

Distributed convex optimization: a consensus-based Newton-Raphson approach

Damiano Varagnolo

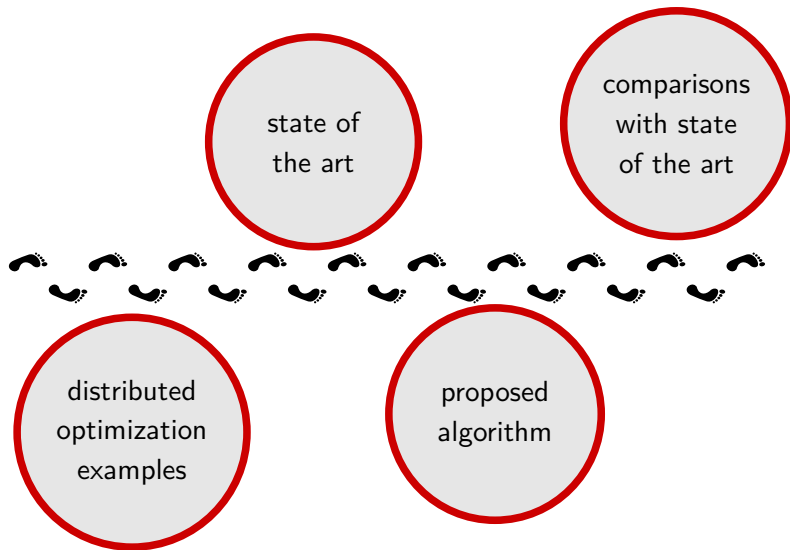
joint work with A. Cenedese, G. Pillonetto, L. Schenato, F. Zanella

Department of Information Engineering - University of Padova

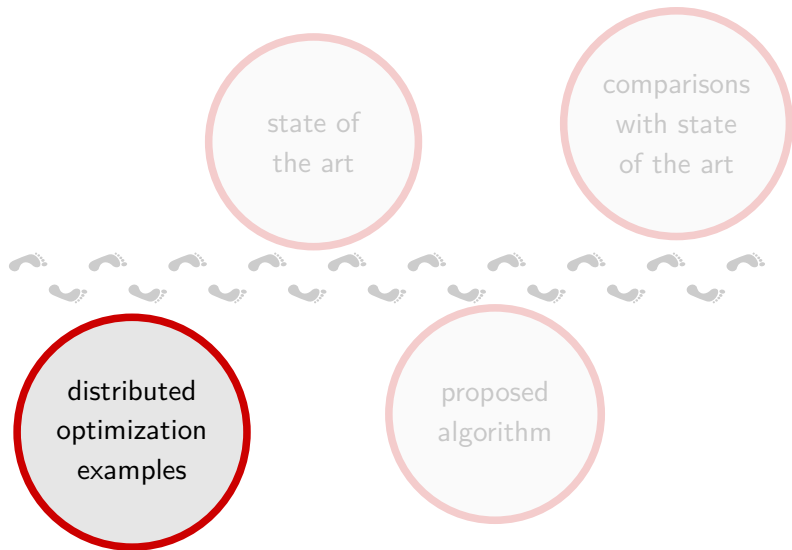
December 14th, 2011 – 50th IEEE CDC



This talk



This talk



Distributed convex optimization and its importance

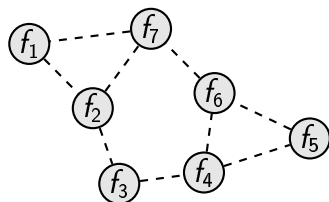
A general problem . . .

$$\begin{aligned} & \text{minimize} && f(x) = \sum_{i=1}^N f_i(x) \\ & \text{subject to} && g(x) \leq 0 \\ & && x \in \mathcal{X} \end{aligned}$$

under
convexity
assumptions

. . . motivated by multi-agents scenarios

Networked system
where neighbors
cooperate to find the
optimum



Distribution optimization - Example 1

Regression in sensor networks

(e.g. when estimation = optimization of a cost function)

Residuals minimization

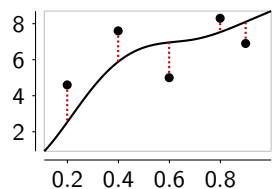
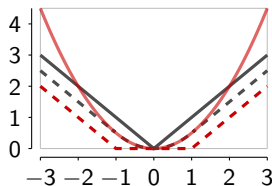
$$\begin{aligned} \min_{\theta} \quad & \sum_{i=1}^N \phi(y_i - \hat{y}_i) \\ \text{s.t.} \quad & \hat{y}_i = \theta^T x_i \end{aligned}$$

$$\phi(r) = |r|^2 \quad (\text{least squares})$$

$$\phi(r) = |r| \quad (\text{least abs. deviations})$$

$$\phi(r) = \begin{cases} 0 & \text{if } |r| < 1 \\ |r| - 1 & \text{otherwise} \end{cases} \quad (\text{Vapnik})$$

$$\phi(r) = \begin{cases} |r|^2 & \text{if } |r| < 1 \\ 2(|r| - 1) & \text{otherwise} \end{cases} \quad (\text{Huber})$$



