+1} = 2B\mathbf{x}_k - B^2\mathbf{x}_{k-1} - \alpha \left( \nabla \mathbf{f}(Y_k^{-1}\mathbf{x}_k) - \nabla \mathbf{f}_{k-1}(Y_{k-1}^{-1}\mathbf{x}_{k-1}) \right)$$

- Iterate does not result in agreement
- The function argument is scaled by the right eigenvector
- Ensures optimality

# Optimization with Row-stochastic weights

- **Algorithm:** arbitrary  $\mathbf{x}_0$ ,  $\tilde{Y}_0 = Y_0 = I_n$ , and  $\mathbf{z}_0 = \nabla f_0$

$$Y_{k+1} = AY_k, \quad Y_\infty = \mathbf{1}_n \boldsymbol{\pi}^\top$$

$$\mathbf{x}_{k+1} = 2A\mathbf{x}_k - A^2\mathbf{x}_{k-1} - \alpha \left( \tilde{Y}_k^{-1} \nabla \mathbf{f}_k - \tilde{Y}_{k-1}^{-1} \nabla \mathbf{f}_{k-1} \right)$$

- **A simple intuitive argument:**
- **Assume each sequence converges to its limit, then**

$$\mathbf{x}_\infty = 2A\mathbf{x}_\infty - A^2\mathbf{x}_\infty - \alpha \left( \tilde{Y}_\infty^{-1} \nabla \mathbf{f}_\infty - \tilde{Y}_\infty^{-1} \nabla \mathbf{f}_\infty \right)$$

$$(I_n - A)^2 \mathbf{x}_\infty = \mathbf{0}$$

$$\mathbf{x}_\infty = c \mathbf{1}_n$$

- **Every agent agrees on  $c$**

# Optimization with Row-stochastic weights

- **Algorithm:** arbitrary  $\mathbf{x}_0$ ,  $\tilde{Y}_0 = Y_0 = I_n$ , and  $\mathbf{z}_0 = \nabla f_0$

$$Y_{k+1} = AY_k, \quad Y_\infty = \mathbf{1}_n \boldsymbol{\pi}^\top$$

$$\mathbf{x}_{k+1} = 2A\mathbf{x}_k - A^2\mathbf{x}_{k-1} - \alpha \left( \tilde{Y}_k^{-1} \nabla \mathbf{f}_k - \tilde{Y}_{k-1}^{-1} \nabla \mathbf{f}_{k-1} \right)$$

- Show that  $c$  is the optimal solution
- Sum the update over  $k$ :

$$\alpha \tilde{Y}_M^{-1} \nabla \mathbf{f}_M = \sum_{k=0}^{M-1} (A - A^2) \mathbf{x}_k + A\mathbf{x}_M + \sum_{k=0}^M (A - I_n) \mathbf{x}_k - \mathbf{x}_{M+1}.$$

$$\begin{aligned} \alpha \boldsymbol{\pi}^\top \tilde{Y}_M^{-1} \nabla \mathbf{f}_M &= \sum_{k=0}^{M-1} \boldsymbol{\pi}^\top (A - A^2) \mathbf{x}_k + \boldsymbol{\pi}^\top A\mathbf{x}_M + \sum_{k=0}^M \boldsymbol{\pi}^\top (A - I_n) \mathbf{x}_k - \boldsymbol{\pi}^\top \mathbf{x}_{M+1} \\ &= \boldsymbol{\pi}^\top A\mathbf{x}_M - \boldsymbol{\pi}^\top \mathbf{x}_{M+1}. \end{aligned}$$

# Optimization with Row-stochastic weights

- **Algorithm:** arbitrary  $\mathbf{x}_0$ ,  $\tilde{Y}_0 = Y_0 = I_n$ , and  $\mathbf{z}_0 = \nabla f_0$

$$Y_{k+1} = AY_k, \quad Y_\infty = \mathbf{1}_n \boldsymbol{\pi}^\top$$

$$\mathbf{x}_{k+1} = 2A\mathbf{x}_k - A^2\mathbf{x}_{k-1} - \alpha \left( \tilde{Y}_k^{-1} \nabla \mathbf{f}_k - \tilde{Y}_{k-1}^{-1} \nabla \mathbf{f}_{k-1} \right)$$

