

# A COMMUNICATION EFFICIENT STOCHASTIC MULTI-BLOCK ALTERNATING DIRECTION METHOD OF MULTIPLIERS

{ HAO YU}  
AMAZON, SEATTLE, WA

## 1. LINEARLY CONSTRAINED STO-OPT

- Linearly constrained stochastic convex programs

$$\min_{\mathbf{x} \in \mathcal{X}_i, \forall i} f(\mathbf{x}) = \sum_{i=1}^N f_i(\mathbf{x}_i) \text{ s.t. } \sum_{i=1}^N \mathbf{A}_i \mathbf{x}_i = \mathbf{b}$$

- $N$  arbitrary; Each  $f_i(\mathbf{x}_i) \triangleq \mathbb{E}_{\xi}[f_i(\mathbf{x}_i; \xi)]$  with expensive true gradient but cheap unbiased stochastic gradient.

### Applications

- Large scale linearly constrained optimization, e.g., linear programs: Too large to store or solve on a single node.
- Distributed machine learning:  $N$  distributed nodes (with possibly non-identical training data) jointly train a common ML model.

## 2. COMMUNICATION EFFICIENT ADMM

- ADMM is effective and popular for distributed optimization, yet suffers significant communication overhead for passing Lagrange multipliers.
- Conventional ADMM involves a communication step following immediately a computation step. (Communication often much more expensive than SGD computation.)
- This paper develops communication efficient multi-block stochastic ADMM to reduce communication rounds w/o. sacrificing convergence.

## 3. OUR ALGORITHM

### Alg1 : Two-Layer Communication Efficient ADMM

- Input: Algorithm parameters  $T, \{\rho^{(t)}\}_{t \geq 1}, \{\nu^{(t)}\}_{t \geq 1}$  and  $\{K^{(t)}\}_{t \geq 1}$ .
- Initialize arbitrary  $\mathbf{y}_i^{(0)} \in \mathcal{X}_i, \forall i$ ,  $\mathbf{r}^{(0)} = \sum_{i=1}^N \mathbf{A}_i \mathbf{y}_i^{(0)} - \mathbf{b}$ ,  $\boldsymbol{\lambda}^{(0)} = \mathbf{0}$ , and  $t = 1$ .
- while  $t \leq T$  do
- Each node  $i$  defines  $\phi_i^{(t)}(\mathbf{x}_i) \triangleq$

$$f_i(\mathbf{x}_i) + \rho^{(t)} \left\langle \mathbf{r}^{(t-1)} + \frac{1}{\rho^{(t)}} \boldsymbol{\lambda}^{(t-1)}, \mathbf{A}_i \mathbf{x}_i - \frac{\mathbf{b}}{N} \right\rangle + \frac{\nu^{(t)}}{2} \|\mathbf{x}_i - \mathbf{y}_i^{(t-1)}\|^2$$

and in parallel updates  $\mathbf{x}_i^{(t)}, \mathbf{y}_i^{(t)}$  using local sub-procedure Alg 2 via

$$(\mathbf{x}_i^{(t)}, \mathbf{y}_i^{(t)}) = \text{STO-LOCAL}(\phi_i^{(t)}(\cdot), \mathcal{X}_i, \mathbf{y}_i^{(t-1)}, K^{(t)})$$

- Each node  $i$  passes  $\mathbf{x}_i^{(t)}$  and  $\mathbf{y}_i^{(t)}$  between nodes or to a parameter server. Update  $\boldsymbol{\lambda}^{(t)}$  and  $\mathbf{r}^{(t)}$

$$\boldsymbol{\lambda}^{(t)} = \boldsymbol{\lambda}^{(t-1)} + \rho^{(t)} \left( \sum_{i=1}^N \mathbf{A}_i \mathbf{x}_i^{(t)} - \mathbf{b} \right)$$

$$\mathbf{r}^{(t)} = \sum_{i=1}^N \mathbf{A}_i \mathbf{y}_i^{(t)} - \mathbf{b}.$$

- Update  $t \leftarrow t + 1$ .

end while

- Output:  $\bar{\mathbf{x}}^{(T)} = \frac{1}{\sum_{t=1}^T \rho^{(t)}} \sum_{t=1}^T \rho^{(t)} \mathbf{x}^{(t)}$

