

## Automatic Performance Estimation for Decentralized Optimization

#### Sébastien Colla, Julien Hendrickx

Mathematical Engineering Department, UCLouvain

Seminar at OR center MIT – September 9, 2022



#### **Decentralized Optimization**



#### Decentralization

- $\succ$  Local function:  $f_i$
- $\blacktriangleright$  Local copy of x:  $x_i$

#### **Iterative algorithm**

- Local computations
- Local communications (W) so that  $x_i = x_j$  (eventually)

#### Decentralized Gradient Descent (DGD)



## Motivations: Decentralized Machine Learning

#### Notations

- Model parameters *x*
- Data set  $\{d \in \mathcal{D}\}$

Model training

$$\min_{x} \sum_{d \in D} \operatorname{Error}(x, d) + \operatorname{regul}(x)$$



#### Decentralization

Part of the data  $\mathcal{D}_i$ Local function  $f_i(x) = \sum_{d \in \mathcal{D}_i} \operatorname{Error}(x, d)$ Local copy of x



**Motivations** Big data – Privacy – Speed Up

#### Other applications



## Decentralized Optimization





#### BUT

Analysis highly complex



Performance bounds: complex and conservative

- Difficult algorithms comparisons
- Difficult parameters tuning









**Impact** for decentralized optimization

- Access to accurate performance of methods
- > Easier **comparison and tuning** of algorithms
- Rapid exploration of new algorithms.

## Outline of the talk

- Performance Estimation Problem (PEP)
- PEP for decentralized optimization
- Analysis of Decentralized Algorithms

Idea: Worst-cases are solutions to optimization problems

$$\max_{f, x^0, \dots, x^K} \quad \text{perf}(f, x^0, \dots, x^K) \stackrel{e.g.}{=} f(x^K) - f(x^*)$$
With  $f \in \text{class of functions}$ 
 $x^0 \quad \text{initial condition}$ 
 $x^k \quad \text{from the algorithm analyzed}$ 

Idea: Worst-cases are solutions to optimization problems

$$\max_{f, x^0, ..., x^K} \quad \text{perf}(f, x^0, ..., x^K) \stackrel{e.g.}{=} f(x^K) - f(x^*)$$

With	f	E	$\mathcal{F}_{\mu,L}$	class of functions
	$  x^{0} $	$-x^* \parallel$	$\leq R$	initial condition
	$x^k =$	= x <sup>k-1</sup>	$-\alpha \nabla f(x^{k-1})$	from the algorithm analyzed

Original idea by	Yoel Drori, and Marc Teboulle (2014)
Further developments by	Adrien B. Taylor, Julien M. Hendrickx, and François Glineur (2017)

Idea: Worst-cases are solutions to optimization problems



[Taylor et al. 2017]



Can be used for tuning, design and proofs.

#### Finite dimensional PEP

Finite dimension

 ${x^k, g^k, f^k}_{k=1...K}$ 





Interpolation conditions for many classical function classes e.g., for  $\mathcal{F}_{\mu}$ ,  $f_j \ge f_i + g_i^T (x_j - x_i) + \frac{\mu}{2} ||x_j - x_i||^2$  for all (i, j) [Taylor et al. 2017]

## SDP formulation of PEP

PEP constraints may be quadratic and non-convex SDP reformulation

Variables  $F = [f^0 \dots f^K]$  $G = P^T P \qquad P = [x^0 \dots x^K g^0 \dots g^K]$ Gram Matrix perf(F,G) PEP max *F, G* **Efficient resolution** With  $G \geq 0$ Interpolation Initial constraints linear in G and FAlgorithm

## Outline of the talk

- Performance Estimation Problem (PEP)
- **PEP** for **decentralized** optimization
- Analysis of Decentralized Algorithms

Idea: Worst-cases are solutions to optimization problems

$\max_{\substack{f_i, x_i^0, \dots, x_i^K, W\\(i = 1 \dots N)}} \operatorname{perf}(f_i, x_i^0, \dots, x_i^K)$					
With	f <sub>i</sub>	E	class of functions		
	$x_i^0$		initial condition		
	$x_i^k$		from the algorithm analyzed		
	W	E	class of network matrices		



How to represent a class of communication network matrices ?

