

# Linear Convergence for Distributed Optimization Without Strong Convexity

Xinlei Yi

Joint work with Shengjun Zhang, Tao Yang, Tianyou Chai, and Karl H. Johansson

December, 2020

School of Electrical Engineering and Computer Science KTH Royal Institute of Technology Stockholm, Sweden



A network of agents cooperatively solve a global optimization problem, where

• each agent i has a local private objective  $f_i(x)$ 

A network of agents cooperatively solve a global optimization problem, where

- each agent i has a local private objective  $f_i(\boldsymbol{x})$
- all agents collaborate together to find the solution to minimize  $f(x) := \frac{1}{n} \sum_{i=1}^{n} f_i(x)$

$$(x) := \frac{1}{n} \sum_{i=1}^{n} f_i(x)$$

 $\min_{x \in \mathbb{R}^p} f$ 





A network of agents cooperatively solve a global optimization problem, where

- each agent i has a local private objective  $f_i(x)$
- all agents collaborate together to find the solution to minimize  $f(x) := \frac{1}{n} \sum_{i=1}^{n} f_i(x)$
- $\bullet$  agents exchange information through the underlying communication network  ${\cal G}$





A network of agents cooperatively solve a global optimization problem, where

- each agent i has a local private objective  $f_i(x)$
- all agents collaborate together to find the solution to minimize  $f(x) := \frac{1}{n} \sum_{i=1}^{n} f_i(x)$
- $\bullet$  agents exchange information through the underlying communication network  ${\cal G}$



- Distributed optimization summarizes many popular machine learning models, e.g., deep learning and federated learning [Dean et al, NeurIPS, 2012]
- Distributed algorithms outperform centralized algorithms in some applications, e.g., training neural networks [Lian et al, NeurIPS, 2017]



Existing algorithms Continuous- and discrete-time distributed algorithms



Existing algorithms

Continuous- and discrete-time distributed algorithms

Existing result

A standard assumption for proving exponential/linear convergence of existing distributed algorithms is strong convexity of the cost functions



Existing algorithms

Continuous- and discrete-time distributed algorithms

Existing result A standard assumption for proving exponential/linear convergence of existing distributed algorithms is strong convexity of the cost functions

Question

Could strong convexity be relaxed?

For example, quadratic functions may be not strongly convex.



Existing algorithms

Continuous- and discrete-time distributed algorithms

Existing result A standard assumption for proving exponential/linear convergence of existing distributed algorithms is strong convexity of the cost functions

Question

Could strong convexity be relaxed?

For example, quadratic functions may be not strongly convex.

Answer in our paper

Yes, it can be relaxed by the Polyak-Łojasiewicz condition.



Polyak–Łojasiewicz (P–Ł) condition The function f(x) satisfies the P–Ł condition with constant  $\nu > 0$  if  $\|\nabla f(x)\|^2 \ge \nu(f(x) - f^*), \ \forall x \in \mathbb{R}^p.$ 

KTH VETRAGE

Polyak–Łojasiewicz (P–Ł) condition The function f(x) satisfies the P–Ł condition with constant  $\nu > 0$  if  $\|\nabla f(x)\|^2 \ge \nu(f(x) - f^*), \ \forall x \in \mathbb{R}^p.$ 

• Examples:  $f(x) = ||Ax||^2$  and  $f(x) = x^2 + 3\sin^2(x)$ .

KTH

Polyak–Łojasiewicz (P–Ł) condition The function f(x) satisfies the P–Ł condition with constant  $\nu > 0$  if  $\|\nabla f(x)\|^2 \ge \nu(f(x) - f^*), \ \forall x \in \mathbb{R}^p.$ 

- Examples:  $f(x) = ||Ax||^2$  and  $f(x) = x^2 + 3\sin^2(x)$ .
- P-Ł does not imply convexity.
- P-Ł implies that every stationary point is a global minimizer.
- Every strongly convex function satisfies P-Ł.
- It is difficult to verify P-L in general.

