# A Newton Tracking Algorithm with Exact Linear Convergence Rate for Decentralized Consensus Optimization

## Jiaojiao Zhang

Department of Systems Engineering and Engineering Management,
The Chinese University of Hong Kong

Joint work with
Qing Ling (SYSU) and Anthony Man-Cho So (CUHK)
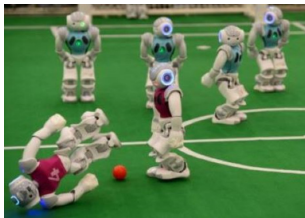
IEEE-CDC 2020

# Outline

# Outline

# Background

Decentralized optimization

## Background

Consider the decentralized convex optimization problem

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

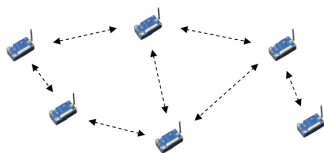- $f_i(x)$ is a convex and twice continuously differentiable function.



Fig 1. Decentralized network

- Data are distributed across a connected network of $n$ nodes.
- Each node is only allowed to send/receive information to/from its neighboring nodes.
- All the nodes cooperate to get a common variable.

# Background

An equivalent <span style="color:red">decentralized</span> formulation is

$$\{x_i^*\}_{i=1}^n := \arg\min_{\{x_i\}_{i=1}^n} \sum_{i=1}^n f_i(x_i)$$
$$\text{s.t. } x_i = x_j, \forall j \in \mathcal{N}_i, \forall i$$

Two key components for decentralized optimization

- consensus - all nodes must agree on the same state, i.e., $x_1^* = \cdots = x_n^*$.
- optimality - the same state should be the minimizer of the original problem, i.e., $x_1^* = \cdots = x_n^* = x^*$.

# Related First-order Work

Primal method

- Gradient methods [Nedic 2009, Yuan 2016]
- Gradient Tracking [Lorenzo 2015, Qu 2017, Nedic 2017,Sun 2019]

Primal-Dual method

- Decentralized Alternating Direction Method of Multipliers (DADMM) [Shi 2014,Chang 2015]
- Decentralized linearized ADMM [Ling 2015]
- Dual Ascent [Maros 2018]

Other method

- EXTRA [Shi 2015]
- NIDS [Li 2019]

# Related Second-order Work

Penalized second-order algorithms converge to a neighborhood of an optimal solution

- Network Newton [Mokhtari 2016]
- Distributed asynchronous Newton-based algorithm [Mansoori 2019]

Primal-dual second-order methods achieve exact convergence with linear rates

- DQM [Mokhtari 2016]
- ESOM [Mokhtari 2016]

Second-order methods with superlinear convergence rates under stricter conditions

- Distributed averaged quasi-Newton method for a master-slave network [Soori 2019]
- Polyak's adaptive Newton method running a finite-time set-consensus inner loop [Zhang 2020]

# Outline

# Algorithm Development

We make the following assumptions

## Assumption 1

*Introduce the mixing matrix $W$ with elements $w_{ij} \geq 0$. $w_{ij} = 0$ if and only if $j \notin \mathcal{N}_i \cup \{i\}$. Further, $W^T = W$, $W 1_{n \times 1} = 1_{n \times 1}$ and $null(I - W) = span(1_{n \times 1})$.*

## Assumption 2

*The local objective functions $f_i(x_i)$ are twice differentiable. Hessians $\nabla^2 f_i(x_i)$ are bounded by*

$$\mu_f I_p \preceq \nabla^2 f_i(x_i) \preceq L_f I_p.$$

# Algorithm Development

Consider the decentralized optimization problem

$$\{x_i^*\}_{i=1}^n := \arg\min_{\{x_i\}_{i=1}^n} \sum_{i=1}^n f_i(x_i)$$

$$\text{s.t. } x_i = x_j, \forall j \in \mathcal{N}_i, \forall i$$

Global negative Newton direction at $\bar{x}^t \triangleq \frac{1}{n}\sum_{i=1}^n x_i^t$ is

$$u^t \triangleq \left(\frac{1}{n}\sum_{i=1}^n \nabla^2 f_i(\bar{x}^t)\right)^{-1} \left(\frac{1}{n}\sum_{i=1}^n \nabla f_i(\bar{x}^t)\right)$$