- **Asymptotically**

$$\alpha \boldsymbol{\pi}^\top \tilde{Y}_\infty^{-1} \nabla \mathbf{f}_\infty = \boldsymbol{\pi}^\top A \mathbf{x}_\infty - \boldsymbol{\pi}^\top \mathbf{x}_\infty$$

$$\alpha \mathbf{1}_n^\top \nabla \mathbf{f}_\infty = 0$$

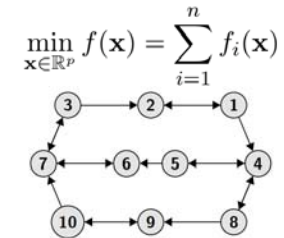
$$\nabla f_1(x_\infty^1) + \nabla f_2(x_\infty^2) + \cdots \nabla f_n(x_\infty^n) = \nabla f_1(c) + \nabla f_2(c) + \cdots \nabla f_n(c) = 0$$

# Optimization with Row-stochastic weights

- **Algorithm:** arbitrary  $\mathbf{x}_0$ ,  $\tilde{Y}_0 = Y_0 = I_n$ , and  $\mathbf{z}_0 = \nabla f_0$   
$$Y_{k+1} = AY_k, \quad Y_\infty = \mathbf{1}_n \boldsymbol{\pi}^\top$$
$$\mathbf{x}_{k+1} = 2A\mathbf{x}_k - A^2\mathbf{x}_{k-1} - \alpha \left( \tilde{Y}_k^{-1} \nabla \mathbf{f}_k - \tilde{Y}_{k-1}^{-1} \nabla \mathbf{f}_{k-1} \right)$$
- We assumed that the sequences reach their limit
- However, under what conditions and at what rate?

# Convergence conditions

- Assume strong-connectivity, Lipschitz-continuous gradients, strongly-convex functions



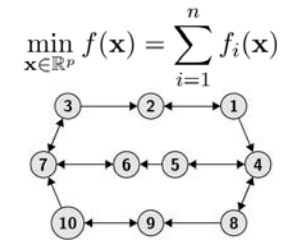
- Consider  $\mathbf{t}_k = \begin{bmatrix} \|\mathbf{x}_k - \hat{\mathbf{x}}_k\| \\ \|\hat{\mathbf{x}}_k - \mathbf{x}^*\|_2 \\ \|\mathbf{z}_k - \hat{\mathbf{z}}_k\| \end{bmatrix}$
- If some norm of  $\mathbf{t}_k$  goes to 0, then each element goes to 0 and the sequences converge to their limits



# Convergence conditions

$$\mathbf{x}_{k+1} = 2A\mathbf{x}_k - A^2\mathbf{x}_{k-1} - \alpha \left( \tilde{Y}_k^{-1} \nabla \mathbf{f}_k - \tilde{Y}_{k-1}^{-1} \nabla \mathbf{f}_{k-1} \right)$$

- Assume strong-connectivity, Lipschitz-continuous gradients, strongly-convex functions



- Theorem 1.** *Let Assumptions A1 and A2 hold. We have*

$$\mathbf{t}_{k+1} \leq G\mathbf{t}_k + H_k \mathbf{s}_k, \quad \forall k.$$

$$G_\alpha = \begin{bmatrix} \sigma & 0 & \alpha \\ \alpha c n l & 1 - \alpha n s & 0 \\ c(\tau + \alpha n l) & \alpha d_1 l n & \sigma + \alpha c \end{bmatrix}$$

- Lemma:  $H_k$  goes to 0 linearly
- Lemma: Spectral radius of  $G$  is less than 1

# Convergence conditions

- **Lemma:** For all values of  $\alpha \in (0, \alpha_1)$ , we have  $\rho(G_\alpha) < 1$ , where

$$\alpha_1 = \frac{\sqrt{\Delta^2 + 4cn^3l(l+s)s(1-\sigma)^2} - \Delta}{2cn^2l(l+s)} \text{ and } \Delta = cns(\tau + 1 - \sigma).$$