KTH

Polyak–Łojasiewicz (P–Ł) condition The function f(x) satisfies the P–Ł condition with constant  $\nu > 0$  if  $\|\nabla f(x)\|^2 \ge \nu(f(x) - f^*), \ \forall x \in \mathbb{R}^p.$ 

- Examples:  $f(x) = ||Ax||^2$  and  $f(x) = x^2 + 3\sin^2(x)$ .
- P-Ł does not imply convexity.
- P-Ł implies that every stationary point is a global minimizer.
- Every strongly convex function satisfies P-Ł.
- It is difficult to verify P-L in general.

#### One special case

If  $g : \mathbb{R}^p \to \mathbb{R}$  is a strongly convex function and  $A \in \mathbb{R}^{p \times p}$  is a matrix, then f(x) = g(Ax) satisfies the P-Ł condition.



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



$$\min_{x\in \mathbb{R}^p} f(x) := rac{1}{n} \sum_{i=1}^n f_i(x)$$
 is equivalent to

$$\min_{\boldsymbol{x}\in\mathbb{R}^{np}}\tilde{f}(\boldsymbol{x}):=\sum_{i=1}^{n}f_{i}(x_{i})$$

n

subject to  $L^{1/2}x = \mathbf{0}_{np}$ ,

where  $\boldsymbol{x} = \operatorname{col}(x_1, \ldots, x_n)$  and  $\boldsymbol{L} = L \otimes \mathbf{I}_p$  with L being the Laplacian.

$$\min_{x\in\mathbb{R}^p}f(x):=rac{1}{n}\sum_{i=1}^n f_i(x)$$
 is equivalent to

$$\min_{\boldsymbol{x}\in\mathbb{R}^{n_p}}\tilde{f}(\boldsymbol{x}):=\sum_{i=1}f_i(x_i)$$

n

subject to  $L^{1/2}x = \mathbf{0}_{np},$ 

where  $x = col(x_1, ..., x_n)$  and  $L = L \otimes I_p$  with L being the Laplacian. Associated augmented Lagrangian:

$$\mathcal{A}(\boldsymbol{x}, \boldsymbol{u}) = \tilde{f}(\boldsymbol{x}) + \frac{lpha}{2} \boldsymbol{x}^{\top} \boldsymbol{L} \boldsymbol{x} + eta \boldsymbol{u}^{\top} \boldsymbol{L}^{1/2} \boldsymbol{x},$$

where  $\boldsymbol{u} \in \mathbb{R}^{np}$  is the dual variable,  $\alpha, \beta > 0$  are parameters.

$$\min_{x\in \mathbb{R}^p} f(x) := rac{1}{n} \sum_{i=1}^n f_i(x)$$
 is equivalent to

$$\min_{\boldsymbol{x}\in\mathbb{R}^{n_p}}\tilde{f}(\boldsymbol{x}):=\sum_{i=1}f_i(x_i)$$

n

subject to  $L^{1/2}x = \mathbf{0}_{np}$ ,

where  $x = col(x_1, ..., x_n)$  and  $L = L \otimes I_p$  with L being the Laplacian. Associated augmented Lagrangian:

$$\mathcal{A}(\boldsymbol{x}, \boldsymbol{u}) = \tilde{f}(\boldsymbol{x}) + \frac{lpha}{2} \boldsymbol{x}^{\top} \boldsymbol{L} \boldsymbol{x} + eta \boldsymbol{u}^{\top} \boldsymbol{L}^{1/2} \boldsymbol{x},$$

where  $u \in \mathbb{R}^{np}$  is the dual variable,  $\alpha, \beta > 0$  are parameters. Minimize  $\mathcal{A}(x, u)$  with a primal-dual algorithm:

$$\begin{split} \boldsymbol{x}_{k+1} = & \boldsymbol{x}_k - \eta \frac{\partial \mathcal{A}(\boldsymbol{x}_k, \boldsymbol{u}_k)}{\partial \boldsymbol{x}_k} = \boldsymbol{x}_k - \eta (\alpha \boldsymbol{L} \boldsymbol{x}_k + \beta \boldsymbol{L}^{1/2} \boldsymbol{u}_k + \nabla \tilde{f}(\boldsymbol{x}_k)), \\ & \boldsymbol{u}_{k+1} = & \boldsymbol{u}_k + \eta \frac{\partial \mathcal{A}(\boldsymbol{x}_k, \boldsymbol{u}_k)}{\partial \boldsymbol{u}_k} = \boldsymbol{u}_k + \eta \beta \boldsymbol{L}^{1/2} \boldsymbol{x}_k, \end{split}$$
where  $\eta > 0$  is a fixed stepsize.