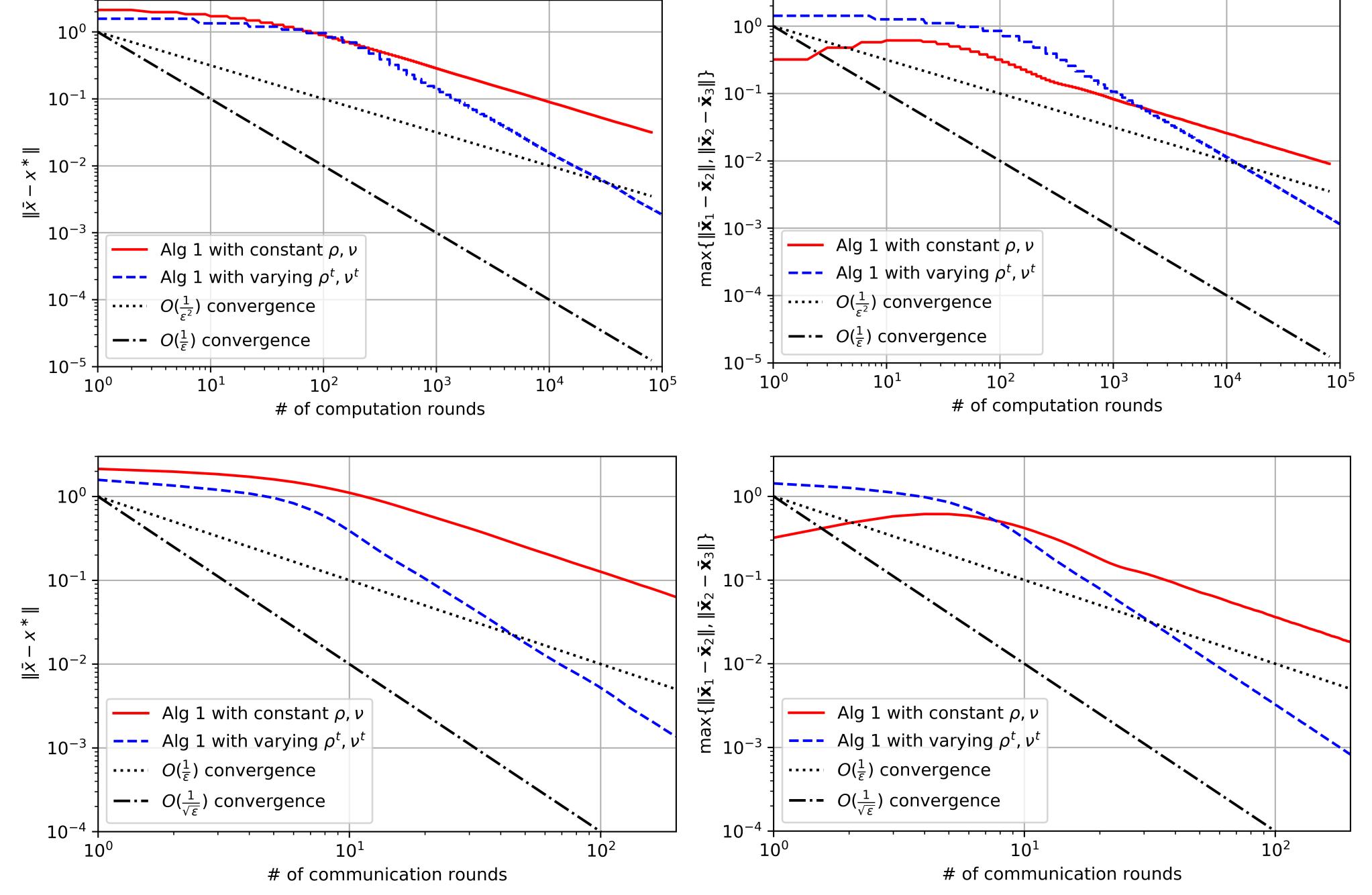
- SGD-LOCAL (Alg 2) is a  $K^{(t)}$  step SGD procedure with designed initialization, step size and averaging rules.