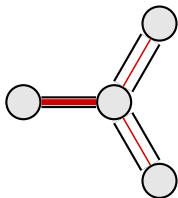
Distribution optimization - Example 2

Resource allocation in wireless systems

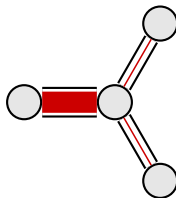
(e.g. when optimal allocation = optimization of a cost function)

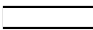

Links capacity allocation [Johansson 2008]

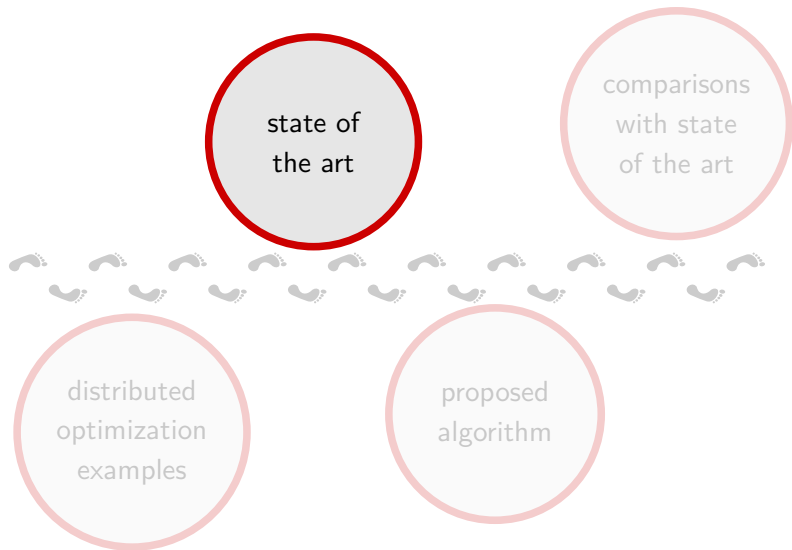
suboptimal allocation



optimal allocation



 's width = allocated link capacity
 's width = data flux



Distributed optimization methods: 3 main categories

- Primal decompositions methods
(e.g. distributed subgradients)
- Dual decompositions methods
(e.g. alternating direction method of multipliers)
- Heuristic methods
(e.g. swarm optimization, genetic algorithms)

Primal decomposition methods (distributed)

Distributed subgradient methods [Nedić Ozdaglar 2009]

$$x_i(k+1) = \mathcal{P}_{\mathcal{X}} \left[\sum_{j=1}^N a_{ij}(k)x_j(k) + \alpha_i(k)g_i(x_i(k)) \right]$$

with

- $\sum_{j=1}^N a_{ij}(k)x_j(k) :=$ aver. consensus step on *local* estimates $x_j(k)$
- $g_i(x_i(k)) :=$ *local* subgradient of *local* cost $f_i(\cdot)$ at $x_i(k)$
- $\alpha_i(k) :=$ *local* stepsize

Convergence properties [Nedić Ozdaglar (2007)]

E.g., for *bounded subgradients* and $\alpha_i(k) = \alpha$ then

$$\liminf_{k \rightarrow +\infty} f(x_i(k)) = f^* + \text{small constant}$$

Dual decomposition methods (distributed)

Alternating Direction Method of Multipliers

[Bertsekas Tsitsiklis 1997]

$$\begin{aligned} & \text{minimize} && f_1(x_1) + f_2(x_2) \\ & \text{subject to} && A_1x_1 + A_2x_2 - b = 0 \end{aligned}$$

Augmented
Lagrangian:

$$\begin{aligned} L_\rho(x_1, x_2, \lambda) := & f_1(x_1) + f_2(x_2) \\ & + \lambda^T (A_1x_1 + A_2x_2 - b) \\ & + \frac{\rho}{2} \|A_1x_1 + A_2x_2 - b\|_2^2 \end{aligned}$$

Algorithm

- 1 $x_1(k+1) = \arg \min_{x_1} L_\rho(x_1, x_2(k), \lambda(k))$
- 2 $x_2(k+1) = \arg \min_{x_2} L_\rho(x_1(k+1), x_2, \lambda(k))$
- 3 $\lambda(k+1) = \lambda(k) + \rho(A_1x_1 + A_2x_2 - b)$