$$\min_{x\in\mathbb{R}^p}f(x):=\frac{1}{n}\sum_{i=1}^n f_i(x) \quad \text{ is equivalent to}$$

$$\min_{\boldsymbol{x}\in\mathbb{R}^{n_p}}\tilde{f}(\boldsymbol{x}) := \sum_{i=1}^n f_i(x_i)$$

n

where  $x = \operatorname{col}(x_1, \dots, x_n)$  and  $L = L \otimes \mathbf{I}_p$  with L being the Laplacian.

Associated augmented Lagrangian:

$$\mathcal{A}(\boldsymbol{x}, \boldsymbol{u}) = \tilde{f}(\boldsymbol{x}) + \frac{\alpha}{2} \boldsymbol{x}^{\top} \boldsymbol{L} \boldsymbol{x} + \beta \boldsymbol{u}^{\top} \boldsymbol{L}^{1/2} \boldsymbol{x},$$

where  $\boldsymbol{u} \in \mathbb{R}^{np}$  is the dual variable,  $\alpha, \beta > 0$  are parameters. Minimize  $\mathcal{A}(\boldsymbol{x}, \boldsymbol{u})$  with a primal-dual algorithm:

$$\begin{aligned} \boldsymbol{x}_{k+1} = & \boldsymbol{x}_k - \eta \frac{\partial \mathcal{A}(\boldsymbol{x}_k, \boldsymbol{u}_k)}{\partial \boldsymbol{x}_k} = \boldsymbol{x}_k - \eta (\alpha \boldsymbol{L} \boldsymbol{x}_k + \beta \boldsymbol{L}^{1/2} \boldsymbol{u}_k + \nabla \tilde{f}(\boldsymbol{x}_k)), \\ \boldsymbol{u}_{k+1} = & \boldsymbol{u}_k + \eta \frac{\partial \mathcal{A}(\boldsymbol{x}_k, \boldsymbol{u}_k)}{\partial \boldsymbol{u}_k} = \boldsymbol{u}_k + \eta \beta \boldsymbol{L}^{1/2} \boldsymbol{x}_k, \end{aligned}$$

where  $\eta > 0$  is a fixed stepsize. With  $v_k = L^{1/2} u_k$ , this algorithm can be rewritten as  $x_{k+1} = x_k - \eta(\alpha L x_k + \beta v_k + \nabla \tilde{f}(x_k)),$  $v_{k+1} = v_k + \eta \beta L x_k.$ 

### Distributed primal-dual gradient descent algorithm



$$\begin{aligned} \boldsymbol{x}_{k+1} = & \boldsymbol{x}_k - \eta(\alpha \boldsymbol{L} \boldsymbol{x}_k + \beta \boldsymbol{v}_k + \nabla \tilde{f}(\boldsymbol{x}_k)), \\ \boldsymbol{v}_{k+1} = & \boldsymbol{v}_k + \eta \beta \boldsymbol{L} \boldsymbol{x}_k. \end{aligned}$$

#### Distributed Primal-Dual Gradient Descent Algorithm

$$x_{i,k+1} = x_{i,k} - \eta \Big( \alpha \sum_{j=1}^{n} L_{ij} x_{j,k} + \beta v_{i,k} + \nabla f_i(x_{i,k}) \Big),$$
  
$$v_{i,k+1} = v_{i,k} + \eta \beta \sum_{j=1}^{n} L_{ij} x_{j,k}.$$

### Distributed primal-dual gradient descent algorithm



$$\begin{split} \boldsymbol{x}_{k+1} = & \boldsymbol{x}_k - \eta(\alpha \boldsymbol{L} \boldsymbol{x}_k + \beta \boldsymbol{v}_k + \nabla \tilde{f}(\boldsymbol{x}_k)), \\ \boldsymbol{v}_{k+1} = & \boldsymbol{v}_k + \eta \beta \boldsymbol{L} \boldsymbol{x}_k. \end{split}$$

#### Distributed Primal-Dual Gradient Descent Algorithm

$$x_{i,k+1} = x_{i,k} - \eta \Big( \alpha \sum_{j=1}^{n} L_{ij} x_{j,k} + \beta v_{i,k} + \nabla f_i(x_{i,k}) \Big),$$
  
$$v_{i,k+1} = v_{i,k} + \eta \beta \sum_{j=1}^{n} L_{ij} x_{j,k}.$$

This algorithm is **single-loop** and communicates **only one** variable.



