

# Automatic Performance Estimation for Decentralized Optimization

Sébastien Colla, Julien Hendrickx



$$\min_{x} f(x) = \sum_{i} f_i(x)$$



#### Decentralization

 $\succ$  Local function:  $f_i$ 

$$\min_{x} f(x) = \sum_{i} f_i(x)$$



#### Decentralization

 $\succ$  Local function:  $f_i$ 

#### **Iterative algorithm**

> Local computations

$$\min_{x} f(x) = \sum_{i} f_i(x)$$



#### Decentralization

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

#### **Iterative algorithm**

Local computations



#### Decentralization

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

#### **Iterative algorithm**

Local computations



#### Decentralization

- $\blacktriangleright$  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  $\{d_k \in \mathcal{D}\}$

Model training

$$\min_{x} \sum_{k} \operatorname{Error}(x, d_{k}) + \operatorname{regul}(x)$$



#### Decentralization

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



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

#### Other applications





Many challenges for better methods





Many challenges for better methods

#### BUT

Analysis highly complex





Many challenges for better methods

BUT

Analysis highly complex



Design: long and complex process

Performance bounds: complex and conservative



Many challenges for better methods

BUT

Analysis highly complex



- Design: long and complex process
- Performance bounds: complex and conservative
- Difficult algorithms comparisons
- Difficult parameters tuning









**Impact** for decentralized optimization

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

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

$$x_0 \quad \text{initial condition}$$

$$x^k \quad \text{from the algorithm analyzed}$$

$$\max_{\substack{f, x^0, \dots, x^K}} \operatorname{perf}(f, x^0, \dots, x^K) \stackrel{e.g.}{=} f(x^K) - f(x^*)$$
With  $f \in \operatorname{class of functions}$ 
Infinite-Dimensional  $x_0$  initial condition
problem  $x^k$  from the algorithm analyzed







## PEP for DGD: network given a priori

## PEP for DGD: network given a priori

$$\max_{\substack{f, x^0, \dots, x^{K}, G \\ y^0, \dots, y^{K-1}}} \operatorname{perf}(f, x^0, \dots, x^K)$$
With  $f \in \operatorname{class of functions}$ 

$$x_0 \qquad \text{initial condition}$$
 $W(G) \qquad \text{given network matrix}$ 
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: network given a priori

$$\max_{\substack{f, x^0, \dots, x^{K}, G \\ y^0, \dots, y^{K-1}}} \operatorname{perf}(f, x^0, \dots, x^K)$$
With  $f \in \operatorname{class of functions}$ 

$$x_0 \qquad \text{initial condition}$$
 $W(G) \qquad \text{given network matrix}$ 
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  $k = 0 \dots K - 1$ 



$$\max_{\substack{f, x^0, \dots, x^K, \mathbf{G} \\ y^0, \dots, y^{K-1}}} \operatorname{perf}(f, x^0, \dots, x^K)$$
With  $f \in \operatorname{class of functions}$ 

$$x_0 \qquad \text{initial condition}$$

$$W(G) \qquad \operatorname{Any \, symmetric \, doubly \, stochastic \, matrix}_{with \, given \, range \, of \, eigenvalues} [\lambda^-, \lambda^+]$$

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

$$\max_{\substack{f, x^0, \dots, x^K, \mathbf{G} \\ y^0, \dots, y^{K-1}}} \operatorname{perf}(f, x^0, \dots, x^K)$$
With  $f \in \operatorname{class of functions}$ 

$$x_0 \qquad \text{initial condition}$$

$$W(G) \qquad \operatorname{Any symmetric doubly stochastic matrix}_{with given range of eigenvalues} [\lambda^-, \lambda^+]$$

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  $k = 0 \dots K - 1$ 

$$\max_{\substack{f, x^0, \dots, x^K, G \\ y^0, \dots, y^{K-1}}} \operatorname{perf}(f, x^0, \dots, x^K)$$
With  $f \in \operatorname{class of functions}$ 

$$x_0 \qquad \text{initial condition}$$

$$W(G) \qquad \operatorname{Any symmetric doubly stochastic matrix}_{with given range of eigenvalues} [\lambda^-, \lambda^+]$$

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

> Search Space for  $x_i^k$  and  $y_i^k$ 

N T

(C1) 
$$y_i^k = \sum_{j=1}^N w_{ij} x_j^k$$
 For each agent  $i = 1 \dots N$ ,  
For each consensus step  $k = 0 \dots K - 1$ 

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

 $\succ$  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  $k = 0 \dots K - 1$   
 $(C1) Y = WX$  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  $[\lambda^-, \lambda^+]$ 

 $\succ$  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  $k = 0 \dots K - 1$   
 $(C1) Y = WX$  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)