**Our idea**: Use a local $u_i^t$ to track the negative global Newton direction such that $u_i^t \approx u^t$.

# Algorithm Development

The proposed **Newton tracking**

$$x_i^{t+1} = x_i^t - u_i^t$$

$$u_i^{t+1} = \left(\nabla^2 f_i\left(x_i^{t+1}\right) + \epsilon I_p\right)^{-1} \left[\left(\nabla^2 f_i\left(x_i^t\right) + \epsilon I_p\right) u_i^t + \nabla f_i\left(x_i^{t+1}\right) - \nabla f_i\left(x_i^t\right)\right.$$

$$\left. + 2\alpha(x_i^{t+1} - \sum_{j \in \mathcal{N}_i} w_{ij} x_j^{t+1}) - \alpha(x_i^t - \sum_{j \in \mathcal{N}_i} w_{ij} x_j^t)\right]$$

with initialization $x_i^0 = 0_p$ and $u_i^0 = \left(\nabla^2 f_i\left(0_p\right) + \epsilon I_p\right)^{-1} \nabla f_i\left(0_p\right)$.

- The global Hessian $\frac{1}{n} \sum_{i=1}^n \nabla^2 f_i\left(\bar{x}^{t+1}\right)$ is replaced by the regularized local Hessian $\nabla^2 f_i\left(x_i^{t+1}\right) + \epsilon I_p$.
- The global gradient $\frac{1}{n} \sum_{i=1}^n \nabla f_i\left(\bar{x}^t\right)$ is replaced by three terms that are locally computable.

# Why $u_i^t \approx u^t$ ?

Rewrite Newton tracking

$$\left(\nabla^2 f_i\left(x_i^{t+1}\right) + \epsilon I_p\right) u_i^{t+1} = [\left(\nabla^2 f_i\left(x_i^t\right) + \epsilon I_p\right) u_i^t + \nabla f_i\left(x_i^{t+1}\right) - \nabla f_i\left(x_i^t\right)$$
$$+ 2\alpha(x_i^{t+1} - \sum_{j \in \mathcal{N}_i} w_{ij} x_j^{t+1}) - \alpha(x_i^t - \sum_{j \in \mathcal{N}_i} w_{ij} x_j^t)]$$

Sum up over $i = 1 \ldots n$ and invoke the double stochasticity of $W$

$$\sum_{i=1}^n \left(\nabla^2 f_i\left(x_i^{t+1}\right) + \epsilon I_p\right) u_i^{t+1} = \sum_{i=1}^n [\left(\nabla^2 f_i\left(x_i^t\right) + \epsilon I_p\right) u_i^t + \nabla f_i\left(x_i^{t+1}\right) - \nabla f_i\left(x_i^t\right)]$$

With $\sum_{i=1}^n \nabla f_i\left(x_i^0\right) = \sum_{i=1}^n \left(\nabla^2 f_i\left(x_i^0\right) + \epsilon I_p\right) u_i^0$, sum up over time $t$

$$\sum_{i=1}^n \left(\nabla^2 f_i\left(x_i^t\right) + \epsilon I_p\right) u_i^t = \sum_{i=1}^n \nabla f_i\left(x_i^t\right)$$

Thus, when $x_i^t$ is close to $\bar{x}^t$, $u_i^t$ tracks a regularized Newton direction.