Drawbacks of the considered algorithms

Primal based strategies

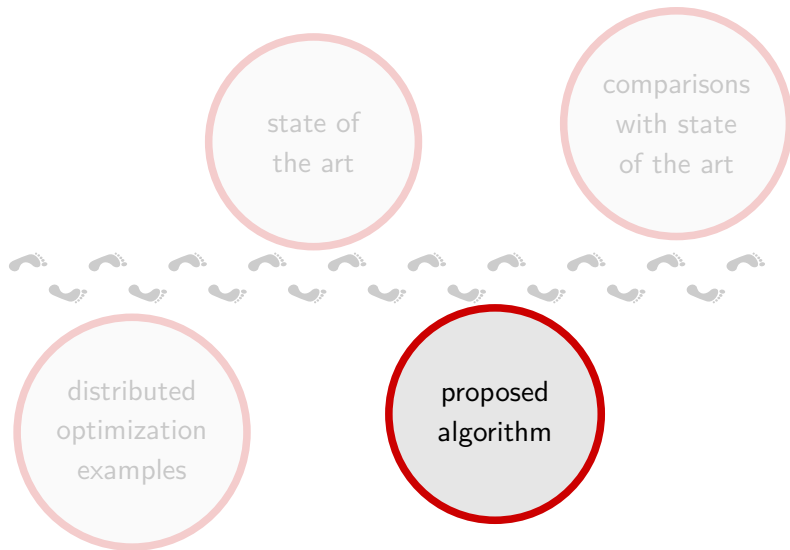
- may be slow
- may not converge to the optimum

Dual based strategies

- may be computationally expensive
- require topological knowledge
- implementation to handle time-varying graphs, time delays, etc. may require effort

The algorithm that we want:

- ① easy to be implemented
- ② with small computational requirements
- ③ does not require synchronization or topology knowledge
- ④ assured to converge to global optimum
- ⑤ inheriting good properties of standard consensus
convergence proofs, robustness, . . .



Our position in literature

How the proposed algorithm relates to other techniques?

- primal decomposition method
- uses second-order approximations

caveat:

- unconstrained convex optimization
- strong assumptions on the cost functions
(all other algorithms can work under our hypotheses)

**our contribute: better convergence speed
for primal methods**

Illustrative example: quadratic local cost functions

Derivation of the algorithm - step 1 on 3

Simplified scalar scenario

$$f_i(x) = \frac{1}{2} a_i (x - b_i)^2 + c_i \quad a_i > 0$$

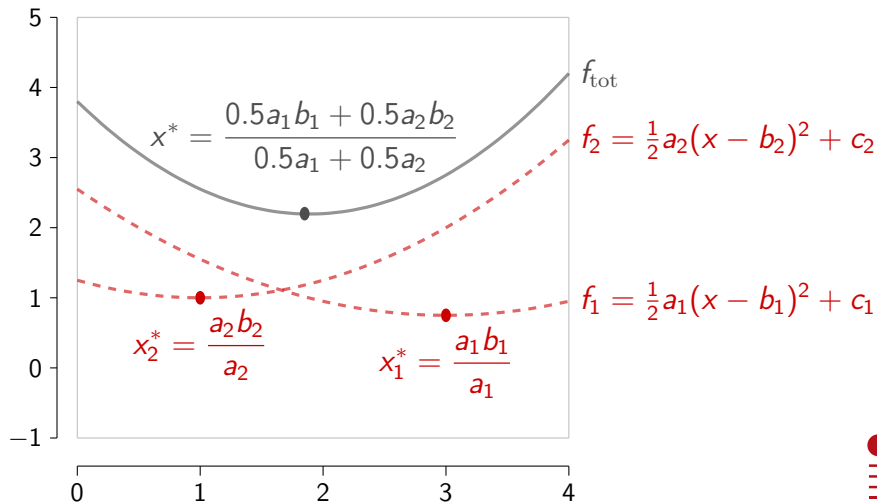
Corresponding global solution

$$x^* := \arg \min_x \sum_i f_i(x) \quad x^* = \frac{\sum_{i=1}^N a_i b_i}{\sum_{i=1}^N a_i} = \frac{\frac{1}{N} \sum_{i=1}^N a_i b_i}{\frac{1}{N} \sum_{i=1}^N a_i}$$

i.e. **parallel of 2 average consensus!**

Illustrative example: quadratic local cost functions

Derivation of the algorithm - step 1 on 3 - graphical interpretation



And for generic convex local cost functions?