#### Distributed Primal-Dual Gradient Descent Algorithm

$$x_{i,k+1} = x_{i,k} - \eta \Big( \alpha \sum_{j=1}^{n} L_{ij} x_{j,k} + \beta v_{i,k} + \nabla f_i(x_{i,k}) \Big),$$
  
$$v_{i,k+1} = v_{i,k} + \eta \beta \sum_{j=1}^{n} L_{ij} x_{j,k}.$$

#### Theorem

If each  $f_i(x)$  is smooth ( $\nabla f_i(x)$  is Lipschitz continuous) and f(x) satisfies P–Ł, then the primal-dual gradient descent algorithm linearly converges to a global optimum: there exists  $\rho \in (0, 1)$  such that

$$\underbrace{\frac{1}{n}\sum_{i=1}^{n}\|x_{i,k}-\bar{x}_{k}\|^{2}}_{\text{Consensus}} + \underbrace{f(\bar{x}_{k})-f^{*}}_{\text{Optimization}} = \mathcal{O}(\rho^{k}).$$



#### Distributed Primal-Dual Gradient Descent Algorithm

$$x_{i,k+1} = x_{i,k} - \eta \Big( \alpha \sum_{j=1}^{n} L_{ij} x_{j,k} + \beta v_{i,k} + \nabla f_i(x_{i,k}) \Big),$$
  
$$v_{i,k+1} = v_{i,k} + \eta \beta \sum_{j=1}^{n} L_{ij} x_{j,k}.$$

#### Theorem

If each  $f_i(x)$  is smooth ( $\nabla f_i(x)$  is Lipschitz continuous) and f(x) satisfies P–Ł, then the primal-dual gradient descent algorithm linearly converges to a global optimum: there exists  $\rho \in (0, 1)$  such that

$$\underbrace{\frac{1}{n}\sum_{i=1}^{n}\|x_{i,k}-\bar{x}_{k}\|^{2}}_{\text{Consensus}} + \underbrace{f(\bar{x}_{k})-f^{*}}_{\text{Optimization}} = \mathcal{O}(\rho^{k}).$$

#### Remark

- Linear convergence is achieved without strong convexity, even without convexity.
- Parameters  $\alpha$ ,  $\beta$ ,  $\eta$  do not have to depend on  $\nu$ .



P-Ł condition: 
$$\|\nabla f(x)\|^2 \ge \nu(f(x) - f^*), \ \forall x \in \mathbb{R}^p.$$

#### Theorem

If each  $f_i(x)$  is smooth and f(x) satisfies P–Ł, there exists  $\rho \in (0,1)$  such that

$$\frac{1}{n}\sum_{i=1}^{n} \|x_{i,k} - \bar{x}_k\|^2 + f(\bar{x}_k) - f^* = \mathcal{O}(\rho^k).$$

#### **Proof sketch:**

Consider the nonnegative potential function:

$$V_{k} = \frac{1}{2} \|\boldsymbol{x}_{k}\|_{\boldsymbol{K}}^{2} + \frac{1}{2} \left\| \boldsymbol{v}_{k} + \frac{1}{\beta} \boldsymbol{g}_{k}^{0} \right\|_{\boldsymbol{Q} + \frac{\alpha}{\beta} \boldsymbol{K}}^{2} + \boldsymbol{x}_{k}^{\top} \boldsymbol{K} \left( \boldsymbol{v}_{k} + \frac{1}{\beta} \boldsymbol{g}_{k}^{0} \right) + n(f(\bar{x}_{k}) - f^{*}),$$
  
where  $\|\boldsymbol{x}_{k}\|_{\boldsymbol{K}}^{2} = \sum_{i=1}^{n} \|x_{i,k} - \bar{x}_{k}\|^{2}$  and  $\boldsymbol{g}_{k}^{0} = \operatorname{col}(\nabla f_{1}(\bar{x}_{k}), \dots, \nabla f_{n}(\bar{x}_{k})).$ 