## 4. PERFORMANCE ANALYSIS

- Under corresponding algorithm parameter rules, to achieve an  $O(\epsilon)$  accuracy solution
  - General Convex: Alg 1 uses  $\tilde{O}(1/\epsilon^2)$  SGD update rounds and  $\tilde{O}(1/\epsilon)$  inter-node communication rounds.
  - Strongly Convex: Alg 1 uses  $\tilde{O}(1/\epsilon)$  SGD update rounds and  $\tilde{O}(1/\sqrt{\epsilon})$  inter-node communication rounds.
- The # of communication rounds is only the square root of that of computation (SGD update) rounds.
- Lowest computation complexity for stochastic convex opt with lower communication complexity than other stochastic ADMM.

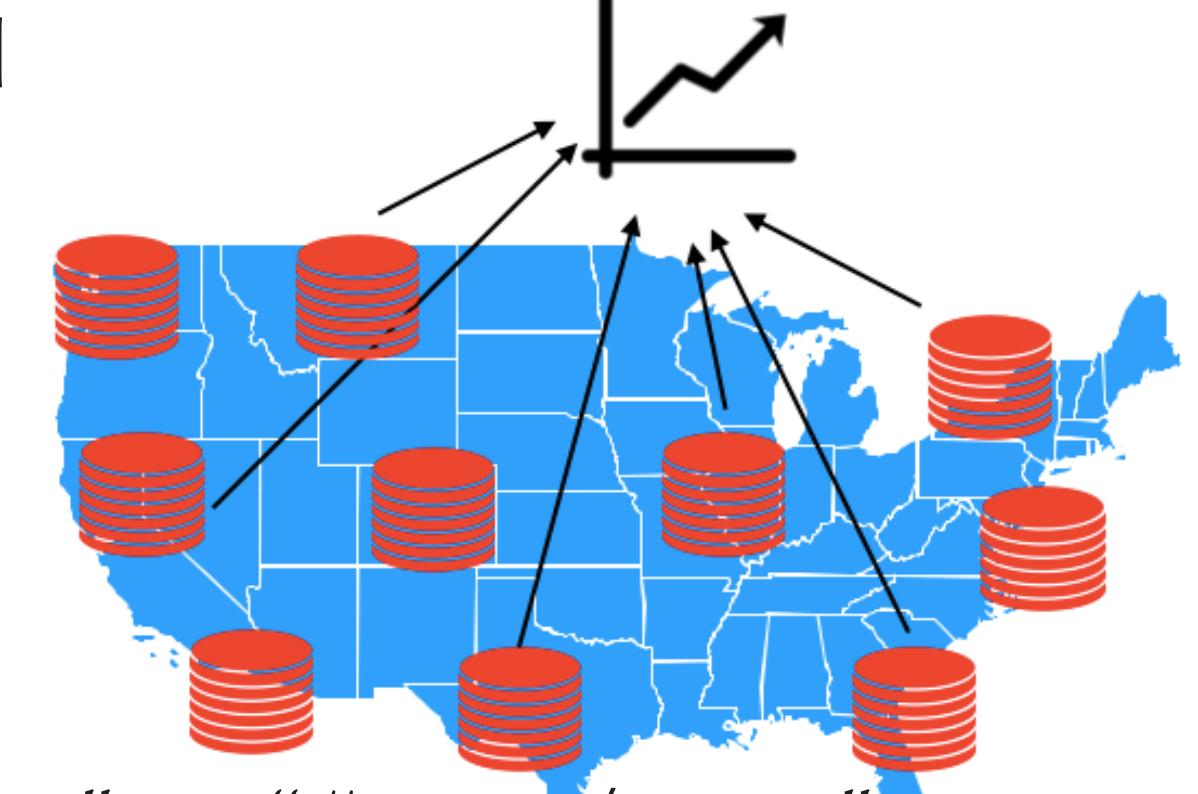
## 5. EXPERIMENTS

- Convergence rate verification: smooth strongly convex



Setting algorithm parameters in line with our theory yields the (better) proven convergence rates

- Distributed logistic regression: 10 distributed nodes with sharded data jointly train a common model
  - RPDBUS ADMM [Gao et.al. 2016]
  - DCS [Lan et.al. 2017]



Compare "loss/consensus" v.s. "# comp/comm"