# Newton Tracking

Newton tracking can be written in a compact form

$$x^{t+1} = x^t - u^t$$
$$u^{t+1} = \left(H^{t+1}\right)^{-1} \left[ H^t u^t + \nabla f\left(x^{t+1}\right) - \nabla f\left(x^t\right) \right.$$
$$\left. + \alpha(I - W)\left(2x^{t+1} - x^t\right) \right]$$

- x (or u) stacks local variables such as $x \triangleq [x_1; \ldots; x_n] \in \mathbb{R}^{np}$
- $W \triangleq W \otimes I_p \in \mathbb{R}^{np \times np}$
- $f(x) = f\left(x_1, \cdots, x_n\right) = \sum_{i=1}^{n} f_i\left(x_i\right)$.
- $\nabla f(x) = \left[\nabla f_1\left(x_1\right); \ldots; \nabla f_n\left(x_n\right)\right] \in \mathbb{R}^{np}$
- $H \in \mathbb{R}^{np \times np}$ is a block diagonal matrix whose $i$-th diagonal block is $\nabla^2 f_i(x) + I_p$

# Connection with Gradient Tracking

Gradient tracking proceeds as

$$x^{t+1} = Wx^t - \alpha y^t$$
$$y^{t+1} = Wy^t + \nabla f\left(x^{t+1}\right) - \nabla f\left(x^t\right)$$

with initialization $y^0 = \nabla f\left(x^0\right)$.

To see the connection with Newton tracking, we rewrite

$$x^{t+1} = x^t - r^t$$
$$r^{t+1} = Wr^t + \alpha\left[\nabla f\left(x^{t+1}\right) - \nabla f\left(x^t\right)\right] + (I - W)\left(x^{t+1} - Wx^t\right)$$

with $r^t = (I - W)x^t + \alpha y^t \in \mathbb{R}^{np}$ and $\sum_{i=1}^{n} r_i^0 = \alpha \sum_{i=1}^{n} \nabla f_i\left(x_i^0\right)$. Sum up over node $i$ and time $t$

$$\sum_{i=1}^{n} r_i^t = \alpha \sum_{i=1}^{n} \nabla f_i\left(x_i^t\right)$$

Thus, when $x_i^t$ is close to $\bar{x}^t$, $r_i^t$ tracks a scaled gradient direction.

## Connection with Primal-dual Algorithms

Under Assumption 1, we have $(I - W)^{\frac{1}{2}}x = 0$ if and only if $x_1 = \cdots = x_n$. Thus, the original problem is equivalent to

$$x^* \triangleq \arg \min_{x} f(x) \quad \text{s.t. } (I - W)^{\frac{1}{2}}x = 0.$$

The augmented Lagrangian is

$$L(x, v) = f(x) + \left\langle v, (I - W)^{\frac{1}{2}}x \right\rangle + \frac{\alpha}{2}x^T(I - W)x.$$

We use a quadratic approximation of $f$ and a linear approximation of $x \mapsto \frac{\alpha}{2}x^T(I - W)x$. The update of $x^{t+1}$ is given by the solution of

$$\min_{x} \left\langle \nabla f\left(x^t\right) + (I - W)^{\frac{1}{2}}v^t + \alpha(I - W)x^t, x - x^t \right\rangle$$
$$+ \frac{1}{2}\left(x - x^t\right)^T \nabla^2 f\left(x^t\right)\left(x - x^t\right) + \frac{\epsilon}{2}\left\|x - x^t\right\|^2.$$

Thus, the updates of $x^{t+1}$ and $v^{t+1}$ are

$$x^{t+1} = x^t - (H^t)^{-1} \left[ \nabla f(x^t) + (I - W)^{\frac{1}{2}} v^t + \alpha(I - W)x^t \right]$$

$$v^{t+1} = v^t + \alpha(I - W)^{\frac{1}{2}} x^{t+1}$$

where we set $x^0 = 0$ and $v^0 = 0$. By manipulation, we get

$$x^{t+1} = x^t - (H^t)^{-1} q^t$$

$$q^{t+1} = q^t + \nabla f(x^{t+1}) - \nabla f(x^t) + \alpha(I - W)(2x^{t+1} - x^t),$$

which is equivalent to Newton tracking in the sense that $u^t = (H^t)^{-1} q^t$.