Derivation of the algorithm - step 2 on 3

For quadratics ...

$$x^* = \frac{\frac{1}{N} \sum_{i=1}^N a_i b_i}{\frac{1}{N} \sum_{i=1}^N a_i} \quad \text{with} \quad \begin{aligned} &\bullet a_i b_i = f_i''(x_i) x_i - f_i'(x_i) \\ &\bullet a_i = f_i''(x_i) \end{aligned}$$

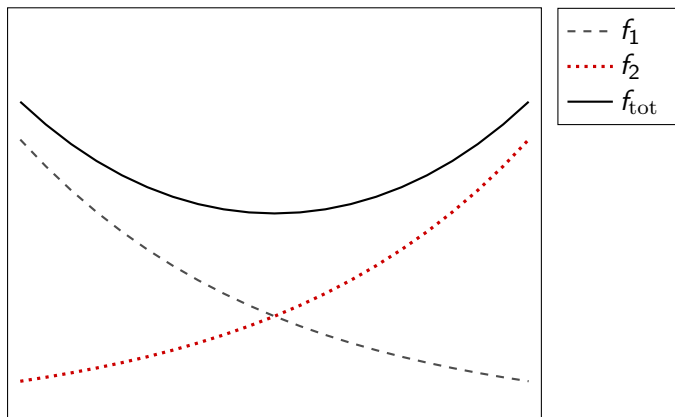
...so let's check

$$x^* \stackrel{?}{=} \frac{\frac{1}{N} \sum_{i=1}^N (f_i''(x_i) x_i - f_i'(x_i))}{\frac{1}{N} \sum_{i=1}^N f_i''(x_i)}$$