 $\overline{X}, \overline{Y}: \text{ agents average vectors}$   $\begin{array}{l}
X_{\perp}, Y_{\perp}: \text{ centered matrices} \\
X_{\perp} = X - \mathbf{1}\overline{X}^{T}, \ Y_{\perp} = Y - \mathbf{1}\overline{Y}^{T} \\
\overline{X} = \overline{Y}$   $\begin{array}{l}
\widehat{X} = \overline{Y} \\
\lambda^{-} X_{\perp}^{T} X_{\perp} \leqslant X_{\perp}^{T} Y_{\perp} \leqslant \lambda^{+} X_{\perp}^{T} X_{\perp}$   $\begin{array}{l}
(1) \\
\chi^{-} X_{\perp}^{T} X_{\perp} \leqslant X_{\perp}^{T} Y_{\perp} \leqslant \lambda^{+} X_{\perp}^{T} X_{\perp}$   $\begin{array}{l}
(2) \\
(Y_{\perp} - \lambda^{-} X_{\perp})^{T} (Y_{\perp} - \lambda^{+} X_{\perp}) \leqslant 0$   $\begin{array}{l}
X_{\perp}, Y_{\perp}: \text{ centered matrices} \\
X_{\perp} = X - \mathbf{1}\overline{X}^{T}, \ Y_{\perp} = Y - \mathbf{1}\overline{Y}^{T}$   $\begin{array}{l}
(1) \\
(2) \\
(3)
\end{array}$ 

 $\succ$  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  $k = 0 \dots K - 1$   
 $(C1) Y = WX$  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)

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

Summary of the constraints for consensus steps between Y and X

$$\overline{X} = \overline{Y}$$

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

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

$$(3)$$

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

$$\max_{\substack{f, x^0, \dots, x^K, G \\ y^0, \dots, y^{K-1}}} \operatorname{perf}(f, x^0, \dots, x^K)$$
With  $f \in \operatorname{class of functions}$ 

$$x_0 \qquad \text{initial condition}$$
W(G) Any symmetric doubly stochastic matrix with given range of eigenvalues  $[\lambda^-, \lambda^+]$ 

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

32



Iterates from DGD

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

K-1

## PEP for DGD: Spectral formulation (Relaxation)

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

$$\underset{y^0, \dots, y^{K-1}}{\operatorname{With}} f \in \operatorname{class of functions}$$

$$x_0 \quad \operatorname{initial condition}$$
Iterates from DGD  $\left\{ x_i^{k+1} = y_i^k - \alpha \nabla f_i(x_i^k) \quad \begin{array}{c} \operatorname{For all} i = 1 \dots N, \\ \operatorname{For all} k = 0 \dots K - 1 \\ \\ \operatorname{Consensus steps} \\ Y = WX \\ W \quad \begin{array}{c} \sum_{\substack{Y = WX \\ \text{doubly stochastic} \\ \lambda(W) \in [\lambda^-, \lambda^+]} \end{array} \right| \left[ \begin{array}{c} \sum_{\substack{X = \overline{Y} \\ Y = X_\perp}} \sum_{\substack{X = \overline{Y} \\ Y = X_\perp}} \\ \left[ \sum_{\substack{X = \overline{Y} \\ Y = X_\perp}} \sum_{\substack{X = \overline{Y} \\ Y = X_\perp}} \right] \\ \end{array} \right]$ 

## PEP for DGD: Spectral formulation (Relaxation)

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

$$\sum_{\substack{y^0, \dots, y^{K-1} \\ With \ f \in class of functions \\ x_0 \qquad \text{initial condition}}$$

$$\max_{\substack{x_0 \\ x_0 \qquad \text{initial condition}}} \operatorname{For all } i = 1 \dots N,$$

$$\operatorname{For all } k = 0 \dots K - 1$$

$$\operatorname{Consensus steps}_{\substack{Y = WX \\ y \ \text{symmetric} \\ \lambda(W) \in [\lambda^-, \lambda^+]}} - \left[ \begin{array}{c} \overline{X} = \overline{Y} & (1) \\ \lambda^- X_{\perp}^T X_{\perp} \leqslant X_{\perp}^T Y_{\perp} \leqslant \lambda^+ X_{\perp}^T X_{\perp} & (2) \\ (Y_{\perp} - \lambda^- X_{\perp})^T (Y_{\perp} - \lambda^+ X_{\perp}) \leqslant 0 & (3) \end{array} \right]$$



Upper bounds for the worst-case performance of DGD

#### Theoretical guarantee for DGD [NOR17]

After K steps of DGD with  $\alpha = \frac{h}{\sqrt{K}}$ , for solving  $\min_{x} f(x) = \sum_{i} f_{i}(x)$ With optimal solution  $x^*$ 

lf

(i)  $f_i$  convex for each *i* and (sub)gradients bounded by *B*;