P-Ł condition: 
$$\|\nabla f(x)\|^2 \ge \nu(f(x) - f^*), \ \forall x \in \mathbb{R}^p.$$

#### Theorem

If each  $f_i(x)$  is smooth and f(x) satisfies P–Ł, there exists  $\rho\in(0,1)$  such that

$$\frac{1}{n}\sum_{i=1}^{n} \|x_{i,k} - \bar{x}_k\|^2 + f(\bar{x}_k) - f^* = \mathcal{O}(\rho^k).$$

#### **Proof sketch:**

Consider the nonnegative potential function:

$$V_{k} = \frac{1}{2} \|\boldsymbol{x}_{k}\|_{\boldsymbol{K}}^{2} + \frac{1}{2} \left\| \boldsymbol{v}_{k} + \frac{1}{\beta} \boldsymbol{g}_{k}^{0} \right\|_{\boldsymbol{Q} + \frac{\alpha}{\beta} \boldsymbol{K}}^{2} + \boldsymbol{x}_{k}^{\top} \boldsymbol{K} \left( \boldsymbol{v}_{k} + \frac{1}{\beta} \boldsymbol{g}_{k}^{0} \right) + n(f(\bar{x}_{k}) - f^{*}),$$
  
where  $\|\boldsymbol{x}_{k}\|_{\boldsymbol{K}}^{2} = \sum_{i=1}^{n} \|x_{i,k} - \bar{x}_{k}\|^{2}$  and  $\boldsymbol{g}_{k}^{0} = \operatorname{col}(\nabla f_{1}(\bar{x}_{k}), \dots, \nabla f_{n}(\bar{x}_{k})).$ 

This function is nonincreasing:

$$V_{k+1} \le V_k - \left[c_1 \|\boldsymbol{x}_k\|_{\boldsymbol{K}}^2 + c_2 \|\boldsymbol{v}_k + \frac{1}{\beta} \boldsymbol{g}_k^0\|_{\boldsymbol{K}}^2 + c_3 n \|\nabla f(\bar{x}_k)\|^2\right]$$



P-Ł condition: 
$$\|\nabla f(x)\|^2 \ge \nu(f(x) - f^*), \ \forall x \in \mathbb{R}^p.$$

#### Theorem

If each  $f_i(x)$  is smooth and f(x) satisfies P–Ł, there exists  $\rho\in(0,1)$  such that

$$\frac{1}{n}\sum_{i=1}^{n} \|x_{i,k} - \bar{x}_k\|^2 + f(\bar{x}_k) - f^* = \mathcal{O}(\rho^k).$$

#### **Proof sketch:**

Consider the nonnegative potential function:

$$V_{k} = \frac{1}{2} \|\boldsymbol{x}_{k}\|_{\boldsymbol{K}}^{2} + \frac{1}{2} \left\| \boldsymbol{v}_{k} + \frac{1}{\beta} \boldsymbol{g}_{k}^{0} \right\|_{\boldsymbol{Q} + \frac{\alpha}{\beta} \boldsymbol{K}}^{2} + \boldsymbol{x}_{k}^{\top} \boldsymbol{K} \left( \boldsymbol{v}_{k} + \frac{1}{\beta} \boldsymbol{g}_{k}^{0} \right) + n(f(\bar{x}_{k}) - f^{*}),$$
  
where  $\|\boldsymbol{x}_{k}\|_{\boldsymbol{K}}^{2} = \sum_{i=1}^{n} \|x_{i,k} - \bar{x}_{k}\|^{2}$  and  $\boldsymbol{g}_{k}^{0} = \operatorname{col}(\nabla f_{1}(\bar{x}_{k}), \dots, \nabla f_{n}(\bar{x}_{k})).$ 

This function is nonincreasing:

$$V_{k+1} \leq V_k - \left[c_1 \|\boldsymbol{x}_k\|_{\boldsymbol{K}}^2 + c_2 \|\boldsymbol{v}_k + \frac{1}{\beta} \boldsymbol{g}_k^0\|_{\boldsymbol{K}}^2 + c_3 n \|\nabla f(\bar{x}_k)\|^2\right].$$