#### PEP for DGD: network given a priori





#### PEP for DGD: class of networks



Iterates from DGD 
$$\begin{cases} y_i^k = \sum_{j} w_{ij} x_j^k \\ x_i^{k+1} = y_i^k - \alpha \nabla f_i(x_i^k) \end{cases}$$
For all  $i = 1 \dots N$ ,  
For all  $k = 0 \dots K - 1$ 

#### PEP for DGD: class of networks



#### Consensus steps in PEP

 $\succ$  Search Space for X and Y

(C1) 
$$y_i^k = \sum_{j=1}^N w_{ij} x_j^k$$
 For each agent  $i = 1 \dots N$ ,  
For each consensus step  $k = 0 \dots K - 1$   
 $Y = WX$  with  $Y_{ik} = y_i^k$ ,  $X_{ik} = x_i^k$ .

(C2)  $W = [w_{ij}]$  is a symmetric and doubly-stochastic matrix with a given range of eigenvalues  $[\lambda^{-}, \lambda^{+}]$ 

Necessary constraints for describing (C1) and (C2)

 $\bar{X}, \bar{Y}: \text{ agents average vectors} \qquad \begin{array}{l} X_{\perp}, Y_{\perp}: \text{ centered matrices} \\ X_{\perp} = X - \mathbf{1}\bar{X}^{T}, \ Y_{\perp} = Y - \mathbf{1}\bar{Y}^{T} \\ \hline \bar{X} = \bar{Y} \qquad (1) \\ \begin{pmatrix} \lambda^{-} X_{\perp}^{T} X_{\perp} \leqslant X_{\perp}^{T} Y_{\perp} \leqslant \lambda^{+} X_{\perp}^{T} X_{\perp} & (2) \\ (Y_{\perp} - \lambda^{-} X_{\perp})^{T} (Y_{\perp} - \lambda^{+} X_{\perp}) \leqslant 0 & (3) \end{array}$ Simplification of (2) and (3) when  $-\lambda^{-} = \lambda^{+} = \lambda$ :  $Y_{\perp}^{T} Y_{\perp} \leqslant \lambda^{2} X_{\perp}^{T} X_{\perp}$ 

#### Consensus steps in PEP

Summary of the constraints for **consensus steps** Y = WX

$$\overline{X} = \overline{Y}$$
(1)  

$$\lambda^{-} X_{\perp}^{T} X_{\perp} \leq X_{\perp}^{T} Y_{\perp} \leq \lambda^{+} X_{\perp}^{T} X_{\perp}$$
(2)  

$$(Y_{\perp} - \lambda^{-} X_{\perp})^{T} (Y_{\perp} - \lambda^{+} X_{\perp}) \leq 0$$
(3)

Advantages of our constraints

- Independent of the algorithm
- Link different consensus steps that use the same matrix
- ✓ Can be incorporated into SDP formulation of PEP, which can be solved efficiently

#### PEP for DGD: class of networks



21

## PEP for DGD: Spectral formulation (Relaxation)

$$\max_{\substack{f_i, x^0, \dots, x^K \\ y^0, \dots, y^{K-1}}} \operatorname{perf}(f_i, x^0, \dots, x^K)$$

$$\underset{y^0, \dots, y^{K-1}}{\operatorname{With}} \quad \begin{array}{l} f_i \in \operatorname{class of functions} \\ x^0 & \operatorname{initial condition} \end{array}$$

$$\operatorname{Iterates from DGD} - \begin{bmatrix} x_i^{k+1} = y_i^k - \alpha \nabla f_i(x_i^k) & \operatorname{For all } i = 1 \dots N, \\ x_i^{k+1} = y_i^k - \alpha \nabla f_i(x_i^k) & \operatorname{For all } k = 0 \dots K - 1 \end{aligned}$$