(ii) Identical starting points:  $x_i^0 = x^0$  for each *i*, s.t.  $||x^0 - x^*||^2 \le R^2$ (iii) Communication matrix W: symmetric, doubly stochastic,

$$\lambda(W) \in [-\lambda, \lambda]$$
 (except for  $\lambda_1(W) = 1$ )

Then

$$f(x_{av}) - f(x^*) \le RB\left(\frac{h^{-1} + h}{2\sqrt{K}} + \frac{2h}{\sqrt{K}(1-\lambda)}\right) \quad \text{where} \quad x_{av} = \frac{1}{K}\sum_k \frac{1}{N}\sum_i x_i^k$$

#### Theoretical guarantee for DGD [NOR17]

After K steps of DGD with  $\alpha = \frac{h}{\sqrt{\kappa}}$ , for solving  $\min_{x} f(x) = \sum_{i} f_{i}(x)$ With optimal solution  $x^*$ 

lf

(i)  $f_i$  convex for each *i* and (sub)gradients bounded by *B*;

(ii) Identical starting points:  $x_i^0 = x^0$  for each *i*, s.t.  $||x^0 - x^*||^2 \le R^2$ 

(iii) Communication matrix W:

symmetric, doubly stochastic,  $\lambda(W) \in [-\lambda, \lambda]$  (except for  $\lambda_1(W) = 1$ )

Then

$$\frac{f(x_{av}) - f(x^*)}{k} \le RB\left(\frac{h^{-1} + h}{2\sqrt{K}} + \frac{2h}{\sqrt{K}(1-\lambda)}\right) \quad \text{where} \quad x_{av} = \frac{1}{K}\sum_k \frac{1}{N}\sum_i x_i^k$$
Performance measure

Same settings in our experiments, with R = 1, B = 1 and h = 1.

37

## DGD – Spectral worst-case evolution with N



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

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



Symmetric range of eigenvalues  $-\lambda \leq \lambda_n(W) \leq \cdots \leq \lambda_2(W) \leq \lambda$ 

For K = 10 iterations, N = 3 agents.

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



For K = 10 iterations, N = 3 agents.



$$W_{1} = \begin{bmatrix} d & c & c \\ c & d & c \\ c & c & d \end{bmatrix}$$
  
with  $c = \frac{1+\lambda}{3}$ ,  $d = 1 - 2c$   
 $\Rightarrow \lambda(W_{1}) = \{1, -\lambda, -\lambda\}$ 

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



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



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



$$W_{1} = \begin{bmatrix} d & c & c \\ c & d & c \\ c & c & d \end{bmatrix}$$
  
with  $c = \frac{1+\lambda}{3}, d = 1 - 2c$   
 $\Rightarrow \lambda(W_{1}) = \{1, -\lambda, -\lambda\}$ 

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



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

#### DGD: Step-Size tuning



For K = 10 iterations, N = 3 agents with  $\lambda = 0.8$  and  $\alpha = \frac{h}{\sqrt{K}}$ 

#### DGD: Step-Size tuning



For K = 10 iterations, N = 3 agents with  $\lambda = 0.8$  and  $\alpha = \frac{h}{\sqrt{K}}$ 

47

#### DGD: Step-Size tuning



For K = 10 iterations, N = 3 agents with  $\lambda = 0.8$  and  $\alpha = \frac{h}{\sqrt{K}}$ 

## **Automatic** tool for accurate **performance estimation** of decentralized optimization methods



## **Automatic** tool for accurate **performance estimation** of decentralized optimization methods



SPECTRAL formulation	EXACT formulation
Spectral class of matrices	Given network matrix W
Relaxation of PEP	ALWAYS exact

## **Automatic** tool for accurate **performance estimation** of decentralized optimization methods



SPECTRAL formulation	EXACT formulation
Spectral class of matrices	Given network matrix W
Relaxation of PEP	ALWAYS exact

For DGD:	$\checkmark$	Independent of N
----------	--------------	------------------

- ✓ Tight with potential negative weights
- ✓ Accurate with nonnegative weights
- ✓ Improve on the literature bound

## **Automatic** tool for accurate **performance estimation** of decentralized optimization methods



PEP idea: worst-cases are solutions of optimization problems

#### Future works

□ Implementation of the formulation in **PESTO toolbox** (in progress)

□ Strong theoretical understanding of our formulation

□ Analyze other decentralized algorithms using our tool

#### References

- [CH21] S. Colla, J. M. Hendrickx, "Automated Worst-Case Performance Analysis of Decentralized Gradient Descent", 2021.
- [Taylor17] A. B. Taylor, "Convex interpolation and performance estimation of firstorder methods for convex optimization", PhD, UCLouvain, Louvain-la-Neuve, Belgium, 2017.
- [NOR17] A. Nedic, A. Olshevsky, and M. G. Rabbat, "Network topology and communication computation tradeoffs in decentralized optimization", 2017.