P-Ł implies

$$\|\nabla f(\bar{x}_k)\|^2 \ge \nu(f(\bar{x}_k) - f^*).$$



#### Theorem

If each  $f_i(x)$  is smooth and f(x) satisfies P–Ł, there exists  $\rho \in (0,1)$  such that

$$\frac{1}{n}\sum_{i=1}^{n} \|x_{i,k} - \bar{x}_k\|^2 + f(\bar{x}_k) - f^* = \mathcal{O}(\rho^k).$$

#### **Proof sketch:**

Combine the nonincreasing property of  $V_k$ 

$$V_{k+1} \le V_k - \left[c_1 \|\boldsymbol{x}_k\|_{\boldsymbol{K}}^2 + c_2 \|\boldsymbol{v}_k + \frac{1}{\beta} \boldsymbol{g}_k^0\|_{\boldsymbol{K}}^2 + c_3 n \|\nabla f(\bar{x}_k)\|^2\right]$$

and P-Ł condition  $\|\nabla f(\bar{x}_k)\|^2 \ge \nu(f(\bar{x}_k) - f^*)$ :

$$V_{k+1} \leq V_k - \left[c_1 \| \boldsymbol{x}_k \|_{\boldsymbol{K}}^2 + c_2 \| \boldsymbol{v}_k + \frac{1}{\beta} \boldsymbol{g}_k^0 \|_{\boldsymbol{K}}^2 + c_3 \nu n (f(\bar{\boldsymbol{x}}_k) - f^*)\right]$$
  
$$\leq (1 - c_4) V_k \leq \cdots \leq (1 - c_4)^{k+1} V_0.$$



#### Theorem

If each  $f_i(x)$  is smooth and f(x) satisfies P–Ł, there exists  $\rho \in (0,1)$  such that

$$\frac{1}{n}\sum_{i=1}^{n} \|x_{i,k} - \bar{x}_k\|^2 + f(\bar{x}_k) - f^* = \mathcal{O}(\rho^k).$$

#### Proof sketch:

Combine the nonincreasing property of  $V_k$ 

$$V_{k+1} \le V_k - \left[c_1 \|\boldsymbol{x}_k\|_{\boldsymbol{K}}^2 + c_2 \|\boldsymbol{v}_k + \frac{1}{\beta} \boldsymbol{g}_k^0\|_{\boldsymbol{K}}^2 + c_3 n \|\nabla f(\bar{x}_k)\|^2\right]$$

and P-Ł condition  $\|\nabla f(\bar{x}_k)\|^2 \ge \nu(f(\bar{x}_k) - f^*)$ :

$$V_{k+1} \leq V_k - \left[c_1 \| \boldsymbol{x}_k \|_{\boldsymbol{K}}^2 + c_2 \| \boldsymbol{v}_k + \frac{1}{\beta} \boldsymbol{g}_k^0 \|_{\boldsymbol{K}}^2 + c_3 \nu n (f(\bar{x}_k) - f^*)\right]$$
  
$$\leq (1 - c_4) V_k \leq \cdots \leq (1 - c_4)^{k+1} V_0.$$

Using  $\frac{1}{n} \sum_{i=1}^{n} ||x_{i,k} - \bar{x}_k||^2 + f(\bar{x}_k) - f^* \le c_5 V_k$  gives the result.



Without P-k condition, the nonincreasing property of  $V_k$  still holds, i.e.,

$$V_{k+1} \le V_k - \left[c_1 \|\boldsymbol{x}_k\|_{\boldsymbol{K}}^2 + c_2 \|\boldsymbol{v}_k + \frac{1}{\beta} \boldsymbol{g}_k^0\|_{\boldsymbol{K}}^2 + c_3 n \|\nabla f(\bar{x}_k)\|^2\right].$$



Without P-Ł condition, the nonincreasing property of  $V_k$  still holds, i.e.,

$$V_{k+1} \leq V_k - \left[ c_1 \| \boldsymbol{x}_k \|_{\boldsymbol{K}}^2 + c_2 \| \boldsymbol{v}_k + \frac{1}{\beta} \boldsymbol{g}_k^0 \|_{\boldsymbol{K}}^2 + c_3 n \| \nabla f(\bar{x}_k) \|^2 \right].$$