underlying idea: use Newton-Raphson approximation

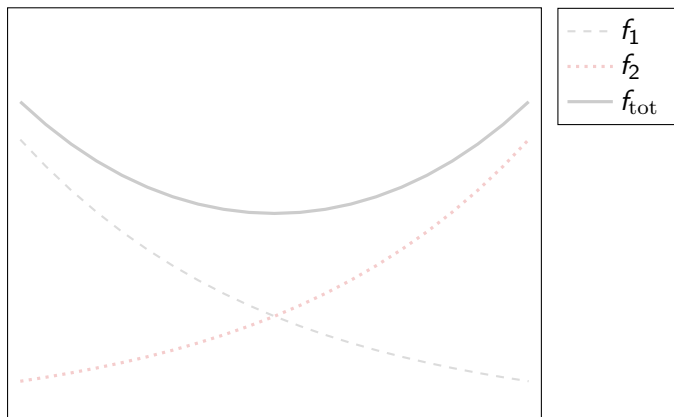
The initial idea

Derivation of the algorithm - step 2 on 3 - graphical interpretation



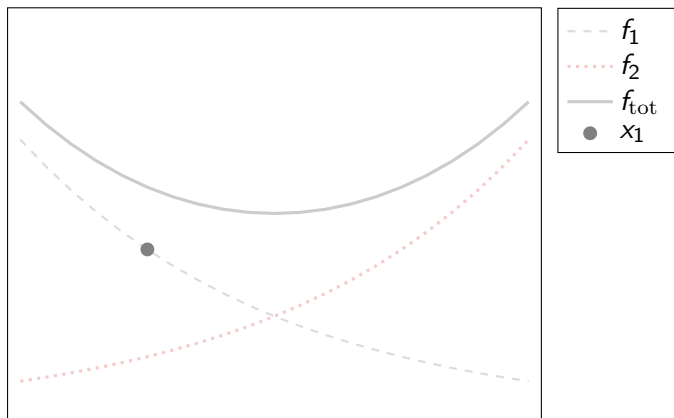
The initial idea

Derivation of the algorithm - step 2 on 3 - graphical interpretation



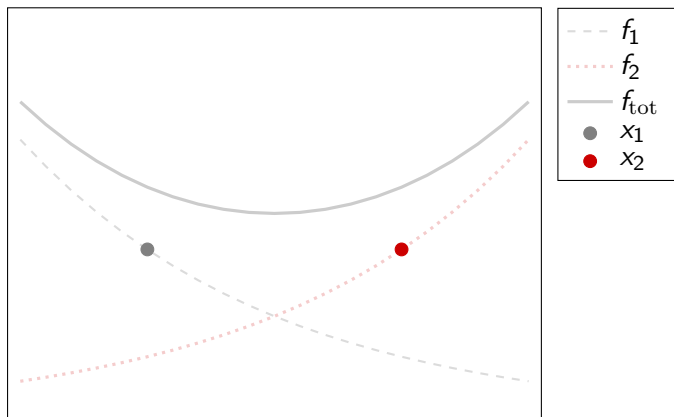
The initial idea

Derivation of the algorithm - step 2 on 3 - graphical interpretation



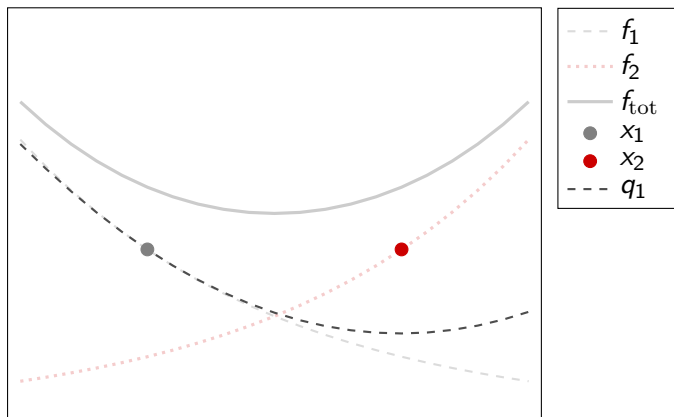
The initial idea

Derivation of the algorithm - step 2 on 3 - graphical interpretation



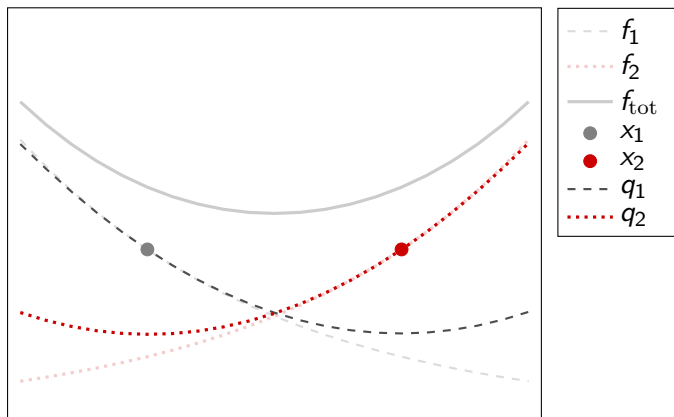
The initial idea

Derivation of the algorithm - step 2 on 3 - graphical interpretation



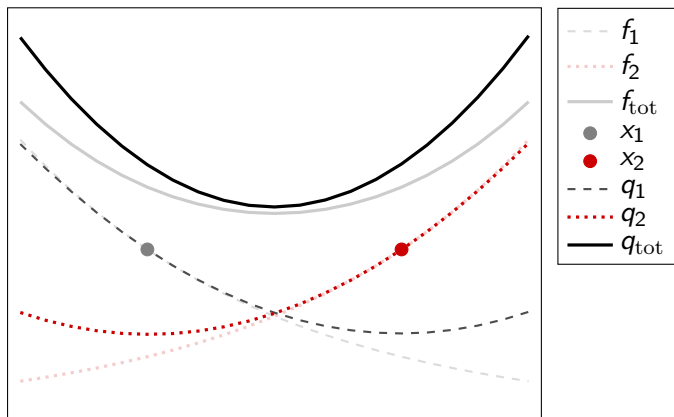
The initial idea

Derivation of the algorithm - step 2 on 3 - graphical interpretation



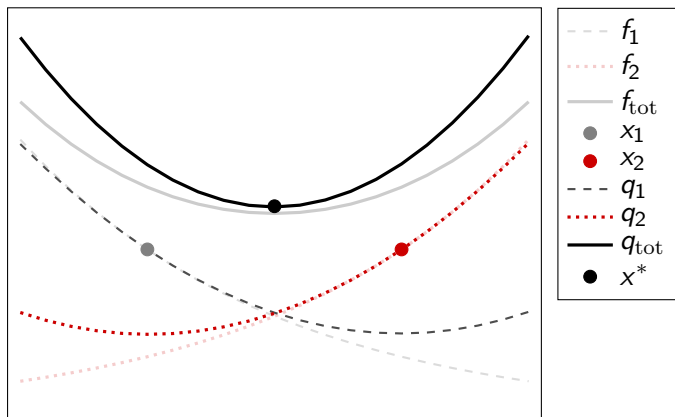
The initial idea

Derivation of the algorithm - step 2 on 3 - graphical interpretation



The initial idea

Derivation of the algorithm - step 2 on 3 - graphical interpretation



$$\text{candidate for } x^* = \frac{\frac{1}{N} \sum_{i=1}^N a_i b_i}{\frac{1}{N} \sum_{i=1}^N a_i} = \frac{\frac{1}{N} \sum_{i=1}^N (f_i''(x_i) x_i - f_i'(x_i))}{\frac{1}{N} \sum_{i=1}^N f_i''(x_i)}$$

The initial idea

Derivation of the algorithm - step 3 on 3 - analysis of the problems

Does it work?