- Recall that

$$G_\alpha = \begin{bmatrix} \sigma & 0 & \alpha \\ \alpha cnl & 1 - \alpha ns & 0 \\ c(\tau + \alpha nl) & \alpha d_1 ln & \sigma + \alpha c \end{bmatrix}, \quad G_0 = \begin{bmatrix} \sigma & 0 & 0 \\ 0 & 1 & 0 \\ c\tau & 0 & \sigma \end{bmatrix}.$$

- Hence,  $\rho(G_0) = 1$  because  $\sigma < 1$ .

# Convergence Rate

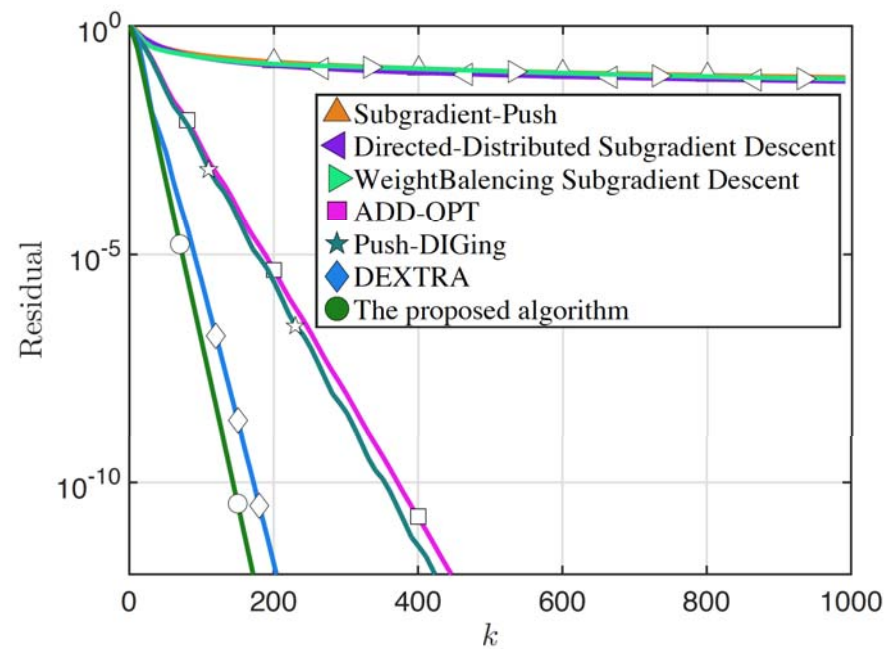
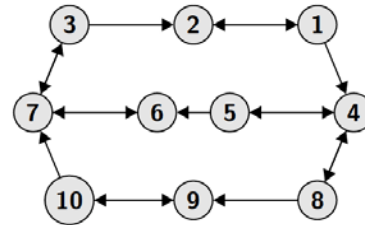
- **Theorem 2.** *With the step-size,  $\alpha \in (0, \alpha_1)$ , the sequence,  $\{\mathbf{x}_k\}$ , converges linearly to the optimal solution,  $\mathbf{x}^*$ , i.e., there exist some constant  $M > 0$  such that*

$$\|\mathbf{x}_k - \mathbf{x}^*\|_2 \leq M(\gamma + \xi)^k, \quad \forall k,$$

*where  $\xi$  is an arbitrarily small constant.*

- **The rate variable  $\gamma$  is the max of fusion rate and the rate at which  $G$  decays**

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



# Conclusions

---

- Optimization with row-stochastic matrices
- Does not require the knowledge of out-neighbors or out-degree
  - Agents require unique identifiers
- Strongly-convex functions with Lipschitz-continuous gradients
- Strongly-connected directed graphs
- Linear convergence

## More Information

---

- My webpage: <http://www.eecs.tufts.edu/~khan/>
- My email: [khan@ece.tufts.edu](mailto:khan@ece.tufts.edu)
- My Lab's YouTube channel:  
<https://www.youtube.com/user/SPARTNatTufts/videos/>