Sum over  $k = 0, \ldots, T - 1$  and rearrange terms:

$$\sum_{k=0}^{T-1} \left[ c_1 \| \boldsymbol{x}_k \|_{\boldsymbol{K}}^2 + c_2 \left\| \boldsymbol{v}_k + \frac{1}{\beta} \boldsymbol{g}_k^0 \right\|_{\boldsymbol{K}}^2 + c_3 n \| \nabla f(\bar{x}_k) \|^2 \right] \le V_0 - V_T \le V_0.$$



Without P-L condition, the nonincreasing property of  $V_k$  still holds, i.e.,

$$V_{k+1} \leq V_k - \left[ c_1 \| \boldsymbol{x}_k \|_{\boldsymbol{K}}^2 + c_2 \| \boldsymbol{v}_k + \frac{1}{\beta} \boldsymbol{g}_k^0 \|_{\boldsymbol{K}}^2 + c_3 n \| \nabla f(\bar{x}_k) \|^2 \right].$$

Sum over k = 0, ..., T - 1 and rearrange terms:

$$\sum_{k=0}^{T-1} \left[ c_1 \| \boldsymbol{x}_k \|_{\boldsymbol{K}}^2 + c_2 \left\| \boldsymbol{v}_k + \frac{1}{\beta} \boldsymbol{g}_k^0 \right\|_{\boldsymbol{K}}^2 + c_3 n \| \nabla f(\bar{x}_k) \|^2 \right] \le V_0 - V_T \le V_0.$$

Division by 
$$nT$$
 gives  $\frac{1}{T} \sum_{k=0}^{T-1} \left[ \frac{1}{n} \sum_{i=1}^{n} \|x_{i,k} - \bar{x}_k\|^2 + \|\nabla f(\bar{x}_k)\|^2 \right] = \mathcal{O}(\frac{1}{T}).$ 

#### KTH vertrikker vertrikker

### Extension: sublinear convergence to stationary points

Without P-L condition, the nonincreasing property of  $V_k$  still holds, i.e.,

$$V_{k+1} \leq V_k - \left[ c_1 \| \boldsymbol{x}_k \|_{\boldsymbol{K}}^2 + c_2 \| \boldsymbol{v}_k + \frac{1}{\beta} \boldsymbol{g}_k^0 \|_{\boldsymbol{K}}^2 + c_3 n \| \nabla f(\bar{x}_k) \|^2 \right].$$

Sum over  $k = 0, \ldots, T - 1$  and rearrange terms:

$$\sum_{k=0}^{T-1} \left[ c_1 \| \boldsymbol{x}_k \|_{\boldsymbol{K}}^2 + c_2 \left\| \boldsymbol{v}_k + \frac{1}{\beta} \boldsymbol{g}_k^0 \right\|_{\boldsymbol{K}}^2 + c_3 n \| \nabla f(\bar{x}_k) \|^2 \right] \le V_0 - V_T \le V_0.$$

Division by 
$$nT$$
 gives  $\frac{1}{T} \sum_{k=0}^{T-1} \left[ \frac{1}{n} \sum_{i=1}^{n} \|x_{i,k} - \bar{x}_k\|^2 + \|\nabla f(\bar{x}_k)\|^2 \right] = \mathcal{O}(\frac{1}{T}).$ 

Corollary (without P-Ł condition)

If each  $f_i(x)$  is smooth, the primal-dual gradient descent algorithm converges to a stationary point sublinearly:

$$\frac{1}{T} \sum_{k=0}^{T-1} \left[ \underbrace{\frac{1}{n} \sum_{i=1}^{n} \|x_{i,k} - \bar{x}_k\|^2}_{\text{Consensus}} + \underbrace{\|\nabla f(\bar{x}_k)\|^2}_{\text{Optimization}} \right] = \mathcal{O}(\frac{1}{T})$$



Without P-L condition, the nonincreasing property of  $V_k$  still holds, i.e.,

$$V_{k+1} \le V_k - \left[ c_1 \| \boldsymbol{x}_k \|_{\boldsymbol{K}}^2 + c_2 \| \boldsymbol{v}_k + \frac{1}{\beta} \boldsymbol{g}_k^0 \|_{\boldsymbol{K}}^2 + c_3 n \| \nabla f(\bar{x}_k) \|^2 \right].$$