1 initialization:

- $y_i(0) := f_i''(x_i(0))x_i(0) - f_i'(x_i(0))$
- $z_i(0) := f_i''(x_i(0))$

2 *average consensus* (in \parallel , P doubly stochastic):

- $\mathbf{y}(k+1) = P\mathbf{y}(k)$
- $\mathbf{z}(k+1) = P\mathbf{z}(k)$

3 local updates: $x_i(k+1) = \frac{y_i(k+1)}{z_i(k+1)}$

The initial idea

Derivation of the algorithm - step 3 on 3 - analysis of the problems

Does it work?

1 initialization:

- $y_i(0) := f_i''(x_i(0))x_i(0) - f_i'(x_i(0))$
- $z_i(0) := f_i''(x_i(0))$

2 **average consensus** (in \parallel , P doubly stochastic):

- $\mathbf{y}(k+1) = P\mathbf{y}(k)$
- $\mathbf{z}(k+1) = P\mathbf{z}(k)$

3 local updates: $x_i(k+1) = \frac{y_i(k+1)}{z_i(k+1)}$

No, **must provide 2 little modifications:**

The initial idea

Derivation of the algorithm - step 3 on 3 - analysis of the problems

Does it work?

① initialization:

- $y_i(0) := f_i''(x_i(0))x_i(0) - f_i'(x_i(0))$
- $z_i(0) := f_i''(x_i(0))$

② **average consensus** (in \parallel , P doubly stochastic):

- $\mathbf{y}(k+1) = P\mathbf{y}(k)$
- $\mathbf{z}(k+1) = P\mathbf{z}(k)$

③ local updates: $x_i(k+1) = \frac{y_i(k+1)}{z_i(k+1)}$

No, **must provide 2 little modifications**:

- x_i changes \Rightarrow must track the changing $f_i'(x_i)$ and $f_i''(x_i)$

The initial idea

Derivation of the algorithm - step 3 on 3 - analysis of the problems

Does it work?

- 1 initialization:
 - $y_i(0) := f_i''(x_i(0))x_i(0) - f_i'(x_i(0))$
 - $z_i(0) := f_i''(x_i(0))$
- 2 **average consensus** (in \parallel , P doubly stochastic):
 - $\mathbf{y}(k+1) = P\mathbf{y}(k)$
 - $\mathbf{z}(k+1) = P\mathbf{z}(k)$
- 3 local updates: $x_i(k+1) = \frac{y_i(k+1)}{z_i(k+1)}$

No, **must provide 2 little modifications:**

- x_i changes \Rightarrow must track the changing $f_i'(x_i)$ and $f_i''(x_i)$
- $x_i(k) = \frac{y_i(k)}{z_i(k)}$ too aggressive!! Should make it milder

The complete algorithm

1 tracking:

- $g_i(k) := f_i''(x_i(k))x_i(k) - f_i'(x_i(k))$
- $h_i(k) := f_i''(x_i(k))$

2 **average consensus** (in \parallel , P doubly stochastic):

- $\mathbf{y}(k+1) = P[\mathbf{y}(k) + \mathbf{g}(k) - \mathbf{g}(k-1)]$
- $\mathbf{z}(k+1) = P[\mathbf{z}(k) + \mathbf{h}(k) - \mathbf{h}(k-1)]$

3 local updates: $x_i(k+1) = (1 - \varepsilon)x_i(k) + \varepsilon \frac{y_i(k+1)}{z_i(k+1)}$

The complete algorithm

1 tracking:

- $g_i(k) := f_i''(x_i(k))x_i(k) - f_i'(x_i(k))$
- $h_i(k) := f_i''(x_i(k))$

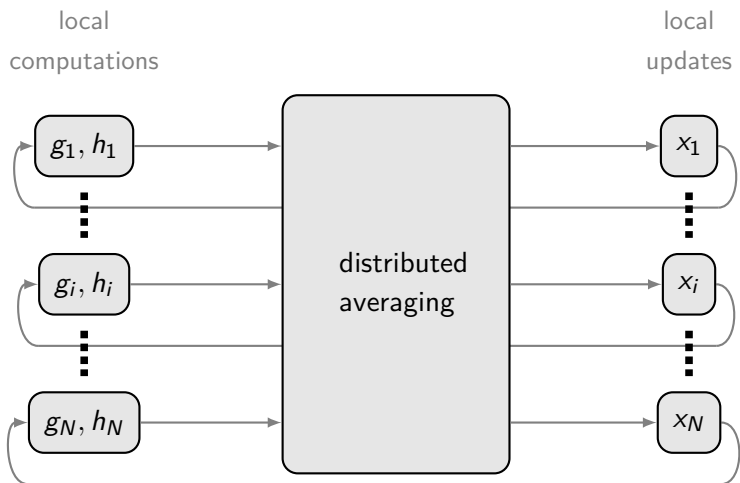
2 **average consensus** (in $\|\cdot\|$, P doubly stochastic):