$$\operatorname{Consensus steps}_{\substack{Y = WX \\ W \text{ doubly stochastic} \\ \lambda(W) \in [\lambda^-, \lambda^+]}} - \begin{bmatrix} \overline{X} = \overline{Y} & & \\ \lambda^- X_{\perp}^T X_{\perp} \leqslant X_{\perp}^T Y_{\perp} \leqslant \lambda^+ X_{\perp}^T X_{\perp} \\ (Y_{\perp} - \lambda^- X_{\perp})^T (Y_{\perp} - \lambda^+ X_{\perp}) \leqslant 0 \end{bmatrix}} \xrightarrow{Notations}_{\substack{X_{ik} = x_i^k \\ Y_{ik} = y_i^k}}$$

Upper bounds for the worst-case performance of DGD

#### Our tool for automatic performance estimation

Apply to any decentralized method using consensus y = Wx



Agent-independent (spectral) formulation

### Agent-Independent formulation

Local view

$$\min_{\substack{x_1, \dots, x_N \\ \text{s.t. } x_i = x_j \quad \forall \ (i, j) \text{ neighbors}}} \sum_{i=1}^N f_i(x_i) \qquad x_i \in \mathbb{R}^d, \qquad i = 1 \dots N$$

**Global view** 

$$\begin{array}{ccc} & \mathbf{x} \in \mathbb{R}^{Nd}, & \mathsf{F} \colon \mathbb{R}^{Nd} \to \mathbb{R} \\ & \underset{\mathbf{X}}{\min} & F(\mathbf{x}) \\ & \text{s.t. } \mathbf{x} \in \mathbb{C} \\ & & \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} \in \mathcal{C} \Leftrightarrow x_1 = \cdots = x_N \end{array}$$

**Change of variable** to decouple C and  $C^{\perp}$ 

## Outline of the talk

- Performance Estimation Problem (PEP)
- **PEP** for **decentralized** optimization
- Analysis of Decentralized Algorithms

#### Results of PEP for DGD

#### Problem

 $\min_{x} f(x) = \frac{1}{N} \sum_{i=1}^{N} f_i(x) \quad \text{with optimal solution } x^*$ 

**DGD Algorithm** 
$$x_i^{k+1} = \sum_{j=1}^N w_{ij} x_j^k - \alpha \nabla f_i(x_i^k)$$

**Settings** *K* steps of DGD with

- Constant step-size:  $\alpha = \frac{1}{\sqrt{K}}$
- Convex local functions *f<sub>i</sub>* with bounded subgradients
- Identical starting points s.t.  $||x^0 x^*||^2 \le 1$
- Symmetric doubly-stochastic network matrix W s.t.  $\lambda(W) \in [-\lambda, \lambda]$  (except for  $\lambda_1(W) = 1$ )

Performance criterion:  $f(x_{av}) - f(x^*)$  where  $x_{av} = \frac{1}{\kappa} \sum_k \frac{1}{N} \sum_i x_i^k$ 

# DGD – Spectral worst-case evolution with N



For K = 5 iterations and  $\lambda(W) \in [-\lambda, \lambda]$ 

#### DGD – Spectral worst-case vs Theoretical bound



For K = 10 iterations, N = 3 agents and  $\lambda(W) \in [-\lambda, \lambda]$ 

#### **Tightness Analysis**



For K = 10 iterations, N = 3 agents and  $\lambda(W) \in [-\lambda, \lambda]$ 

#### **Tightness Analysis**



For K = 10 iterations, N = 3 agents and  $\lambda(W) \in [-\lambda, \lambda]$ 

# DGD – Spectral worst-case evolution with alpha



For K = 10 iterations and  $\lambda(W) \in [-\lambda, \lambda]$ 

# DGD – Spectral worst-case evolution with alpha



For K = 10 iterations and  $\lambda(W) \in [-\lambda, \lambda]$ 

#### Problem

$$\min_{x} f(x) = \frac{1}{N} \sum_{i=1}^{N} f_i(x)$$

with optimal solution  $x^*$ 

#### DIGing Algorithm

gradient tracking technique

$$x_{i}^{k+1} = \sum_{j=1}^{N} w_{ij} x_{j}^{k} - \alpha s_{i}^{k}$$
$$s_{i}^{k+1} = \sum_{j=1}^{N} w_{ij} s_{j}^{k} + \nabla f_{i}(x_{i}^{k+1}) - \nabla f_{i}(x_{i}^{k})$$