# Outline

# Convergence

## Theorem 1

*Under Assumptions 1 and 2, suppose that the parameters $\epsilon$ and $\alpha$ satisfy $\epsilon - \alpha\lambda_{\max}(I - W) > \frac{4L_f^2}{\mu_f}$. Then Newton tracking starting with $x_i^0 = 0_p$ and $u_i^0 = \left(\nabla^2 f_i(0_p) + \epsilon I_p\right)^{-1} \nabla f_i(0_p)$ satisfies*

$$\left\|\zeta^{t+1} - \zeta^*\right\|_G^2 \leq \frac{1}{1 + \delta'} \left\|\zeta^t - \zeta^*\right\|_G^2,$$

*where $\delta' > 0$.*

- Define $\zeta^t = \begin{bmatrix} x^t \\ v^t \end{bmatrix}, \zeta^* = \begin{bmatrix} x^* \\ v^* \end{bmatrix}, G = \begin{bmatrix} \epsilon I - \alpha(I - W) & 0 \\ 0 & \frac{1}{\alpha}I \end{bmatrix}.$

- Theorem 1 shows that the sequence $\{\left\|\zeta^{t+1} - \zeta^*\right\|_G^2\}_t$ converges linearly with the factor $\frac{1}{1+\delta'}$.

# Outline

# Numerical Experiments

Decentralized logistic regression problem

$$\min_{x \in \mathbb{R}^p} \frac{\rho}{2} \|x\|^2 + \sum_{i=1}^{n} \sum_{j=1}^{m_i} \ln \left( 1 + \exp \left( - \left( o_{ij}^T x \right) p_{ij} \right) \right)$$

- Node $i$ has access to $m_i$ training samples, $(o_{ij}, p_{ij}) \in \mathbb{R}^p \times \{-1, +1\}$.
- Relative error $\|x^t - \hat{x}^*\| / \|x^0 - \hat{x}^*\|$.

# Comparison with Related Methods

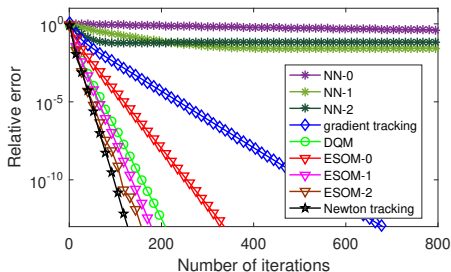- $n = 10$, connectivity ratio$=0.5$, $m_i = 12$, $p = 8$, $\rho = 0.001$



Fig 2. Relative errors of Newton tracking, gradient tracking, NN-K, ESOM-K, and DQM versus number of iterations.
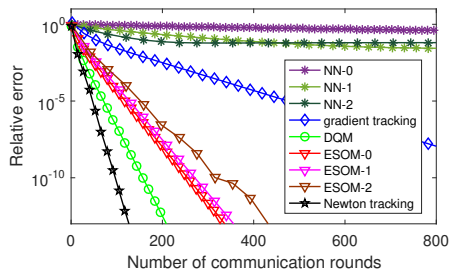


Fig 3. Relative errors of Newton tracking, gradient tracking, NN-K, ESOM-K, and DQM versus number of communication rounds.

# Effect of Network Topology

- Four different topologies including line graph, cycle graph, random graphs, and complete graph
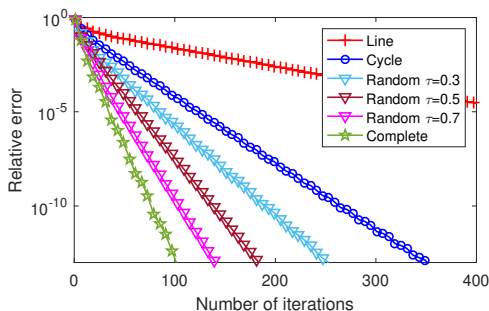- $n = 10$, $m_i = 12$, $p = 8$, connectivity ratio$=0.5$, $\rho = 0.001$



Fig 4. Relative errors of Newton tracking for different topologies

# Outline

To summarize

- We propose Newton tracking, in which each node updates its local variable along a modified local Newton direction.
- Newton tracking employs a fixed step size and yet can still be proven to converge to an exact solution.
- We give the connections between Newton tracking and several existing methods, including gradient tracking and primal-dual algorithms.
- Newton tracking is linearly convergent under the strong convexity assumption.

*Thank you !*