- $\mathbf{y}(k+1) = P[\mathbf{y}(k) + \mathbf{g}(k) - \mathbf{g}(k-1)]$
- $\mathbf{z}(k+1) = P[\mathbf{z}(k) + \mathbf{h}(k) - \mathbf{h}(k-1)]$

3 local updates: $x_i(k+1) = (1 - \varepsilon)x_i(k) + \varepsilon \frac{y_i(k+1)}{z_i(k+1)}$

(numerical) remark: step 2 may be substituted with asymptotical average consensus algorithms

The complete algorithm – Block schematic representation



$$g_i(k) = f_i''(x_i(k))x_i(k) - f_i'(x_i(k))$$
$$h_i(k) = f_i''(x_i(k))$$

$$x_i(k+1) = (1 - \varepsilon)x_i(k) + \varepsilon \frac{y_i(k+1)}{z_i(k+1)}$$

Convergence properties

Hypotheses

- $f_i \in \mathcal{C}^2(\mathbb{R})$
- f_i' and f_i'' bounded
- f_i strictly convex
- $x^* \neq \pm\infty$
- **null initial conditions** (for g_i, h_i, y_i, z_i)

Thesis

- there exists a positive $\bar{\varepsilon}$ s.t. if $\varepsilon < \bar{\varepsilon}$ then

$$\lim_{k \rightarrow +\infty} \mathbf{x}(k) = x^* \mathbf{1}$$

(convergence \propto as Newton-Raphson strategies over \bar{f})

**importance of the proof:
gives insights on key properties**

- 1 transform the algorithm in a continuous-time system
- 2 recognize the existence of a two-time scales dynamical system
- 3 analyze separately fast and slow dynamics
(singular perturbation methods [Khalil (2002)])

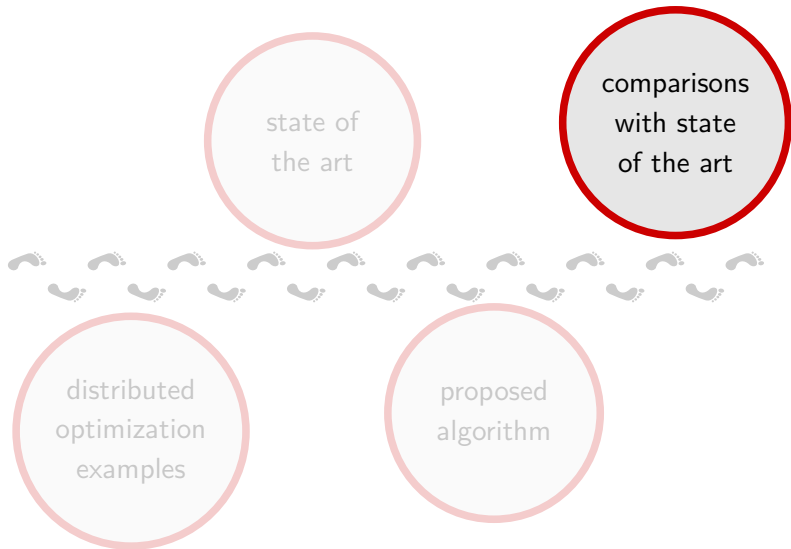
Good qualities

- easy to be implemented
- “small” computational requirements

Bad qualities

Up to now, requires strong assumptions:

- $f_i \in \mathcal{C}^2(\mathbb{R})$
- f_i strictly convex
- f_i' and f_i'' bounded

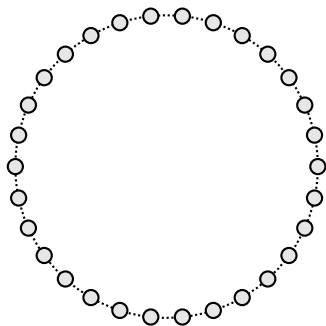


Experiments description

- circulant graph, $N = 30$

- $P = \begin{bmatrix} 0.5 & 0.25 & & & 0.25 \\ 0.25 & 0.5 & 0.25 & & \\ & \ddots & \ddots & \ddots & \\ & & 0.25 & 0.5 & 0.25 \\ 0.25 & & & 0.25 & 0.5 \end{bmatrix}$

- $f_i = \text{sum of exponentials}$



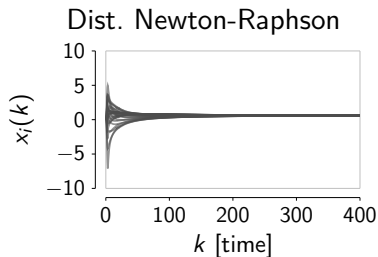
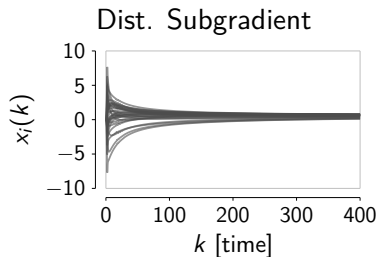
Comparisons with a Distributed Subgradient