**Settings** •  $f_i$  are *L*-smooth and  $\mu$ -strongly convex

• Initial: 
$$\frac{1}{N} \sum_{i=1}^{N} \left\| x_i^0 - x^* \right\|^2 \le 1$$
 and  $\frac{1}{N} \sum_{i=1}^{N} \left\| s_i^0 - \overline{\nabla f_i^0} \right\|^2 \le 1$ 

■ Symmetric doubly-stochastic network matrix W s.t.  $\lambda(W) \in [-\lambda, \lambda]$  (except for  $\lambda_1(W) = 1$ )

Performance criterion:

$$\frac{1}{N} \sum_{i=1}^{N} \left\| x_i^K - x^* \right\|^2$$

Spectral PEP formulation

Independent of the number of agents N

Same worst-case matrix than DGD

$$W_1 = \Pi + \lambda(\Pi - I)$$
 with  $\Pi = \frac{1}{N} \mathbf{1} \mathbf{1}^T$ 

$$\Rightarrow \quad \lambda(W_1) = \{1, -\lambda, \dots, -\lambda\}$$

**Exact** for spectrally doubly-stochastic matrices



For 
$$\mu = 0.1$$
,  $L = 1$  and  $\lambda(W) \in [-0.9, 0.9]$   
Computed for  $N = 2$ .

[NOS16] A. Nedic, A. Olshevsky, and W. Shi, "Achieving geometric convergence for distributed optimization over time-varying graphs," SIAM Journal on Optimization, 2016.



For 
$$\mu = 0.1$$
,  $L = 1$  and  $\lambda(W) \in [-0.9, 0.9]$   
Computed for  $N = 2$ .

[NOS16] A. Nedic, A. Olshevsky, and W. Shi, "Achieving geometric convergence for distributed optimization over time-varying graphs," SIAM Journal on Optimization, 2016.

#### Algorithms comparison



For K = 10 iterations,  $\mu = 0.1$ , L = 1 and  $\lambda(W) \in [-\lambda, \lambda]$ .

## Conclusion





## Numerical tool for **automatic performance computation** of decentralized optimization methods

#### PEP idea: worst-cases are solutions of optimization problems

	SPECTR	AL formulation	EXACT formulation	
	Spectral	class of matrices	Given network matrix W	
	r Relax	ation of PEP	ALWAYS exact	
For DGD a	nd DIGing:	<ul> <li>✓ Independent of N</li> <li>✓ Tight when negative weights are allowed</li> <li>✓ Improve on the literature bounds</li> </ul>		

We can answer a large diversity of (new) questions

#### **Future works**

- □ Other class of networks (any suggestion?)
- □ Agent-independent PEP formulation

### References

[CH21] S. Colla, J. M. Hendrickx, "Automatic Performance Estimation for Decentralized Optimization", preprint 2022.

- [Taylor et al.] A. B. Taylor, J. M. Hendrickx, F. Glineur, "Exact worst-case performance of first-order methods for composite convex optimization," SIAM Journal on Optimization, 2015
- [NOR17] A. Nedic, A. Olshevsky, and M. G. Rabbat, "Network topology and communication computation tradeoffs in decentralized optimization", 2017.
- [NOS16] A. Nedic, A. Olshevsky, and W. Shi, "Achieving geometric convergence for distributed optimization over time-varying graphs," SIAM Journal on Optimization, 2016.