Sum over  $k = 0, \ldots, T - 1$  and rearrange terms:

$$\sum_{k=0}^{T-1} \left[ c_1 \| \boldsymbol{x}_k \|_{\boldsymbol{K}}^2 + c_2 \left\| \boldsymbol{v}_k + \frac{1}{\beta} \boldsymbol{g}_k^0 \right\|_{\boldsymbol{K}}^2 + c_3 n \| \nabla f(\bar{x}_k) \|^2 \right] \le V_0 - V_T \le V_0.$$

Division by 
$$nT$$
 gives  $\frac{1}{T} \sum_{k=0}^{T-1} \left[ \frac{1}{n} \sum_{i=1}^{n} \|x_{i,k} - \bar{x}_k\|^2 + \|\nabla f(\bar{x}_k)\|^2 \right] = \mathcal{O}(\frac{1}{T}).$ 

Corollary (without P-Ł condition)

If each  $f_i(x)$  is smooth, the primal-dual gradient descent algorithm converges to a stationary point sublinearly:

$$\frac{1}{T} \sum_{k=0}^{T-1} \left[ \underbrace{\frac{1}{n} \sum_{i=1}^{n} \|x_{i,k} - \bar{x}_k\|^2}_{\text{Consensus}} + \underbrace{\|\nabla f(\bar{x}_k)\|^2}_{\text{Optimization}} \right] = \mathcal{O}(\frac{1}{T})$$

Remark: guaranteed convergence rate also for nonconvex cost functions.

X.L. Yi et al | CDC 2020 | Convergence Analysis

# Nonconvex distributed regularized logistic regression problem



Each component function:

$$f_i(x) = \frac{n}{m} \sum_{l=1}^{m_i} (y_{il} \log(1 + \exp(-x^\top z_{il})) + (1 - y_{il}) \log(1 + \exp(x^\top z_{il}))) + \sum_{s=1}^p \frac{\lambda \mu[x]_s^2}{1 + \mu[x]_s^2}$$

Compared algorithms:

- Distributed primal-dual gradient descent algorithm (DPD-GDA) [Our paper]
- Distributed gradient descent algorithm (D-GDA)
- Distributed gradient tracking algorithm (D-GTA)
- Distributed proximal primal-dual algorithm (Prox-PDA)
- Distributed xFILTER algorithm (**xFILTER**)

[Zeng & Yin, TSP, 2018] [Qu & Li, TCNS, 2018] [Hong et al, ICML, 2017] [Sun & Hong, TSP 2019]

# Nonconvex distributed regularized logistic regression problem





• Our DPD-GDA gives the best performance in general.

X.L. Yi et al | CDC 2020 | Simulations

Conclusions



• **Problem**: 
$$\min_{x \in \mathbb{R}^p} f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x)$$
, where each  $f_i$  could be nonconvex

• Assumptions: each  $f_i$  is smooth, f satisfies P-Ł, and  $\mathcal{G}$  is connected

P-Ł condition :  $\|\nabla f(x)\|^2 \ge \nu(f(x) - f^*), \ \forall x \in \mathbb{R}^p$ 

- Method: primal-dual gradient descent algorithm ( $\nu$  is not used)
- Result:
  - Linear convergence with relaxed assumptions
  - Sublinear convergence for nonconvex cost functions

Conclusions



• **Problem**: 
$$\min_{x \in \mathbb{R}^p} f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x)$$
, where each  $f_i$  could be nonconvex

• Assumptions: each  $f_i$  is smooth, f satisfies P-Ł, and  $\mathcal{G}$  is connected

P-Ł condition :  $\|\nabla f(x)\|^2 \ge \nu(f(x) - f^*), \ \forall x \in \mathbb{R}^p$ 

- Method: primal-dual gradient descent algorithm ( $\nu$  is not used)
- Result:
  - Linear convergence with relaxed assumptions
  - Sublinear convergence for nonconvex cost functions
- Extensions:
  - Stochastic gradient descent algorithm
  - Zeroth-order algorithm

# Thank you for your attention!