Nedić Ozdaglar *Dist. subgr. meth. for multi-agent opt.* (2009)

① $\mathbf{x}^{(c)}(k) = P\mathbf{x}(k)$ (consensus step)

② $x_i(k+1) = x_i^{(c)}(k) - \frac{\rho}{k} f'_i \left(x_i^{(c)}(k) \right)$ (local gradient descent)

Numerical comparison

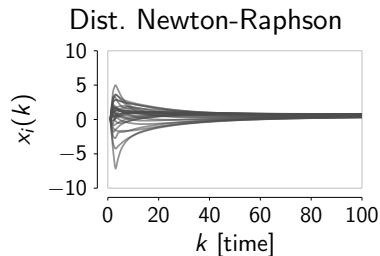
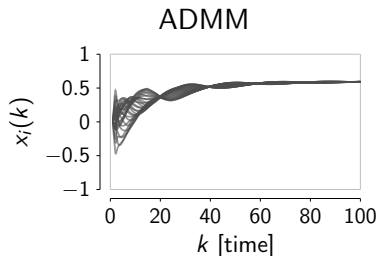


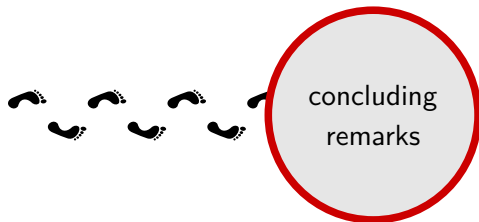
Comparisons with (an) ADMM

Bertsekas Tsitsiklis, *Parall. and Dist. Computation* (1997)

$$L_\rho := \sum_i \left[f_i(x_i) + y_i^{(\ell)}(x_i - z_{i-1}) + y_i^{(c)}(x_i - z_i) + y_i^{(r)}(x_i - z_{i+1}) + \frac{\delta}{2} |x_i - z_{i-1}|^2 + \frac{\delta}{2} |x_i - z_i|^2 + \frac{\delta}{2} |x_i - z_{i+1}|^2 \right]$$

Numerical comparison





The algorithm we proposed ...

- is a distributed Newton-Raphson strategy (+)
- requires minimal network topology knowledge (+)
- requires minimal agents synchronization (+)
- is simple to be implemented (+)
- converges to global optimum under convexity and smoothness assumptions (+ / -)
- is numerically faster than subgradients (+) but slower than ADMM (-)

Currently working on (or already performed)

- extension to multi-dimensional problems
- extension to modified Newton strategies
- analytical characterization of the convergence speed for quadratic functions and specific graphs
(with comparisons to other methods)
- relax the assumptions
(strict convexity, C^2 , ...)
- find automatic stepsizes tuning strategies
- propose quasi-Newton strategies



K. C. Kiwiel (2004)

Convergence of approximate and incremental subgradient methods for convex optimization

SIAM Journal on Optimization



D. P. Bertsekas (1982)

Constrained Optimization and Lagrange Multiplier Methods

Academic Press



D. P. Bertsekas and J. N. Tsitsiklis (1997)

Parallel and Distributed Computation: Numerical Methods

Athena Scientific



A. Nedić and A. Ozdaglar (2009)

Distributed subgradient methods for multi-agent optimization

IEEE Transactions on Automatic Control



B. Johansson (2008)

On Distributed Optimization in Networked Systems

Ph.D. Thesis, KTH



A. Nedić and A. Ozdaglar (2007)

On the Rate of Convergence of Distributed Subgradient Methods for Multi-agent Optimization

CDC



S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein (2010)

Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers

Foundations and Trends in Machine Learning



M. Zargham, A. Ribeiro, A. Ozdaglar, A. Jadbabaie (2011)

Accelerated Dual Descent for Network Optimization

ACC



D. P. Bertsekas (2011)

Centralized and Distributed Newton Methods for Network Optimization and Extensions

Technical Report LIDS 2866



H. K. Khalil (2002)

Nonlinear Systems

Prentice Hall

Distributed convex optimization: a consensus-based Newton-Raphson approach

Damiano Varagnolo

joint work with A. Cenedese, G. Pillonetto, L. Schenato, F. Zanella

Department of Information Engineering - University of Padova

December 14th, 2011 – 50th IEEE CDC

varagnolo@dei.unipd.it
www.dei.unipd.it/~varagnolo/
google: damiano varagnolo